



UNIVERSIDAD DE LAS PALMAS  
DE GRAN CANARIA

PROGRAMA DE DOCTORADO

SIMULACIÓN NUMÉRICA EN CIENCIAS Y TECNOLOGÍA

TESIS DOCTORAL

EXTRACCIÓN AUTOMÁTICA

DE

NEXOS LÉXICOS

AUTORA: ISABEL SÁNCHEZ BERRIEL

LAS PALMAS DE GRAN CANARIA, MAYO DE 2015





PROGRAMA DE DOCTORADO  
SIMULACIÓN NUMÉRICA EN CIENCIAS Y TECNOLOGÍA  
DEPARTAMENTO DE MATEMÁTICAS

TESIS DOCTORAL  
EXTRACCIÓN AUTOMÁTICA  
DE  
NEXOS LÉXICOS

AUTORA: ISABEL SÁNCHEZ BERRIEL

DIRECTOR:  
DR. OCTAVIO SANTANA SUÁREZ

CO-DIRECTOR:  
DR. JOSÉ RAFAEL PÉREZ AGUIAR

LAS PALMAS DE GRAN CANARIA, MAYO DE 2015



# Agradecimientos

Este trabajo de tesis supone el punto de arribo gracias a los esfuerzos realizados por un numeroso grupo de personas que han persistido en apoyar la trayectoria seguida desde que emprendí mis primeros pasos con las colocaciones hasta que comencé a “vislumbrar la luz al final del túnel”. A todos ellos quiero agradecer su aportación, cada uno en su especialidad y vocación universitaria.

En primer lugar quiero referirme, al Dr. Octavio Santana Suárez, director de esta tesis, quien no sólo abrió la puerta de su ciencia e interés por el tema, sino que la ha mantenido de par en par sin ambages, alimentando la obstinación, el tesón y la ilusión de las que me he armado para recorrer este camino; quisiera agradecerle el esfuerzo, todo el tiempo dedicado, el conocimiento compartido y los recursos facilitados sin los que no hubiese podido llevar a buen puerto este proyecto. Al Dr. José Pérez Aguiar le agradezco sinceramente todo el apoyo como co-director, y cada una de sus orientaciones que han contribuido a la consecución de este trabajo.

Reservo un lugar para mis compañeros de la Universidad de La Laguna: Luzma, Jesús y Macu, que alguna que otra vez les ha tocado sustituirme en mis estancias en Las Palmas. Inestimable ha sido la ayuda de José Luis Roda, que siempre que lo he necesitado ha aportado la tecnología que tantas palabras parecían devorar. A Félix García, mi maestro, mil gracias por además ayudar en todo lo que estaba en su mano para que yo pudiera desempeñar esta tarea.

Un recuerdo especial para mis amigos Bea y Robert, que con sus risas y todo lo demás siempre han estado ahí.

Por último a todas esas personas que, como en el cuento del zapatero y los duendes, silenciosamente han propiciado cada uno de los pasos dados. Mis padres, mis hijas, Cristo, Vicky, Rafa, Sayo, mi familia, todos ellos mis duendes.



## Resumen

En el presente trabajo se presentan análisis y soluciones al problema de la extracción automática de colocaciones léxicas. Con independencia de la lengua de que se trate, a partir de la propuesta de Sinclair —estudia las frecuencias de aparición en corpus textuales—, las técnicas resuelven el problema desde el punto de vista estadístico: una combinación de palabras se considera colocación si su aparición conjunta queda motivada en lugar de provenir del puro azar. Aunque este enfoque sea demasiado simplista para el lingüista, aporta la herramienta básica para resolver el problema desde la perspectiva computacional. Bajo la premisa de implementar una solución que requiera el menor número posible de recursos léxicos, se conjuga el conocimiento lingüístico con un planteamiento estadístico. Las conclusiones y aportaciones que se derivan dan respuesta a la extracción de colocaciones de un corpus textual sea cual sea su volumen.

Se maneja un corpus textual de 11 000 textos literarios y no literarios con un número aproximado de 300 000 000 palabras de la lengua española. A partir de esta amplia muestra del español, se procede a la extracción de colocaciones mediante técnicas basadas en las frecuencias de palabras; se analizan los resultados que se obtienen y se contrastan con las combinaciones del Diccionario Combinatorio del Español Contemporáneo, Redes (DCECR). Se dan soluciones a la inestabilidad que producen las marcadas diferencias entre las frecuencias de uso de las distintas palabras en el corpus, que se ve acentuada por el uso de un corpus con tan amplio volumen. Por una parte, se diseña una metodología de análisis centrada en la palabra, y por otra se proponen nuevos indicadores, basados en las técnicas de detección de *outliers*, que tratan de captar las características de uso preferentes de las uniones presentes en este tipo de combinaciones. Las propuestas se orientan de modo que sean válidas en corpus menos extensos: se utiliza como corpus reducido las novelas de D. Benito Pérez Galdós.

Se resuelve también el problema de extracción automática de nexos semánticos motivados por la selección de argumentos por parte de los predicados; se enfoca el estudio hacia la determinación de clases semánticas entre los elementos que combinan con una determinada palabra, se lleva a cabo un repaso de las técnicas que permiten extraer tales grupos de forma automática. En este caso, se opta por adaptar el indicador *asociación de selección* (Resnik, 1997), que mide el vínculo entre un grupo de argumentos y su predicado, utilizando para evaluarlos diccionarios ideológicos del español como clasificaciones de palabras presentes en el corpus. Las pruebas se realizaron con el Diccionario Ideológico Vox (DIV); en concreto, el trabajo se centra en los casos: *verbo + grupos de sustantivos*, *adjetivo + grupos de sustantivos* y *adverbio + grupos de verbos*.

Por último, se aporta una aplicación software que permite consultar las propiedades combinatorias de una palabra de la lengua española a través de Internet, así como la extracción automática de las colocaciones en textos en español de forma amigable para el usuario.



# Índice

<b>1. INTRODUCCIÓN.....</b>	<b>1</b>
1.1. INTRODUCCIÓN .....	1
1.2. ENFOQUE ESTADÍSTICO.....	2
1.3. ENFOQUE LÉXICO-SEMÁNTICO.....	4
1.4. LAS COLOCACIONES EN ESPAÑOL.....	6
1.5. CARACTERÍSTICAS.....	8
1.5.1. <i>Flexibilidad</i> .....	8
1.5.2. <i>Determinación</i> .....	10
1.5.3. <i>Frecuencia, Preferencia</i> .....	13
1.5.4. <i>Composicionalidad</i> .....	14
1.5.5. <i>Gradualidad</i> .....	15
1.6. CLASIFICACIÓN.....	16
1.6.1. <i>Categoría gramatical de los colocados</i> .....	16
1.6.2. <i>Grado de restricción de los colocados</i> .....	22
1.6.3. <i>Semántica</i> .....	23
1.7. CARACTERÍSTICAS POR CATEGORÍAS.....	26
1.7.1. <i>Sustantivo + Verbo</i> .....	26
1.7.2. <i>Verbo deslexicalizado + sustantivo<sub>CD</sub></i> .....	28
1.7.3. <i>Verbo deslexicalizado + (sustantivo<sub>CD</sub>) + preposición + sustantivo</i> .....	31
1.7.4. <i>Selección de una acepción especial en el verbo</i> .....	32
1.7.5. <i>Colocaciones léxicas</i> .....	33
1.7.6. <i>Sustantivo<sub>sujeto</sub> + verbo</i> .....	34
1.7.7. <i>Sustantivo<sub>1</sub> + de + Sustantivo<sub>2</sub></i> .....	35
1.7.8. <i>Sustantivo + Adjetivo</i> .....	35
1.7.9. <i>Verbo + Adverbio</i> .....	39
1.7.10. <i>Adverbio + Adjetivo</i> .....	40
1.7.11. <i>Verbo + Adjetivo</i> .....	40
1.7.12. <i>Verbo + como + sustantivo</i> .....	41
1.8. LAS COLOCACIONES EN LOS DICCIONARIOS.....	41
1.8.1. <i>Colocaciones en el DEA</i> .....	41
1.8.2. <i>Colocaciones en el DUE</i> .....	44
1.8.3. <i>Diccionario VOX</i> .....	44
1.9. ELEMENTOS AJENOS A LAS COLOCACIONES.....	45
<b>2. EXTRACCIÓN AUTOMÁTICA DE COLOCACIONES UTILIZANDO CRITERIOS ESTADÍSTICOS.....</b>	<b>47</b>
2.1. MEDIDAS DE ASOCIACIÓN .....	47
2.1.1. <i>Frecuencia Relativa</i> .....	48

---

2.1.2.	<i>Información Mutua</i> .....	48
2.1.3.	<i>z-score</i> .....	49
2.1.4.	<i>t-score</i> .....	49
2.1.5.	<i>Test de Dunning</i> .....	49
2.1.6.	<i>Test de Poisson</i> .....	50
2.2.	BASE DE DATOS DE INFORMACIÓN COMBINATORIA. ....	50
2.2.1.	<i>Características del Corpus</i> .....	51
2.3.	LA BASE DE DATOS.....	59
2.4.	ANÁLISIS DE LOS RESULTADOS.....	63
2.4.1.	<i>Frecuencia Relativa e Información Mutua</i> .....	64
2.4.2.	<i>Test estadísticos</i> .....	70
2.4.3.	<i>Análisis por tipos de combinación según la categoría gramatical</i> .....	74
2.4.4.	<i>Combinaciones sustantivo + adjetivo</i> .....	76
2.4.5.	<i>Combinaciones verbo + adverbio</i> .....	78
2.5.	INFLUENCIA DEL TAMAÑO DEL CORPUS.....	79
2.6.	LAS COLOCACIONES COMO COMBINACIONES ATÍPICAS.....	80
2.6.1.	<i>Método basado en la desigualdad de Chebyshev</i> .....	82
2.6.2.	<i>Métodos robustos</i> .....	84
2.6.3.	<i>Método basado en el diagrama de cajas y bigotes</i> .....	85
2.6.4.	<i>Método de Hampel</i> .....	90
2.7.	ESTRATEGIA PARA LA EXTRACCIÓN DE INFORMACIÓN COLOCACIONAL GENERADA POR UN CORPUS EXTENSO.....	93
2.7.1.	<i>Resultados obtenidos sobre el corpus completo</i> .....	96
2.7.2.	<i>Resultados obtenidos sobre el corpus de Galdós</i> .....	110
<b>3.</b>	<b>DETECCIÓN DE NEXOS CON GRUPOS SEMÁNTICOS</b> .....	<b>113</b>
3.1.	ANÁLISIS CUALITATIVO DE LOS DATOS Y SELECCIÓN LÉXICA.....	113
3.2.	SIMILITUD SEMÁNTICA.....	120
3.2.1.	<i>Técnicas para determinar la similitud semántica</i> .....	120
3.3.	OBTENCIÓN DE CLASES LÉXICAS.....	122
3.3.1.	<i>Agrupamiento mediante técnicas estadísticas</i> .....	123
3.3.2.	<i>Obtención de grupos semánticos a partir de la relación predicado-argumentos</i> 133	
3.3.3.	<i>Diccionarios Ideológicos, bases de conocimientos léxicas</i> .....	134
3.3.4.	<i>Extracción de preferencias de selección</i> .....	140
3.4.	EVALUACIÓN DE LA ASOCIACIÓN DE SELECCIÓN EN EL DCECR.....	146
<b>4.</b>	<b>APLICACIÓN DESARROLLADA</b> .....	<b>165</b>
4.1.	INTRODUCCIÓN.....	165
4.2.	ARQUITECTURA DEL SISTEMA.....	166
4.3.	ANÁLISIS DE CASOS DE USO.....	168

---

4.4.	COLEXWEB: APLICACIÓN WEB CLIENTE .....	169
4.4.1.	<i>Análisis de palabra</i> .....	171
4.4.2.	<i>Análisis de texto</i> .....	174
4.5.	ELEMENTOS PARA ACELERAR EL ANÁLISIS DE TEXTOS .....	177
<b>5.</b>	<b>CONCLUSIONES</b> .....	<b>181</b>
<b>6.</b>	<b>BIBLIOGRAFÍA</b> .....	<b>185</b>
	<b>ANEXO</b> .....	<b>191</b>



# 1. Introducción

## 1.1. Introducción

El término colocación es utilizado en lingüística para designar un tipo de combinaciones de palabras usadas con relativa frecuencia en una determinada lengua que, sin llegar a ser combinaciones libres, tampoco constituyen locuciones o expresiones idiomáticas. Ejemplos de colocaciones en español:

*tener apetito*  
*afrontar riesgos*  
*dar (una) voltereta*  
*poner (un) telegrama*  
*sufrir (un) desengaño*  
*vomitara insultos*  
*accidente aparatoso*  
*agenda apretada*  
*catarro impresionante*  
*charla animada*  
*comida succulenta*  
*enemigo acérrimo*  
*noticia terrible*  
*precio razonable*  
*registro civil*  
*vacaciones escolares*  
*izar (la) bandera*  
*enhebrar (la) aguja*  
*arribar (un) barco*  
*cardar (la) lana*  
*cariarse (los) dientes*  
*ensillar (un) caballo*  
*banco de peces*  
*acceso de tos*  
*rebanada de pan*  
*comer frugalmente*  
*asesinar brutalmente*  
*competir duramente*  
*disgustarse seriamente*

En este capítulo se presentan distintos elementos que definen y caracterizan las colocaciones léxicas en español. El objetivo del mismo radica en la recopilación de información

sobre el tema como punto de partida para la construcción de un reconocedor automático de colocaciones. El estudio se estructura de la forma que se expone a continuación. En los dos primeros puntos se hace una revisión de las diversas corrientes sobre el tema: el enfoque estadístico, y la interpretación léxico-semántica del fenómeno, respectivamente. En el tercer punto se tratan las características de las colocaciones en español extraídas de los trabajos recopilados. A continuación se recogen las diversas clasificaciones de colocaciones propuestas en la bibliografía, tanto desde el punto de vista de las categorías gramaticales implicadas como desde los aspectos semánticos. Una vez establecida la taxonomía colocacional, en el siguiente punto se exponen las características específicas de cada categoría. Por último, se aborda la cuestión de la explotación de los diccionarios de uso –Diccionario del Español Actual (DEA)<sup>1</sup>, el Diccionario de Uso del Español (DUE)<sup>2</sup>, Diccionario General de la Lengua Española de Vox y el Nuevo Diccionario de Voces de Uso Actual– como fuente de colocaciones y cómo las presentan los autores en los artículos con objeto de obtener material de partida que permitan detectar vínculos o relaciones semánticas entre los elementos que forman una colocación o un grupo de ellas, como el caso:

*“dar se utiliza con sustantivos de golpe”*

## **1.2. Enfoque estadístico.**

Los primeros estudios sobre las colocaciones según los autores consultados se deben a Firth (1957), quien introduce el término **collocation** para referirse a la coaparición de palabras. Sin embargo, Gloria Corpas apunta trabajos previos de autores como Saussure, Bally o Porzig que ya tratan cuestiones estrechamente relacionadas con este fenómeno (Corpas, 1996). A partir de los trabajos de Firth surge la llamada escuela sistémica británica a la que pertenecen Halliday, Sinclair, Jones. Para tal corriente, las colocaciones se reducen a combinaciones frecuentes de palabras. Por utilizar el criterio de la frecuencia como rasgo principal de las colocaciones, en sus trabajos se propone el uso de técnicas estadísticas para la exploración de textos en busca de combinaciones de palabras cuya frecuencia de aparición pueda considerarse superior a la de otras combinaciones. Los cálculos se realizan a partir de las llamadas **líneas de concordancia**. Una línea de concordancia es un determinado fragmento en el que aparece una palabra dada en un texto o corpus textual. Por ejemplo, si se considera el texto:

“**vio** los sagaces ojos que bajo el ala de ancho sombrero de terciopelo viejo resplandecían; **vio** la mano morena y acerada que empuñaba una vara verde, y el ancho pie que, al moverse, hacía sonajear el hierro de la espuela...”(Doña Perfecta, Benito Pérez Galdós).

---

<sup>1</sup> *Diccionario del Español Actual*, M. Seco y otros, Madrid, Aguilar, 1999. Dos volúmenes.

<sup>2</sup> *Diccionario de Uso del Español*, [M. Moliner], 2 vols., Gredos, Madrid, 1966.

Las líneas de concordancia o contextos de la palabra **vio** son:

**vio** los sagaces ojos que bajo

sombrero de terciopelo viejo resplandecían; **vio** la mano morena y acerada

Cada línea de concordancia corresponde con una aparición de la palabra en el texto –permite el conteo para cálculos de frecuencias. Para extraerlas es necesario fijar la palabra bajo estudio y cuántas palabras antes y después de la misma conformarán el contexto –en el ejemplo se fijó en 5. En general, se considera una cantidad fija,  $n$ , que recibe el nombre de **distancia colocacional (collocation span)**. Sinclair y Jones estudiaron qué valor se le debe asignar y propusieron para  $n$  el valor de 4 (Koike, 2001). Sin embargo, en este mismo trabajo Koike presenta un ejemplo en el que se muestra que la distancia colocacional no tiene por qué ser 4:

Ambos juristas son expertos en la materia y han viajado a Madrid a fin de que la **petición** que en su momento pueda **cursar** el magistrado español sea acorde a la legislación británica (EPI, 27/10 al 2/1198:2).

En este caso la colocación *cursar petición* está separada por una cantidad de palabras mayor que 4. Otros autores han estudiado esta cuestión, Rogghe propone la distancia de 4 –salvo que sean adjetivos que será de 2–, Haskel, +/-3, Al Madi +/-1, pero que puede llegar a 5 ó 6 según los casos, Smadja propone 5 y Miall considera una distancia media de 9, aunque en ocasiones es necesario llegar hasta 15. (Corpas, 1996).

Las técnicas descritas en las líneas anteriores forman parte de la disciplina conocida como Lingüística del Corpus. Además de la distancia colocacional, en esta disciplina también son objeto de estudio las medidas de lo estrechamente relacionadas que estén los dos elementos de la colocación, como por ejemplo: **frecuencia de aparición de A con B**, **z-score** e **información mutua**. Aunque no son los únicos indicadores, sí son los más comunes: miden la frecuencia de aparición conjunta de palabras, pero no son capaces de reflejar aspectos más complejos presentes en las colocaciones. Existen trabajos que echan mano de técnicas estadísticas más avanzadas como el análisis factorial para intentar captar tales intimidades, pero requieren que se apliquen a una gran cantidad de textos para garantizar la fiabilidad de los resultados.

El ejemplo de la Tabla 1-1 (Koike, 2001), recoge los datos sobre la coocurrencia de *medida* y *drástica* en un corpus e intenta averiguar si *medida drástica* es una colocación. Indica que *drástica* aparece con *medida*, pero *medida* con *drástica* no.

<i>medida drástica</i>	
	<b>Frecuencia</b>
medida	250
drástica	8
coocurrencia	4
drástica con medida	$\frac{4}{8} \times 100 = 50\%$
medida con drástica	$\frac{4}{250} \times 100 = 1'6\%$

Tabla 1-1. Ejemplo del cálculo de la frecuencia de coocurrencia.

A lo que Koike llama **frecuencia de aparición** de *A* con *B*, no tiene por qué presentar el mismo valor que la frecuencia de aparición de *B* con *A*; en realidad, es un indicador del porcentaje de veces que una palabra aparece con otra respecto al número total de veces que aparece. Si este valor es elevado, quiere decir que si aparece *A* es muy posible que también aparezca *B*. Para Koike es suficiente indicio de colocación que el valor sea superior al 20%.

Las medidas del **z-score** e **información mutua** constituyen valores más complejos y precisos utilizados para la extracción de colocaciones en Lingüística de Corpus. Este tema será tratado en profundidad en el capítulo 2, donde se analizarán estos y otros indicadores que permiten detectar colocaciones de forma automática en textos electrónicos.

### 1.3. Enfoque léxico-semántico.

Trabajos posteriores en la materia –Benson, Haussman, Mitchell, Coseriu, Cruse, Haensch o Mel'čuk– tratan las colocaciones como un fenómeno más complejo que la mera combinación usual de palabras. Apuestan porque las colocaciones ponen de manifiesto relaciones sintagmáticas o semánticas entre palabras.

Mitchell considera que tales vínculos se producen entre lexemas (Koike, 2001). Benson distingue entre 1) **colocaciones gramaticales** que se componen de una palabra dominante seguida de una preposición o estructura sintáctica –*consistir en*– y 2) **colocaciones léxicas** que se forman exclusivamente con nombres, adjetivos, verbos y adverbios –entre sus aportaciones a la teoría colocacional figura la primera clasificación establecida en la materia (Carballo, 1998), (Alonso, 1997) y (Koike, 2001).

Coseriu define las **solidaridades léxicas** que coinciden parcialmente con el concepto de colocación. Una solidaridad léxica es una relación orientada en la que existe un **lexema determinante** y un **lexema determinado** –en el sentido de que el determinado implica al determinante–, pero no a la inversa. Por ejemplo, *caballo alazán* es una solidaridad léxica, porque *alazán* implica a *caballo*, pero no el recíproco. Una **solidaridad léxica unilateral** es del



tipo *morder dientes*, *lamer lengua* –\*\*?no suelen coocurrir. Sin embargo, *caballo alazán*, *ladrar perro*... que pertenecen al conjunto de las solidaridades léxicas multilaterales sí son colocaciones. Por otra parte, *radicalmente opuesto* es una colocación pero no es una solidaridad, puesto que no hay ninguna implicación en los lexemas que la forman, por ello no se puede afirmar que las colocaciones coincidan plenamente con las solidaridades léxicas, sino que algunas colocaciones coinciden con las solidaridades léxicas multilaterales (Corpas 1996, Koike 2001).

De especial importancia es la aportación de Hausman: llevó a cabo una taxonomía de las colocaciones y construyó una teoría sobre el estatus semántico de los elementos que intervienen en las colocaciones –definió **base** y **colocativo**. Para referirse a los dos lexemas que coocurren en una colocación, los lingüistas que tratan del tema en la actualidad usan el par de términos. Cada uno juega un papel distinto en la colocación, la base selecciona una acepción especial en el colocativo, que no puede ser descodificado sin la presencia de la base. Este rol diferenciado de los elementos que conforman la colocación se ilustra en la Tabla 1-2

Colocación	base	colocativo	Acepción del colocativo seleccionada
<i>levantar veda</i>	<i>veda</i>	<i>levantar</i>	<b>19. tr.</b> Hacer que cesen ciertas penas, prohibiciones o vejámenes impuestos por autoridad competente. <i>Levantar el entredicho, el destierro, el arresto, el embargo</i>
<i>café fuerte</i>	<i>café</i>	<i>fuerte</i>	<b>9. adj.</b> Dicho de un color o de un sabor: <b>intenso</b> .
<i>lógica aplastante</i>	<i>lógica</i>	<i>aplastante</i>	<b>2. adj.</b> Abrumador, terminante, definitivo.

Tabla 1-2 Semántica de los elementos de la colocación

Mel'čuk desarrolla sus trabajos desde la perspectiva de las **funciones léxicas** definidas para representar las posibilidades combinatorias de un determinado lexema que expresan un cierto sentido. Al igual que una función matemática, una función léxica asocia un conjunto de valores al elemento al que se aplica o **elemento llave** y que se corresponde con la base de Hausmann –el valor que toma la función es el colocativo o un conjunto de ellos. El vínculo establecido entre los elementos que forman una colocación determina la función que se debe escoger; se permiten utilizar funciones compuestas como combinación de varias funciones léxicas para expresar relaciones más complejas. Algunas funciones léxicas con los sentidos que expresan se recogen en la Tabla 1-3.

<b>Función</b>	<b>Sentido</b>	<b>Ejemplo</b>
<i>Son</i>	Emitir sonido típico	<i>Son(gato) = maullar</i>
<i>Magn</i>	Muy, Intenso	<i>Magn(enemigo) = acérrimo</i>
<i>Mult</i>	Conjunto de	<i>Mult(oveja) = rebaño</i>
<i>Fact<sub>i</sub></i>	Realizarse, Llevarse a cabo	<i>Fact<sub>0</sub>(sospecha)=confirmarse, corroborarse</i>
<i>Bon</i>	Bueno	<i>Bon(clima) = benigno</i>
<i>Ver</i>	Tal como debe ser	<i>Ver(cuchillo) = afilado</i>
<i>Sign</i>	Porción de	<i>Sign(pan) = tostada</i>
<i>Func<sub>i</sub>(C)</i>	Verbo semánticamente vacío que toma C como sujeto gramatical e i como su primer complemento	<i>Func<sub>i</sub>(viento) = soplar</i>
<i>Oper<sub>i</sub>(C)</i>	Verbo semánticamente vacío que toma C como sujeto gramatical e i como su primer complemento	<i>Oper<sub>i</sub>(atención) = prestar</i>
<i>Incep</i>	Comenzar	<i>IncepFunc<sub>0</sub>(epidemia)=declararse</i>

Tabla 1-3. Ejemplos de funciones léxicas.

Las funciones léxicas constituyen una herramienta útil a la hora de representar las relaciones comunes que conectan distintos lexemas presentes en las colocaciones; sin embargo, entre las críticas que reciben se mencionan aquí las siguientes:

1. Pueden producir combinaciones que no lo son, como en el caso de: *Sing(biblioteca) = libro*.
2. Existen colocaciones que no pueden ser descritas con las funciones léxicas, como en el caso de *discusión bizantina* (Koike, 2001).

#### 1.4. Las colocaciones en español.

Los estudios sobre colocaciones en español se producen a partir de la década de los 90, si bien a finales de los 70 Seco introduce el concepto de **contorno de la definición lexicográfica** que hace referencia al término colocación (Corpas, 1996), (Penadés, 2001) y (Castillo, 1998).

Corpas las considera un tipo de **unidades fraseológicas**<sup>3</sup>.

Bosque las supone una manifestación de **selección léxica**<sup>4</sup>, en las que los colocativos que actúan como predicados seleccionan sus argumentos, no de forma individual, sino grupal

<sup>3</sup> Una unidad fraseológica es una construcción formada por al menos dos palabras ortográficas —puede ser una oración compuesta— con las siguientes características potenciales: uso habitual, presenta cierta fijación, idiomática —su significado no es composicional—, presenta variantes. Todos estos aspectos no aparecen en la misma medida en los distintos tipos de unidad fraseológica: colocaciones, locuciones, y enunciados fraseológicos (paremias y fórmulas rutinarias), sino que se dan en distinto grado (Corpas, 1996).

<sup>4</sup> Las colocaciones constituyen una manifestación de selección léxica. Los predicados —verbales, adjetivales, adverbiales o preposicionales— seleccionan sus argumentos, y al hacerlo restringen el

---

–del mismo, o de varios campos semánticos– que conforman lo que se denomina **clase léxica**. De esta manera, lo realmente importante radica en detectar la clase que selecciona un determinado colocativo; entre los ejemplos que presenta en su trabajo aparecen:

*supino* {ignorancia, incompetencia, inutilidad, necedad, desconocimiento, estupidez, ridiculez, imbecilidad, irresponsabilidad, egoísmo, cinismo}.

*terminantemente* {prohibir, negar, desmentir, excluir, rechazar, oponerse}

En algunos casos establece la relación entre el colocativo y algún rasgo definitorio de la clase, como en el caso:

*“universalmente + verbos que denotan aceptación”*

Zuluaga se refiere a las colocaciones como un tipo de combinación de al menos dos lexemas, en las que puede aparecer algún elemento gramatical y de las que no se puede decir que sean combinaciones libres, ni tampoco unidades fraseológicas, pero que comparten propiedades de unas y otras.

Por último citamos a Alonso que estudia las colocaciones en el marco de la teoría Sentido-Texto de Mel’čuk.

Si bien se aprecian distintas interpretaciones del fenómeno colocacional entre los lingüistas españoles, aparece como un denominador común de sus trabajos el que descartan la interpretación exclusivamente estadística de las colocaciones por presentar los siguientes problemas:

1. No todas las combinaciones frecuentes cuentan como colocaciones.
2. Puede que no se considere colocación aquella que forman unas palabras de baja frecuencia.
3. El uso exclusivo de la frecuencia no puede manifestar fenómenos semánticos presentes en las mismas.

---

conjunto de entidades que pueden denotar en función de unos rasgos semánticos — a veces muy abiertos y en otras considerablemente restringidos. Frente a lo que parece sugerir una buena parte de la bibliografía, en muchísimos casos las unidades seleccionadas por los predicados no son piezas léxicas aisladas, sino clases léxicas de mayor o menor intensidad. (Bosque, 2001).

## 1.5. Características.

En los siguientes epígrafes se ahonda en los rasgos específicos que resultan fundamentales a la hora de discriminar si una combinación de lexemas constituye o no una colocación. En general, existe una serie de características propias de las colocaciones, algunas de ellas diferenciadoras de otros tipos de combinaciones y otras que comparten al menos en parte debido a la transición gradual **combinación libre**→**colocación**→**locución** (Zuluaga, 2002), (Blasco, 2002), (Penadés, 2001) y (Corpas, 1996) –he aquí una de las causas por las que es difícil decidir si una determinada construcción es o no una colocación. Cabe destacar las discrepancias entre autores a la hora de catalogar como colocaciones determinadas construcciones; por ejemplo, pueden ser controvertidos los siguientes casos (Bosque, 2001):

*ejecutar un castigo*  
*empantanarse un proyecto*  
*enfilarse la calle principal*  
*disminución progresiva...*

Penadés advierte que el Diccionario fraseológico del español moderno (Varela, Kubarth, 1994) existen numerosas colocaciones incluidas como locuciones; entre otras:

*dar corte*  
*dar el cambiazco*  
*dar el espectáculo*  
*dar el golpe.*

### 1.5.1. Flexibilidad

En contraposición con las estructuras por lo general fijas de las locuciones (Koike, 2001), la primera propiedad que se trata en este apartado se refiere a la **flexibilidad formal** de las colocaciones. Admiten<sup>5</sup>:

1. Cambios de categoría gramatical:

*comida frugal / comer frugalmente*

2. Modificación adjetival: En las colocaciones es lícito incluir un adjetivo que modifique al sustantivo, la colocación:

*hacer un aterrizaje*

también se puede usar:

*hacer un aterrizaje forzoso*

---

<sup>5</sup> Todos los ejemplos se obtuvieron de (Koike, 2001).

sin embargo, esto no es válido en las locuciones, por ejemplo, la locución:

*quemarse las pestañas*

no admite ser utilizada en la forma:

*quemarse las pestañas largas*

3. Transformación a pasiva. Podemos transformar a pasiva una colocación sustantivo + verbo:

*trasplantar órganos*

en pasiva:

*el órgano fue transplantado*

lo que no es posible en el caso de las locuciones:

*escurrir el bulto*

no se utiliza :

*el bulto fue escurrido*

4. Relativización:

*seguir la línea / Este libro marca la línea que deben seguir sus partidarios.*

*echar el ojo (a algo) / El ojo que acabo de echar a ese vestido.*

5. Nominalización: Es posible en las colocaciones como se observa en el ejemplo:

*trasplantar un órgano / el trasplante de órganos*

pero no en las locuciones:

*escurrir el bulto*

no se usa:

*el bulto que fue escurrido*

6. Pronominalización:

Asumió el **cargo** de alcalde, pero su repentina enfermedad le impidió **desempeñarlo** (desempeñar cargo)

De aquí que se afirme que la combinación realmente se da entre los lexemas<sup>6</sup> y no entre las lexías simples. Tal aseveración explica las colocaciones que se forman con una unidad léxica simple y una locución, como: *dar un golpe de Estado*. Por tanto, la exploración de un texto para la extracción de colocaciones no puede llevarse a cabo sobre palabras gráficas.

Por otra parte, los elementos que intervienen pueden sustituirse en general por otros relacionados, bien porque:

1. Pertenecen al mismo campo semántico:

*desempeñar* {un cargo, una función o un papel} (Corpas, 1996)

2. Son sinónimos:

{mantener, sostener} una *conversación*

*armar* {alboroto, algarabía, barullo}

3. Son hiperónimos, hipónimos:

*cariarse* {un diente, la dentadura}

*escamar* {pescado, besugo}

En virtud de esta característica se pueden describir propiedades combinatorias mediante definiciones del tipo: “*dar combina con sustantivos de golpe*”, o bien “*dar selecciona sustantivo de golpe*” que expresamos en forma de reglas:

“*dar + sustantivos de golpe*”

“*universalmente + verbos que denotan aceptación*”

### 1.5.2. Determinación.

Los elementos que forman una colocación reciben el nombre de colocados, uno de ellos goza de autonomía semántica y recibe el nombre de **base**, el significado del otro depende precisamente de ésta, está determinado por ella y se denomina **colocativo**. Como se mencionó con anterioridad, esta nomenclatura fue introducida por Hausmann en la terminología colocacional y en la actualidad su uso se ha generalizado entre todos los estudiosos del tema consultados. Respecto a las características señaladas de los dos colocados, se dice que la base es una palabra **autosemántica** –no necesita del colocativo para ser definida– y el colocativo es **sinsemántica**

---

<sup>6</sup> **lexema**. 1. m. Ling. Unidad mínima con significado léxico que no presenta morfemas gramaticales, por ejemplo, sol; o que, poseyéndolos, prescinde de ellos por un proceso de segmentación; por ejemplo, terr en enterráis. (DRAE).

–su definición necesita ser completada por la base. Esta característica es fundamental en las siguientes definiciones de colocación que aportan Alonso (1994) y García Platero (2002):

Una **colocación** o “semi-frasema” (FL)  $AB$  es una combinación de dos lexemas  $A$  y  $B$ , de tal forma que su significante es la suma regular de los significantes de los lexemas constituyentes y su significado incluye el significado del lexema  $A$  y un significado  $C$  que es:

a) Bien  $C \Leftrightarrow B$  y

- i.  $C$  es vacío: el lexema  $B$  es un auxiliar que se usa para sostener una configuración sintáctica

*dar un paseo...*

- ii.  $C$  no es vacío, pero el lexema  $B$  expresa  $C$  solo en combinación con  $A$  o con otros pocos lexemas similares

*odio mortal*

*interés vivo...*

b) Bien  $C = B$  y

El lexema  $B$  se selecciona de modo restringido: en combinación con  $A$  no se puede reemplazar por otro sinónimo

*café fuerte...*

c)  $C$  incluye el sentido de  $A$

*pelo rubio*

*nariz aguileña...*

Todos los casos posibles se pueden explicar mediante alguna de las posibilidades:

1. La base selecciona una acepción especial en el colocativo; generalmente se usa en sentido figurado y se reconoce su significado gracias a la base con la que aparece. Los tres primeros casos en la definición de Alonso se ajustan a este tipo. Es el caso de colocaciones como:

*banco de peces*

*precio astronómico*

*levantar sesión*

*locamente enamorado*

*dar envidia...*

Se han subrayado los colocativos. Se incluyen las definiciones que proporciona el DRAE para la acepción en que se utilizan tales términos:

**banco. 5. m.** Conjunto de peces que van juntos en gran número.

**astronómico. 2. adj. coloq.** Que se considera desmesuradamente grande. Sumas, distancias astronómicas.

**sesión. levantar la ~.**

**1.** Concluirla.

**locamente fr. 2. adv. m.** Excesivamente, sin prudencia ni moderación

**dar. 23. tr.** Causar, ocasionar, mover. *Dar gusto, gana.*

Se advierte en estas definiciones que los significados de *banco*, *astronómico*, *levantar*, *locamente* y *dar* están determinados por cada uno de los otros colocados; si no estuvieran presentes, no se sabría con qué acepción se utiliza el colocativo respectivo –tal comportamiento refleja la precisión semántica que sustenta la colocación.

2. No se puede definir el colocativo sin la base, pero en este caso lo que existe es una relación típica entre los dos elementos colocados. Los colocativos son términos que especifican propiedades, acciones o procesos específicos de las bases; por eso se dice que existe un vínculo típico entre ambos. Valen como ejemplos:

*triángulo isósceles*

*querochar abeja*

*pelo lacio*

*zarpar barco*

*trinchar carne...*

Así, *isósceles* es un tipo de triángulo y no existe ambigüedad respecto a este término. En tales ocasiones, por estar la base contenida en el colocativo, se consideran las colocaciones como casos de *solidaridades léxicas*.

**isósceles. V. triángulo isósceles**

**querochar. 1. intr.** Dicho de las abejas y de otros insectos: Poner la querocha.

Uno de los significados de *lacio* en el DRAE manifiesta que es una propiedad específica del *pelo*.

**lacio, cia. 3.** Dicho del cabello: Que cae sin formar ondas ni rizos.

Además de existir una relación típica entre los colocados, se observa que el significado de *lacio* queda determinado por la base.



El caso de *zarpar barco* es similar al anterior, puesto que la segunda acepción de este verbo requiere para su definición la presencia de la base de la colocación, y ésta a su vez selecciona una de las acepciones posibles, pero no en un sentido figurado.

**zarpar. 2.** intr. Dicho de un barco o de un conjunto de ellos: Salir del lugar en que estaban fondeados o atracados.

Lo mismo sucede con el último ejemplo:

**trinchar. 1.** Partir en trozos la comida para servirla.

Si los dos elementos que configuran una determinada combinación son vocablos autosemánticos estamos ante una combinación libre, en cambio, si una unidad es autosemántica y la otra sinsemántica se está ante una colocación. Los elementos que forman una locución son palabras extrasemánticas, porque el significado de la misma no proviene de los significados de sus constituyentes. (García Platero, 2002).

### 1.5.3. Frecuencia, Preferencia.

Las colocaciones son combinaciones usuales de palabras que los hablantes nativos de una lengua repiten tan habitualmente en situaciones comunicativas similares que las han convertido en estereotipos. Aunque teóricamente se podría escoger como colocación cualquier combinación posible, en muchas ocasiones ¿no son usuales o de uso preferente; así que no todas las posibilidades de elección –más o menos amplia– de los colocativos deberían considerarse como colocaciones, sino las que han sido sancionadas por la comunidad.

Así, *fervorosamente* y *fervientemente* se combinan de forma preferente, o casi exclusivamente con *rezar* (u *orar*), aunque los diccionarios suelen dar como sinónimos *vivamente*, *ardientemente*, o *vehementemente*. De este modo, *fervientemente* y *fervorosamente* se especializan en su empleo con *rezar*, en detrimento de otros adverbios de manera (García-Page, 2001).

En *lanzarse al ataque*, *lanzarse* ha sido preferido a *arrojarse*, y *ataque* ha sido preferido a *lucha*. (Zuluaga, 2002). En *guerra mundial*, *mundial* ha sido preferido a *internacional* (Zuluaga, 2002).

Tal perspectiva supone un problema a la hora de buscar generalizaciones sobre las propiedades combinatorias de algún lexema dado.

A pesar de que la elevada frecuencia de uso de una combinación respecto a otras posibles no sea un criterio determinante para asegurar que forme una colocación, constituye una

característica que se debe considerar. Por otra parte, es una consecuencia, en parte, del uso de clichés o combinaciones preferentes.

En los lenguajes de especialidad es normal el uso de combinaciones típicas de términos que derivan en colocaciones; por ejemplo

en informática:

*disco duro*  
*cargar el fichero*  
*eliminar la carpeta*  
*implementar un programa...*

o en derecho:

*{interponer, cursar, presentar} un recurso*  
*{abrir, instruir, tramitar} un proceso...*

En tales casos, los vínculos léxicos son más fácilmente identificables, porque el problema aparece delimitado; más difícil resultaría detectarlos en la lengua general.

#### **1.5.4. Composicionalidad.**

En cuanto al significado de una colocación se da un abanico de posibilidades que abarca desde la **composicionalidad** a la **seudocomposicionalidad** como se explica a continuación:

**Composicionalidad total:** sucede tanto en *triángulo isósceles* como en *querochar abeja*. El significado se obtiene de la suma de los significados de los dos colocados. La colocación se puede descodificar fácilmente y se dice que es totalmente transparente semánticamente. Nótese que tal característica aparece en las mismas colocaciones que se incluyen como ejemplos en el punto 2 del apartado Determinación.

**Seudocomposicionalidad o composicionalidad parcial.** En las colocaciones en las que la base selecciona un uso figurado del colocativo se necesita un conocimiento específico por parte del hablante para poder asignarles el significado; por tanto se dice que presentan cierta *opacidad semántica* (García Platero, 2002) y no son tan fácilmente descodificables.

Esta absoluta transparencia que atribuye al significado de las partes, no siempre se da con la misma intensidad; por ejemplo en *dinero negro*, si al hablante no nativo se le da ya hecha la combinación no podrá conocer exactamente su significado, si previamente no se le ha informado de que la base *dinero* en presencia de *negro* asigna al colocativo una acepción distinta de la que posee ('ilegal') (Castillo Carballo, 1998).



puesto que estos vínculos se podrían representar como:

*“triángulo + tipo de triángulo” {isósceles, escaleno, rectángulo, etc.}*

mientras que la relación en la segunda colocación corresponde a la estructura:

*“dar + sustantivo de estado físico o de ánimo” {hambre, sueño, celos, etc.}*

Debido a tal gradualidad, resulta difícil delimitar dónde terminan las combinaciones libres y dónde empiezan las colocaciones, o dónde terminan éstas y empiezan las locuciones; por tanto, no es extraño que se catalogue alguna combinación en un grupo erróneo, lo cual no quiere decir que las características de cada grupo no aparezcan bien claras y precisas (Penadés, 2001), (Bosque, 2001) y (Zuluaga, 2002).

## **1.6. Clasificación.**

En cuanto a las variedades existentes se pueden clasificar las colocaciones según distintos criterios: la categoría gramatical de los colocados, el grado de restricción entre los colocados y el significado de la colocación.

### **1.6.1. Categoría gramatical de los colocados.**

Varias clasificaciones se pueden citar en este apartado: la de Hausmann y la de Benson, o las clasificaciones de colocaciones en español propuestas por Corpas, Castillo Carballo, Penadés y K. Koike, con diferencias mínimas.

La primera taxonomía propuesta se debe a Hausmann, quien la estableció basándose en la función gramatical de los elementos exclusivamente, se presenta en la Tabla 1-4. La tipología de Benson de la Tabla 1-5 se define sobre las **colocaciones léxicas** que se forman con: *nombres, adjetivos, verbos y adverbios* a diferencia de las **colocaciones gramaticales** que siguen el esquema:

*(nombre, adjetivo, adverbio) + preposición + estructura sintáctica*

Se puede apreciar que para delimitar alguno de los grupos se apoya también en características semánticas de alguno o ambos colocados.

<b>Tipo</b>	<b>Base</b>
sustantivo + adjetivo	sustantivo
sustantivo + verbo	sustantivo
verbo + sustantivo <sub>objeto</sub>	sustantivo
verbo + adverbio	verbo
adjetivo + adverbio	adjetivo
sustantivo + preposición + sustantivo	sustantivo

Tabla 1-4. Clasificación de Hausmann.

	<b>Tipo</b>	<b>Base</b>
<b>Gramaticales</b>		
	verbo (‘creación’ o ‘activación’) + nombre / pronombre	sustantivo
	verbo (‘erradicación’, ‘anulación’) + nombre	sustantivo
	adjetivo + nombre	sustantivo
<b>Léxicas</b>	nombre + de + nombre (uno representa a la unidad y otro el conjunto).	sustantivo
	adverbio + adjetivo	adjetivo
	verbo + adverbio	verbo

Tabla 1-5. Clasificación de Benson.

Las siguientes propuestas que se recogen se refieren al español. Nótese que en las tablas correspondientes se resaltan en negrita las diferencias; por ejemplo, en la clasificación de Corpas de la Tabla 1-6 aparece la clase:

*verbo + sustantivo<sub>objeto</sub>*

que Carballo en la Tabla 1-7 la muestra más detallada:

*verbo + sustantivo<sub>objeto</sub>*

*verbo + preposición + sustantivo*

también se marca la entrada:

*adjetivo + adverbio*

puesto que las dos autoras difieren en el orden

*adjetivo + adverbio* para Corpas

*adverbio + adjetivo* para Carballo

En lo esencial, ambas se pueden considerar la misma clasificación (Koike, 2001) –se aprecia que están basadas en las de Hausmann y Benson.

<b>Tipo</b>	<b>Base</b>	<b>Ejemplos</b>
sustantivo <sub> sujeto</sub> + verbo	sustantivo	<i>correr un rumor</i> <i>estallar una guerra</i> <i>zarpar un barco</i>
verbo + sustantivo <sub> objeto</sub>	sustantivo	<i>desempeñar un cargo</i> <i>zanjar un desacuerdo</i> <i>conciliar el sueño</i> <i>dar comienzo</i>
adjetivo + sustantivo	sustantivo	<i>fuentes fidedignas</i> <i>enemigo acérrimo</i> <i>éxito fulgurante</i> <i>oído fino</i>
sustantivo + sustantivo	sustantivo	<i>hombre clave</i>
sustantivo+preposición+sustantivo	sustantivo	<i>tableta de chocolate</i> <i>enjambre de abejas</i> <i>ronda de negociaciones</i>
verbo + adverbio	verbo	<i>caer pesadamente</i> <i>negar rotundamente</i> <i>fracasar estrepitosamente</i>
adjetivo + adverbio	adjetivo	<i>firmemente convencido</i> <i>estrechamente ligado</i> <i>rematadamente loco</i> <i>diametralmente opuesto</i>

Tabla 1-6. Clasificación de Corpas.

<b>Tipo</b>	<b>Base</b>	<b>Ejemplos</b>
sustantivo <sub>sujeto</sub> + verbo	sustantivo	<i>correr un rumor</i> <i>estallar una guerra</i> <i>zarpar un barco</i>
verbo + sustantivo <sub>objeto</sub>	sustantivo	<i>desempeñar un cargo</i> <i>dar comienzo</i> <i>tener repercusión</i>
verbo + preposición + sustantivo	sustantivo	<i>poner en cuestión</i>
adjetivo + sustantivo	sustantivo	<i>fuentes fidedignas</i>
sustantivo + sustantivo		<i>enemigo acérrimo</i> <i>éxito fulgurante</i> <i>oído fino</i> <i>hombre clave</i>
sustantivo <sub>1</sub> +preposición+sustantivo <sub>2</sub>	sustantivo	<i>tableta de chocolate</i> <i>enjambre de abejas</i> <i>ronda de negociaciones</i>
verbo + adverbio	verbo	<i>caer pesadamente</i> <i>negar rotundamente</i> <i>fracasar estrepitosamente</i>
adverbio + adjetivo	adjetivo	<i>firmemente convencido</i> <i>estrechamente ligado</i> <i>rematadamente loco</i> <i>diametralmente opuesto</i>

Tabla 1-7. Clasificación de Castillo Carballo.

En la Tabla 1-8, Penadés revisa las colocaciones en las que interviene un verbo; al igual que Koike de la Tabla 1-9 incluye las de:

*verbo + adjetivo*

	<b>Subtipo</b>	<b>Base</b>	<b>Ejemplos</b>
verbo + SN	verbo + nombre	sustantivo	<i>correr un rumor</i> <i>estallar una guerra</i> <i>zarpar un barco</i>
	verbo + determinante + nombre		<i>desempeñar un cargo</i> <i>dar comienzo</i>
verbo + SP	verbo+preposición + determinante+ nombre	sustantivo	<i>poner en cuestión</i>
verbo + adverbio		verbo	<i>caer pesadamente</i> <i>negar rotundamente</i> <i>fracasar</i> <i>estrepitosamente</i>
verbo + adjetivo		adjetivo	<i>resultar ileso</i> <i>salir malparado</i>
Estructura comparativa			<i>dormir como un lirón</i> <i>dar vueltas como un trompo</i>

Tabla 1-8. Clasificación de Penadés.

y propone considerar las que reproducen una estructura comparativa, que para Koike es una colocación compleja del tipo:

*verbo + locución adverbial*

La propuesta de Koike recoge las anteriores pero detalla algunos subgrupos; también abarca las colocaciones complejas no contempladas en otras tipologías, al menos específicamente. Como novedad, en esta clasificación se considera un único tipo

*sustantivo + verbo*

pero que se subdivide en:

*sustantivo<sub>sujeto</sub> + verbo, verbo + sustantivo<sub>CD</sub>*

y en:

*verbo + preposición + sustantivo.*



	<b>Tipo</b>	<b>Subtipo</b>	<b>Base</b>	<b>Ejemplos</b>
<b>Simples</b>	sustantivo+verbo	sustantivo <sub>sujeto</sub> + verbo	sustantivo	<i>correr un rumor</i> <i>estallar una guerra</i> <i>zarpar un barco</i>
		verbo + sustantivo <sub>CD</sub>		<i>desempeñar un cargo</i> <i>dar comienzo</i>
		verbo+preposición+ sustantivo		<i>poner en cuestión</i>
	adjetivo+ sustantivo		sustantivo	<i>fuelle fidedigna</i> <i>enemigo acérrimo</i> <i>oído fino</i>
	sustantivo+de+ sustantivo		sustantivo	<i>tableta de chocolate</i> <i>enjambre de abejas</i>
	verbo + adverbio		verbo	<i>caer pesadamente</i> <i>negar rotundamente</i>
<b>Complejas</b>	adverbio+ adjetivo / participio		adjetivo	<i>firmemente convencido</i> <i>rematadamente loco</i>
	verbo + adjetivo		adjetivo	<i>resultar ileso</i> <i>salir malparado</i>
	verbo+locución nominal			<i>dar un golpe de estado</i> <i>levantar castillos en el aire</i>
	locución verbal+SN			<i>poner en juego un recurso</i> <i>pasar a limpio un borrador</i>
	sustantivo+ locución adjetival			<i>dinero contante y sonante</i>
	verbo+locución adverbial			<i>reírse a mandíbula batiente</i> <i>llorar a moco tendido</i>
adjetivo+locución adverbial			<i>bruto como un arado</i> <i>sano como una manzana</i>	

Tabla 1-9. Clasificación de Koike

### 1.6.2. Grado de restricción de los colocados.

Según el grado de restricción Cowie diferencia entre colocaciones libres, restringidas, estables y categoría puente (Corpas, 1998). El grado de restricción va aumentando gradualmente en el orden citado. El significado en las **colocaciones libres** es totalmente composicional y pueden aparecer en numerosas colocaciones. En las **colocaciones restringidas** se está ante una acepción especial del colocativo que selecciona la base. El campo colocacional en las **estables** es mínimo, se limita a uno o dos colocados. Por último, el uso figurado del colocativo y el que quepa la sustitución sinonímica señalan las características de la **categoría puente**<sup>7</sup> (Tabla 1-10).

<b>Libres</b>	<i>provocar una pelea / guerra / discusión</i> <i>empezar una pelea / guerra / discusión</i> <i>ganar una pelea / guerra / discusión</i>
<b>Restringidas</b>	<i>correr peligro / suerte / aventuras</i>
<b>Estables</b>	<i>conciliar el sueño</i>
<b>Categoría Puente</b>	<i>levantar una calumnia / una mentira</i> <sup>8</sup>

Tabla 1-10. Clasificación de Cowie según el grado de restricción.

En la misma línea, Koike establece dos categorías: **colocaciones amplias** y **colocaciones estrechas** como aparecen en la Tabla 1-11. En las primeras, el colocativo tiene la capacidad de combinarse con numerosas bases debido a que su significado es general y polisémico en la mayoría de las ocasiones. En las segundas, el vínculo entre los dos lexemas es fuerte, las bases obligatoriamente tendrán que ser compatibles con el significado específico del colocativo<sup>9</sup>. Las colocaciones amplias abarcan las categorías libres y restringidas de Cowie y las estrechas recogen la estable y la categoría puente.

<b>Colocaciones estrechas</b>	<b>Colocaciones amplias</b>
<i>escrutar votos</i>	<i>recontar {votos, reses, presos}</i>
<i>extirpar quiste</i>	<i>quitar {quiste, muela}</i>
<i>asestar puñetazo</i>	<i>dar {puñetazo, consejo, paseo}</i>

Tabla 1-11. Colocaciones amplias y estrechas de Koike.

Una distinción similar propone García Platero al hablar de colocaciones libres y restringidas. En las primeras existe una gran cantidad de elementos que pueden combinarse con la base, mientras que en las restringidas como:

<sup>7</sup> Ejemplos recogidos del Manual de Fraseología

<sup>8</sup> Levantar una calumnia es una colocación, al contrario que levantar una mentira que no lo es.

<sup>9</sup> Ejemplos recogidos de (Koike, 2001).

*fruncir el ceño*

el colocativo es difícilmente sustituible.

### 1.6.3. Semántica.

El significado de la colocación también lleva a Koike a proponer dos categorías diferenciadas dentro de los tipos:

*sustantivo + verbo*

y

*sustantivo + adjetivo*

que son: las **funcionales** y las **léxicas**.

**Colocaciones funcionales:** el colocativo pierde su significado y pasa a tener un valor funcional. En el caso de las *sustantivo + verbo* de la Tabla 1-12, el verbo habilita al sustantivo para funcionar como un verbo complejo, aporta la información de número, persona y tiempo, mientras que el sustantivo aporta el significado. En general, existe un verbo simple relacionado morfológicamente con el sustantivo, aunque no siempre.

<b>Funcional verbo+ sustantivo</b>	<b>Verbo relacionado</b>
<i>tener miedo</i>	<i>temer</i>
<i>hacer una aclaración</i>	<i>aclarar</i>
<i>dar un golpe</i>	<i>golpear</i>
<i>albergar esperanza</i>	<i>esperar</i>

Tabla 1-12. Colocaciones funcionales: verbo + sustantivo.

En las *sustantivo+adjetivo*, el adjetivo intensifica cuantitativa o cualitativamente al sustantivo como se muestra en la Tabla 1-13.

<b>Funcionales sustantivo+adjetivo</b>	
<i>error garrafal</i>	cuantitativa
<i>módico precio</i>	cuantitativa
<i>dolor atroz</i>	cuantitativa
<i>comida atroz</i>	cualitativa
<i>tiempo horroroso</i>	cualitativa
<i>chico fenomenal</i>	cualitativa

Tabla 1-13. Colocaciones funcionales: adjetivo + sustantivo.

**Colocaciones léxicas:** Los elementos que conforman la colocación conservan su significado léxico plenamente, de tal forma que en tal caso se está ante colocaciones totalmente composicionales, como aparecen en la Tabla 1-14.

Verbo + sustantivo	Sustantivo + adjetivo
<i>moler café</i>	<i>pelo lacio</i>
<i>cruzarse de brazos</i>	<i>poder adquisitivo</i>
<i>zarpar barco</i>	<i>barrio céntrico</i>
<i>amasar fortuna</i>	<i>cuchillo afilado</i>
<i>afrontar riesgo</i>	<i>gobierno autoritario</i>

Tabla 1-14. Colocaciones léxicas.

En las colocaciones *sustantivo + verbo* hay que incluir las **aspectuales**; se llaman así porque su significado indica alguno de los siguientes aspectos: *incoativo* (comienzo), *durativo-reiterativo* (duración), o *terminativo-resultativo* (final). A diferencia de las funcionales, el verbo tiene en tal caso una cierta carga semántica; por ejemplo, en

*asaltar duda*

*asaltar* indica acometer de pronto algo, en este caso una *duda*.

Aspecto	Colocación
Incoativo	<i>asaltar duda</i>
durativo-reiterativo	<i>sostener conversación</i>
terminativo-resultativo	<i>cerrar sesión</i>

Tabla 1-15. Colocaciones aspectuales.

Sobre la base de estos grupos, y recurriendo a la definición de Alonso se consideran los siguientes casos:

**Colocativo con valor plenamente funcional.** Se está ante un verbo deslexicalizado<sup>10</sup> y un sustantivo generalmente deverbal, o bien ante un adjetivo que juega un papel meramente intensificador en alguno de los siguientes sentidos de: muy grande, magnífico, extraordinario, etc. Se puede considerar dentro de este grupo el caso ‘a). i’ de la definición de Alonso.

*dar un golpe* | *golpear*

**DRAE:** *golpear*. Dar un golpe o golpes repetidos.

<sup>10</sup> Un verbo deslexicalizado contribuye mínimamente al significado de la colocación, su función principal es aportar el tiempo y número a la construcción que conjuntamente funciona como un verbo cuyo significado proviene del significado del sustantivo.

*tener miedo* | *temer*

**DRAE. temer.** Tener a una persona o cosa por objeto de temor.

*tener esperanzas* | *esperar*

**DRAE. esperar.** Tener esperanza de conseguir lo que se desea.

*hacer una aclaración* | *aclarar*

**DRAE. aclarar. 10.** Hacer clara, perceptible, manifiesta o inteligible alguna cosa, ponerla en claro, explicarla.

*soberana paliza*

**DEA. soberano –na. 2.** Superior o extraordinario. Frec con intención ponderativa.

*dolor atroz*

**DEA. atroz. 2.** Muy grande o extraordinario.

*inteligencia colosal*

**DEA. colosal. 2.** Magnífico o extraordinario.

**Vínculo típico.** El colocativo expresa alguna propiedad, acción, o proceso característico de la base. Es el caso que aparece en el punto ‘c)’ de la definición de Alonso.

*zarpar barco*

**DEA. zarpar. 1.** Salir [un barco] del lugar en que estaba fondeado o atracado.

*libar abeja*

**DEA. libar. 1.** Chupar [un insecto el néctar, u otra sustancia azucarada, de las flores].

*enebrar aguja*

**DEA. enebrar. 1.** Pasar la hebra por el ojo [de la aguja (cd)].

*dilapidar fortuna*

**DEA. dilapidar** Derrochar o despilfarrar.

*iglesia católica*

**DEA. católico. 1. b)** De los católicos o de la Iglesia católica.

*pelo lacio*

**DEA. lacio. 3.** [Cabello] que cae liso y blando.

**Selección de acepción en el colocativo.** Abarca los casos ‘a).ii’ y ‘b)’ de la definición de Alonso.

*albergar esperanza*

**DEA. albergar. 4.** Tener de manera estable [un sentimiento].

*administrar anestesia*

**DEA. administrar. 4.** Dar o aplicar [algo, esp. medicamentos].

*asaltar la duda*

**DEA. asaltar. 3. (col)** Acometer repentinamente [a alguien una idea o un sentimiento].

*perder el tren*

**DEA. perder. 7.** No llegar a tiempo de poder utilizar [un medio de transporte (*cd*)].

*tocar el piano*

**DEA. tocar. 6. b)** Hacer sonar [un instrumento musical o un dispositivo de señales acústicas].

*amasar una fortuna*

**DEA. amasar. 3.** Acumular o reunir [bienes o dinero].

*calar melón*

**DEA. calar. 2.** Cortar un pedazo [de un melón o de otra fruta (*cd*)] para probarlos.

*color crudo*

**DEA crudo. 8.** [Color]. blanco o amarillento.

*alimentación equilibrada*

**DEA. equilibrada. 2.** Que tiene equilibrio [4, 5 y 6].

**equilibrio. 5.** Justa proporción de los elementos [de algo].

*cuestión vital*

**DEA. vital. 2. b)** De suma importancia.

## **1.7. Características por categorías.**

En los apartados siguientes se exponen las características propias de los distintos tipos de colocaciones que se recopilaron de la bibliografía analizada. Los epígrafes se corresponden con cada una de las categorías según la función gramatical de los colocados, si bien en cada uno se alude a las categorías establecidas según la semántica, en caso de ser necesario.

### **1.7.1. Sustantivo + Verbo.**

La información sobre las colocaciones

*sustantivo + verbo*

es más numerosa que las restantes, debido a que son las más abundantes y también las más estudiadas del español.

Los verbos que forman parte de colocaciones se denominan colocacionales, los que no, se llaman no colocacionales; figuran en la Tabla 1-16 (Koike, 2001). Los colocacionales<sup>11</sup>, a su vez se dividen en:

**Funcionales:** No mantienen su significado léxico. Por ejemplo, *hacer* o *sacar* en:

*hacer una aclaración*

*sacar la conclusión...*

**Léxicos:** Mantienen su significado léxico. Por ejemplo

*afinar el piano...*

Según con qué se combine un mismo verbo puede tener valor funcional o léxico. Por otra parte, tanto unos como otros pueden ser generales o específicos.

**Generales:** Los verbos generales tienen un número elevado de acepciones y forman numerosas colocaciones. El campo colocacional es amplio.

**Específicos:** Forman colocaciones estrechas que reflejan relaciones típicas entre el verbo y el sustantivo.

<b>Colocacionales</b>	<b>Funcionales</b>	Generales	<i>dar, tener, hacer, sentir, poner, echar</i>
		Específicos	<sup>12</sup> <i>abrigar, albergar, fraguar, brindar, tejer, tramar, etc.</i>
	<b>Léxicos</b>	Generales	<i>tocar, gastar</i>
		Específicos	<i>tañer, dilapidar</i>
<b>No colocacionales</b>	<i>querer, pensar, indicar, afirmar, preguntar, verbos de decir, asegurar, advertir, vender, comprar, agradecer, obligar, invitar, gustar, parecer, interesar, considerar, creer, hallar, poder, saber, querer, intentar, etc.</i>		

Tabla 1-16. Clasificación de los verbos.

En cuanto a los sustantivos, cabe la división entre concretos y abstractos de la Tabla 1-17 (Koike, 2001), ya que los primeros son los que intervienen en las colocaciones léxicas y los segundos los que seleccionan los sentidos figurados de los verbos.

<sup>11</sup> Koike, 2001

<sup>12</sup> Koike incluye aquí blandir, pero si es blandir espada, arma, etc. es léxico.

<b>Concretos</b>	instrumentos musicales	<i>guitarra, piano, cuerdas</i>	<b>Colocaciones léxicas</b>
	utensilios	<i>plancha, mordaza, tornillo</i>	
	medios de transporte	<i>tren, barco, guagua</i>	
	partes del cuerpo humano	<i>corazón, pelo, barba</i>	
	prendas de vestir	<i>sombrero, gafas, descosido</i>	
	alimentos	<i>café, pan, plato</i>	
	animales	<i>balidos, paloma, ganado</i>	
	otros		
<b>Abstractos</b>	de acción	<i>discusión, vigilancia, pago</i>	<b>Colocaciones funcionales.</b> Selección de acepción especial en el colocativo.
	de situación o estado de cosas	<i>crisis, norma, reputación</i>	
	de cualidad	<i>paciencia, habilidad, sensibilidad</i>	
	de evento	<i>boda, aniversario, congreso</i>	
	de sentimiento o sensación	<i>cariño, abatimiento, desesperación</i>	
	otros	<i>rumor, proyecto, indicio</i>	

Tabla 1-17. Clasificación de los sustantivos.

### 1.7.2. Verbo deslexicalizado + sustantivo<sub>CD</sub>

Compuesta por verbos funcionales generales, donde el sustantivo generalmente es deverbal:

“*dar + sustantivo de golpe*”:

*dar golpe*  
*dar bofetón*  
*dar empujón*  
*dar puñetazo*  
*dar tropezón...*

“*dar + sustantivo terminado en -azo excepto los aumentativos y algunas palabras primitivas como cazo o cedazo*”:

*dar bocinazo*  
*dar cambiazo*  
*dar frenazo...*

“*dar + sustantivo de voz o sonido*”:

*dar berrido...*



*dar:*

*dar albergue*  
*dar alcance*  
*dar aprobación*  
*dar asentimiento*  
*dar asilo*  
*dar cobijo...*

*dar un:*

*dar un abrazo*  
*dar un aviso*  
*dar un barrido...*  
*darse un(a):*  
*darse un atracón (de)*  
*darse un batacazo*  
*darse una ducha...*

*darse:*

*darse autobombo*  
*darse maquillaje*  
*darse postín...*

*“tener + sustantivo referido a cualidad; generalmente terminado en el sufijo –“dad”:*

*tener piedad*  
*tener curiosidad*  
*tener habilidad*  
*tener conocimiento*  
*tener costumbre...*

*“tener + sustantivo referido a sentimiento”:*

*tener amor*  
*tener cariño*  
*tener angustia*  
*tener pasión...*

*“tener + sustantivo referido a estado físico”:*

*tener jaqueca*  
*tener sueño*  
*tener tos*  
*tener fatiga*  
*tener apetito...*

*tener:*

*tener un altercado (con)*

*tener la ambición (de)*

*tener anhelo (de)*

*tener un apuro*

*tener atraso*

*tener confianza...*

*tener + su:*

*tener su inicio*

*tener su reflejo*

*tener su origen*

*tener su causa*

*tener su apogeo*

*tener su auge*

*tener su esplendor...*

“*hacer + sustantivo abstracto terminado en -ción, -sión; que son derivados de* *debebarles*”:

*hacer aclaración*

*hacer alusión*

*hacer reclamación*

*hacer adquisición...*

*hacer:*

*hacer acopio (de)*

*hacer alarde*

*hacer una afrenta*

*hacer una advertencia*

*hacer aspavientos*

*hacer boicot*

*hacer fortuna*

*hacer un gesto (de)*

*hacer pis*

*hacer un viaje...*

“*sentir + sustantivo de estado anímico o físico*”:

*sentir alegría*

*sentir calor*

*sentir fatiga*

*sentir inquietud*

*sentir júbilo...*

“*sentir + sustantivo de sentimiento*”:

*sentir cariño*

*sentir celos*

*sentir odio*

*sentir pudor...*

“*echar + sustantivo que indican acto de decir*”:

*echar bronca*

*echar perorata*

*echar reprimenda...*

“*echar + sustantivo relacionado con la vista*”:

*echar vistazo*

*echar mirada*

*echar ojeada...*

“*recibir, sufrir, llenarse, padecer, tener, obtener, encontrar, experimentar, pasar, seguir, cosechar, ganarse y algunos otros*” + *sustantivo abstracto forman colocaciones con valor pasivo:*

*recibir un achuchón*

*sufrir un atentado*

*tener una conmoción*

*padecer un desengaño*

*seguir un tratamiento...*

“*causar, producir, provocar, etc. + sustantivo de estado emotivo*”; *forman colocaciones con valor causativo:*

*acarrear un disgusto*

*infundir aliento*

*imponer condiciones*

*proporcionar alegría...*

### **1.7.3. Verbo deslexicalizado + (sustantivo<sub>CD</sub>) + preposición + sustantivo.**

Verbo intransitivo + sintagma preposicional forman colocaciones con valor pasivo:

*caer en el olvido*

*incurrir en desprecio*

*estar fuera de uso...*

*llenar, sumir, llevar, reducir, etc.* forman colocaciones con valor causativo. *abismar* (a alguien) en el dolor, *dejar* (a alguien) en la ruina, *incurrir* en el enojo.

### 1.7.4. Selección de una acepción especial en el verbo.

Verbo funcional específico + sustantivo<sub>CD</sub>:

*abrigar esperanzas*  
*albergar propósitos*  
*fragar idea*  
*brindar ayuda*  
*tejer una trama*  
*tramar un engaño...*

En las definiciones se aprecian clases léxicas para las bases.

Verbo aspectual<sup>13</sup> + sustantivo:

Sustantivo<sub>SUJETO</sub> de estado de ánimo + verbo incoativo:

*asaltar duda*  
*invadir melancolía*  
*entrar fiebre...*

Verbo de inicio de algo + sustantivo<sub>CD</sub>:

*abrir una crisis*  
*cobrar cariño*  
*emprender una marcha...*

{*caer, empeñarse, entrar, estallar, montar, ponerse, prorrumpir, romper*} + *en* + sustantivo:

*caer en el pecado*  
*empeñarse en una discusión*  
*estallar en llanto*  
*romper en sollozos...*

Las colocaciones de aspecto terminativo-resultativo se presentan en los tres esquemas sintácticos.

*pasarse el enfado*

<sup>13</sup> Los verbos aspectuales son aquellos que indican alguno de los aspectos que se recogen en la siguiente tabla; generalmente, cuando se usan en sentido figurado.

Aspecto	Valor	Ejemplo
<i>incoativo</i>	<i>comienzo</i>	<i>estallar una guerra</i>
<i>terminativo-resultativo</i>	<i>final</i>	<i>zanjar una cuestión</i>
<i>durativo-reiterativo</i>	<i>duración</i>	<i>encontrarse en un dilema</i>
<i>intensificativo</i>	<i>intensificación</i>	<i>corroer la envidia</i>

*desvanecerse el dolor*  
*abandonar una actitud*  
*levantar la sesión*  
*zanjar un asunto*  
*dejarse de bromas*  
*quedar en la impunidad*  
*distraer (a alguien) de preocupación...*

La mayoría de las colocaciones aspectuales durativas-reiterativas siguen el esquema:

*verbo + preposición + sustantivo*

y el verbo suele ser *andar* o *estar*.

*andar con chismes*  
*andar de caza*  
*estar de broma*  
*estar en circulación*  
*seguir en coma...*

Las de aspecto intensificativo también siguen principalmente el esquema:

*verbo + preposición + sustantivo*

También las hay en forma negativa. El sustantivo suele ser de sentimiento o de estado de ánimo.

*anegarse en llantos*  
*arder en deseos*  
*consumirse de envidia*  
*rebosar de alegría*  
*angustiar temor*  
*clavar la mirada (en)*  
*no ahorrar esfuerzos*  
*no escatimar medios...*

### **1.7.5. Colocaciones léxicas.**

Estas colocaciones están formadas por un verbo léxico y un sustantivo concreto entre los que hay una relación típica, como por ejemplo:

*tocar la guitarra*  
*afinar el piano*  
*perder el tren*

*moler café*  
*balarse la oveja*  
*apuntillar un toro*  
*tejer una alfombra...*

El significado es totalmente composicional. Son muy pocos los casos en que aparecen sustantivos abstractos, aunque también existen:

*amasar fortuna*  
*zanjar discusión...*

son algunas de las presentadas en los ejemplos por Koike(2001); sin embargo, los esquemas

verbo léxico general + sustantivo concreto

y

*verbo + sustantivo abstracto*

se encuadran en las de selección de una acepción especial.

#### **1.7.6. Sustantivo<sub>sujeto</sub> + verbo.**

Los verbos suelen ser intransitivos o pronominales. Las construcciones pronominales impersonales se incluyen dentro de este tipo de colocaciones:

*declararse una epidemia*  
*desatarse una polémica...*

Entre las colocaciones aspectuales existen numerosos ejemplos de éstos. Los verbos que intervienen en este tipo de colocaciones:

Denotan una acción característica de la persona o cosa designada por el sustantivo:

*palpitar corazón*  
*ensortijarse el pelo*  
*zarpar el barco...*

Indican fenómenos meteorológicos:

*despuntar el día*  
*caer la lluvia*  
*soplar el viento*  
*encapotarse el cielo...*

Se refieren a sonidos emitidos por animales:

*aullar el lobo*  
*cacarear el gallo*

*relinchar el caballo...*

Todos los casos se corresponden con colocaciones léxicas.

### 1.7.7. Sustantivo<sub>1</sub> + de + Sustantivo<sub>2</sub>.

En estas colocaciones, el sustantivo<sub>2</sub> es la base y el sustantivo<sub>1</sub> es el colocativo que expresan alguno de los significados siguientes:

El sustantivo<sub>1</sub> indica grupo o colección del sustantivo<sub>2</sub> si no puede expresar otros significados:

*banco de peces*

*ovillo de hilo*

*nube de polvo*

*bandada de aves...*

El sustantivo<sub>1</sub> significa porción del sustantivo<sub>2</sub> en combinación con éste:

*rebanada de pan*

*jarra de cerveza*

*copo de nieve*

*tableta de chocolate...*

El sustantivo<sub>1</sub> indica impulso violento o repentino, manifestación brusca del sentimiento especificado mediante el sustantivo<sub>2</sub>:

*arrebato de ira*

*acceso de histeria*

*ataque de tos...*

El sustantivo<sub>1</sub> en los dos primeros casos es cuantificador, y puede ser: **acotador** (*pedazo, trozo, tajada, etc.*), de **medida** (*kilo, litro, galón, etc.*), o de **grupo** (*grupo, serie, manada, etc.*).

### 1.7.8. Sustantivo + Adjetivo.

Los adjetivos también pueden ser colocacionales o no colocacionales, según intervengan o no en la formación de colocaciones. Los que pertenecen a la primera categoría son solo los calificativos, pero de éstos se pueden excluir los situacionales Tabla 1-18 puesto que los únicos que presentan restricciones combinatorias con sustantivos son los relacionales y los cualitativos (Koike, 2001).

<i>Relacionales</i>		<i>Nacionalidad, religión, francés, republicano, instituciones, clase carnívoro, musical, etc. social, etc.</i>
<i>Cualitativos</i>	<b>Cualidades físicas</b>	<i>Tamaño, altura, alto, bajo, largo, corto, ancho, anchura, volumen, estrecho, gordo, delgado, rojo, color, temperatura, blanco, caliente, frío, pesado, ligero, rápido, veloz, joven, peso, velocidad, tiempo, viejo, efímero, insípido, amargo, gusto, tacto, vista, claro, oscuro, fragante, olfato, oído, sonido, etc. ruidosos, silencioso.</i>
	<b>Cualidades no físicas</b>	<i>espirituales, vicios y audaz, temeroso, orgulloso, virtudes, estado perezoso, contento, aburrido, anímico, afecto, etc. alegre, enemigo, indiferente</i>
<i>Situacionales</i>		<i>cuasi determinativos, siguiente, próximo, pasado, referido</i>
		<i>adverbiales modales posible, probable, presunto,</i>
		<i>marcadores de la único, propio, principal intensidad</i>
		<i>circunstanciales antiguo, próximo, actual</i>
		<i>aspectuales frecuente, constante</i>

Tabla 1-18. Clasificación de los adjetivos calificativos, (Koike, 2001).

El colocativo en este caso es el adjetivo y al igual que en la construcción

*sustantivo + verbo*

o bien adquiere un valor funcional (adjetivos funcionales) o su significado léxico (adjetivos léxicos). Los adjetivos funcionales intensifican cuantitativa o cualitativamente el sustantivo; cuentan con dos posibilidades: que el sentido literal del adjetivo sea intensificador (adjetivos funcionales), o bien que el adjetivo se convierta en intensificador cuando se usa en sentido figurado (adjetivos funcionalizados). Un mismo adjetivo puede actuar como léxico o como funcional según el sustantivo con el que se combina. En la Tabla 1-19 se recogen ejemplos de ambos tipos de adjetivos colocacionales.



<b>Funcional</b>	Recto	<i>Dimensión física</i>	<i>grande</i>	<i>pena grande</i>
			<i>profundo</i>	<i>profundo dolor</i>
		<i>Cantidad</i>	<i>mucho</i>	
			<i>bastante</i>	
			<i>poco</i>	
		<i>Intensidad</i>	<i>intenso</i>	<i>genio fuerte</i>
			<i>fuerte</i>	<i>carácter débil</i>
			<i>débil</i>	
	Figurado		<i>mortal</i>	<i>susto mortal</i>
			<i>terrible</i>	<i>prisa terrible</i>
			<i>atroz</i>	<i>actuación atroz</i>
<b>Léxico</b>	Recto		<i>afilado</i>	<i>cuchillo afilado</i>
			<i>elocuente,</i>	<i>orador elocuente</i>
			<i>mortal</i>	<i>veneno mortal</i>

Tabla 1-19. Clasificación de los adjetivos colocacionales, (Koike, 2001).

Las combinaciones

*sustantivo abstracto + adjetivo funcional*

como

*pena grande*

son colocaciones, el adjetivo se usa con valor intensificador; sin embargo,

*sustantivo concreto + adjetivo funcional*

como

*casa grande*

no pasa de ser una combinación libre. Las colocaciones funcionales en las que interviene un adjetivo funcional son amplias, mientras que en las que se usa un adjetivo léxico en sentido figurado son estrechas. En los ejemplos siguientes se presentan las definiciones de algunos adjetivos funcionalizados; como se aprecia, alguno de los significados son: *muy grande*, *extraordinario*, *muy bueno*, etc. –en general se está ante una intensificación positiva superlativa.

**horroroso –sa**

**DRAE:** Que causa horror. **2.** fam. Muy feo.

**DEA 2.** (*col*) Muy grande o extraordinario.

**atroz**

**DRAE: 2.** Enorme, grave.

**3.** fam. Muy grande o desmesurado.

**DEA: 2.** Muy grande o extraordinario.

**fenomenal**

**DRAE: 3.** Fam. Tremendo, muy grande.

**4.** fig. Estupendo, admirable.

**DEA: 1.** Muy bueno o extraordinario.

**2.** Muy grande o tremendo.

**3.** Muy bien.

Al contrario de lo que sucedía en las colocaciones

*sustantivo + verbo*

en este caso sí que se pueden tener sustantivos referidos a persona como en:

*tonto redomado*

*lector empedernido*

*soldado bisoño...*

En la definición de los adjetivos suele hacerse referencia a las bases con las que puede combinarse; en general, tales bases que no son únicas, sino que el vínculo se puede establecer con un conjunto de sustantivos pertenecientes al mismo campo semántico: *fino* se puede combinar con sustantivos que denotan sentidos, como *oído, vista u olfato*. (Corpas, 1996).

Por último se señala que existe un tipo de construcción:

*sustantivo + sustantivo*

que autores como Corpas o Castillo incluyen dentro de esta categoría, porque uno de los sustantivos actúa en realidad con función de adjetivo dentro de la misma, como en:

*hombre clave*

*viaje relámpago...*

sin embargo, para Koike es un tipo de compuesto.

### 1.7.9. Verbo + Adverbio.

El adverbio es el colocativo en estas combinaciones, en la mayoría de los casos es de modo –terminado en *mente*–; funciona como un intensificador de la acción denotada por el verbo como se ve en los siguientes ejemplos:

*llover torrencialmente*  
*desear ardientemente...*

Algunas excepciones a esta característica son:

*hablar claro*  
*pisar firme...*

La mayoría de las colocaciones de este grupo presentan una variante en la forma:

*sustantivo + adjetivo*

de manera que el sustantivo de éstas es deverbal y se corresponde con el verbo de aquéllas y el adverbio–terminado en *mente*– derivado del adjetivo como se aprecia en los ejemplos de la Tabla 1-20.

<b>sustantivo + adjetivo</b>	<b>verbo +adverbio</b>
<i>lluvia torrencial</i>	<i>llover torrencialmente</i>
<i>comida opípara</i>	<i>comer opíparamente</i>
<i>cierre hermético</i>	<i>cerrar herméticamente</i>

Tabla 1-20. Algunas colocaciones *sustantivo + adjetivo* y la correspondiente *verbo + adverbio*

Sin embargo, las colocaciones *sustantivo + adjetivo* que se citan a continuación no tienen correspondencia con las *verbo + adverbio*:

*diferencia abismal*  
*cuestión capital*  
*sueldo ridículo*  
*error garrafal*  
*fe ciega...*

ni tampoco:

*diametralmente opuesto*  
*clínicamente muerto*  
*guardar celosamente...*

presentan la versión *sustantivo + adjetivo*.

Por último incluimos en este apartado el caso: *verbo + gerundio* con valor adverbial como:

*salir zumbando*

*salir pitando.*

### 1.7.10. Adverbio + Adjetivo.

En este caso también el adverbio es el colocativo y actúa como intensificador, pero de un adjetivo –en muchas ocasiones es un participio con valor adjetival–, como en:

*locamente enamorado*

*firmemente convencido*

*profundamente dormido*

*materialmente imposible...*

aunque hay casos en que la base es un verdadero adjetivo como en:

*locamente enamorado*

*mundialmente famoso...*

Nótese que en el caso

*adverbio + participio*

suele existir la correspondiente versión:

*verbo del que proviene el participio + adverbio*

o bien:

*sustantivo deverbal + adjetivo*

como se observa en la Tabla 1-21.

Se pone de manifiesto la flexibilidad formal en las colocaciones, en las que se puede dar el cambio de categoría gramatical de los colocados.

<b>sustantivo + adjetivo</b>	<b>verbo + adverbio</b>	<b>adverbio + adjetivo</b>
<i>relación estrecha</i>		<i>relacionado estrechamente</i>
	<i>dormir profundamente</i>	<i>profundamente dormido</i>
	<i>afectar visiblemente</i>	<i>visiblemente afectado</i>

Tabla 1-21. Correspondencia entre las colocaciones adverbio + adjetivo y otras estructuras sintácticas

### 1.7.11. Verbo + Adjetivo.

Son muy pocas las colocaciones en este grupo, incluso hay autores que no contemplan tal caso en sus clasificaciones; ejemplos de ellas son:

*{salir, resultar} ileso (de)*

*{salir, quedar} incólume (de)*

*{resultar, salir} redondo (el negocio)*

*andar {ajetreado, liado}*

*caer simpático, dejar frío...*

### 1.7.12. Verbo + como + sustantivo.

Este tipo de colocaciones las introdujo Penadés en la tipología colocacional; reproducen una estructura comparativa:

*dar vueltas como un trompo*

*dormir como un tronco...*

## 1.8. Las colocaciones en los diccionarios.

En español se dispone del DICE (Diccionario de Colocaciones del Español) de Alonso en el que se describen las restricciones combinatorias de un lexema mediante las funciones léxicas de Mel'čuk y los diccionarios Redes y Práctico de Bosque en el que se recogen las clases léxicas cuyos miembros coaparecen con un colocativo dado. Hay que mencionar también el Diccionario Fraseológico del Español Moderno (Varela, Kubarth, 1994), pero que solo recoge locuciones, complejos fraseológicos con casillas vacías y expresiones. Las dos primeras son locuciones, mientras que las expresiones coinciden con las fórmulas rutinarias de la clasificación de Corpas. En definitiva, tampoco recoge las colocaciones.

Además de los diccionarios específicos de colocaciones se puede explotar la información presente en los diccionarios tradicionales, especialmente en los de uso –Diccionario de Uso del Español (DUE, María Moliner) y el DEA (Diccionario del Español Actual, Seco)– sobre las restricciones combinatorias de las palabras.

### 1.8.1. Colocaciones en el DEA.

En las definiciones de este diccionario se especifica el **contorno**<sup>14</sup> de la misma, delimitado por corchetes y formado, en muchas ocasiones, por una lista de unidades léxicas que forman colocación con la unidad léxica que se define.

**panorámico -ca I adj 1** De(l) panorama. | FQuintana-Velarde *Política* 59: Tales son los rasgos básicos y generales .. Es preciso ahora obtener, además de la visión panorámica general que ya poseemos, algunos detalles concretos. **b)** [Imagen] que representa una vasta extensión de terreno. *Tb n f.* | \* Es una fotografía panorámica del pueblo. **c)** [Autocar o vagón] que permite ver

<sup>14</sup> La definición está constituida por dos clases de elementos, el contenido –abarca el significado en sí– y el contorno –recoge el contexto.

cómodamente el paisaje en todas direcciones. | J. Oyarzun *Abc* 15.10.70, 11: Estas viejecitas de trajes floreados y sombreros de paja que a veces vemos circular por Madrid en grandes autobuses panorámicos .. esconden bajo su aparente fragilidad una gran entereza. **d)** (Cine) [Pantalla] de superficie cóncava y más ancha que la normal, destinada a películas en cinemascope y similares. | \* Este cine tiene pantalla panorámica.

**paranormal** *adj* (*Psicol*) [Cosa, esp. fenómeno] que no tiene explicación científica. | *Ya* 13.2.75, 43: La autora estudia el tema agrupándolo así: creencias fantásticas y supersticiones; los muertos; .. facultades y fenómenos paranormales.

**sentir**<sup>1</sup> (*conjug* 60) A tr **1** Experimentar [una sensación o sensaciones]. *Tb abs.* | *Arce Testamento* 20: Comenzaba a sentir sed. *Gambra Filosofía* 135: Hay algunos que defienden el conocimiento en los vegetales, aduciendo a su favor ciertos hechos que parecen probar que los vegetales sienten. **b)** Experimentar [un estado afectivo o de ánimo]. | *CNavarro Perros* 92: Pensó en lo que sería la vida íntima de un ser como Fidel y sintió una lástima inmensa por su mujer. **c)** Ser afectado [por algo (*cd*)]. | *Ama casa* 1972 91: Los rododendros. Es un árbol hermoso .. Se cultiva en tierra de brezo y siente mucho el frío.

A la vista de los ejemplos se observa que en algunos casos se dan rasgos generales para especificar el contorno como “*una sensación o sensaciones*” en la acepción 1 de sentir, “*estado afectivo o de ánimo*”, en el b). En el caso de paranormal se tiene un elemento general como “*cosa*” y uno específico como “*fenómeno*” y en panorámico solo se dan elementos específicos en el contorno como “*imagen*” en b), “*autocar o vagón*” en c) y “*pantalla*” en d). Obsérvese que en todos estos casos se puede decir que se está ante colocaciones formadas por el lexema de la entrada y los elementos del contorno; aunque esto no sucede siempre, sí son numerosas las oportunidades en las que se da. Nótese que en algunos casos se obtiene una restricción combinatoria entre la entrada y una clase léxica:

“*sentir + sustantivo de estado afectivo o de ánimo*”<sup>15</sup>

mientras que en los restantes aparecen colocaciones específicas:

*imagen panorámica*  
*autocar panorámico*  
*vagón panorámico*  
*pantalla panorámica*  
*fenómeno paranormal*

<sup>15</sup> Koike también resalta esta característica para indicar las posibilidades de combinación del verbo sentir: “el verbo sentir puede combinarse con sustantivos que expresan ‘estado anímico o físico’ y ‘sentimiento’”.

Si bien “*cosa paranormal*” no es una colocación, en el ejemplo aparece

*facultad paranormal*

que sí lo es.

En general, cuando el lexema que se está definiendo es el colocativo, las bases de las colocaciones se recogen en el contorno, aunque en el caso de que el colocativo sea un adjetivo también se pueden encontrar en las explicaciones en cursiva que utilizan expresiones del tipo:

“*Dicho normalmente de, Referido esp. a, etc.*”.

Por otra parte, si en el contorno se habla de rasgos generales que definen una clase léxica, no quiere decir que todos los elementos que presenten ese rasgo formen colocación, puede haber excepciones. Además, como se aprecia en el apartado b) de panorámico solo hay un elemento en el contorno: “imagen”; sin embargo,

*fotografía panorámica*

también es colocación, si bien fotografía es un vocablo que se refiere a un tipo de imagen. En definitiva, a pesar de que en algunos casos si se consideran estrictamente los contornos, se puede incurrir en un exceso o defecto de producción de las bases que combinan con un colocativo, no supone un obstáculo para afirmar que la información contenida en el contorno es muy valiosa a la hora de caracterizar las restricciones combinatorias de los colocativos.

En algunos casos, también se dispone de los verbos con los que se combina un determinado sustantivo, como ejemplifica la Tabla 1-22; en muchos casos, el marcador *frec con el v* en la definición de algunos sustantivos en el DEA se corresponde con colocaciones.

<b>Sustantivo</b>	<b>Frec con el v</b>
<i>gustirrinín</i>	<i>dar</i>
<i>leñazo</i>	<i>dar</i>
<i>golpe</i>	<i>dar</i>
<b><i>golletazo</i></b>	<b><i>dar</i></b>

Tabla 1-22. Ejemplos del DEA.

### 1.8.2. Colocaciones en el DUE.

En las bases es donde se presentan las posibilidades combinatorias de los lexemas –en forma de enlaces frecuentes. Para los sustantivos se dan: los adjetivos y verbos; para los verbos, los adverbios y las preposiciones con que se construyen.

En el caso de los adjetivos, las colocaciones o las solidaridades léxicas pueden aparecer señaladas por “*se aplica a*”, “*aplicado a*”, “*dicese de*” o “*se dice (de)*”. También se introducen en los ejemplos que aparecen en letras simples y entrecomilladas, en el apartado de palabras afines o bien en el de frases y modismos, en cursiva y al final del artículo (Corpas, 1992).

### 1.8.3. Diccionario VOX

Este diccionario se puede considerar como paradigma de la precisión semántica que aportan las colocaciones a la lengua, abundando su uso en las definiciones de los lemas que recoge. Algunos ejemplos en los que se manifiesta este fenómeno

**animar:** Infundir ánimo y energía

**ilusionar:** Despertar esperanzas, gralte. atractivas

**truncar:** Quitar [a uno] las ilusiones o esperanzas

**confiar:** Dar confianza o esperanza [a uno]

**desahuciar:** Quitar [a uno] las esperanzas de conseguir lo que desea

**hechizar:** Cautivar el ánimo, embelesar

Del mismo modo que en el DEA, se observan colocaciones en los sustantivos marcados como contornos de los verbos:

**arriar:** Bajar [una vela o bandera que estaba izada]

**enarbolar:** Levantar en alto [estandarte, bandera, etc.]

**largar:** Desplegar [la bandera, las velas, etc.]

**vibrar:** Agitar en el aire [la pica, la lanza, etc.]; arrojar con ímpetu y violencia [una cosa que vibre]

**suspender:** especialmente, figurado. Privar temporalmente [a uno del sueldo o empleo] que tenía

**levantar:** Vigorizar [el ánimo, la moral, etc.]; engrandecer, ensalzar; impulsar hacia cosas altas [el pensamiento, el corazón, etc.]

Aparecen colocaciones funcionales cuando se definen verbos relacionados con un sustantivo:



**desesperanzar:** Quitar la esperanza [a uno]

**esperanzar:** Dar esperanza [a uno]

**esperanzar:** Tener esperanza

**apasionar:** Causar, excitar alguna pasión [a uno]

**apasionar:** Llenarse de pasión

**respetar:** Tener respeto

En cuanto a los adjetivos, las colocaciones se observan en el contorno:

**acéfalo:** [sociedad, secta, etc.] Que no tiene jefe.

**acéfalo:** [feto] Sin cabeza o sin parte considerable de ella.

**abisal:** [zona del mar profundo] [zona del mar profundo] Que se extiende más allá del talud continental, y corresponde a profundidades mayores de 2000 m

**frío:** por extensión. [color] Que produce efectos sedantes como el azul, el verde, etc.

**imberbe:** [joven] Que no tiene barba.

**interactivo:** [programa] Que permite una interacción, a modo de diálogo, entre el ordenador y el usuario.

**recurrente:** [fenómeno] Que vuelve a su punto de partida.

**visceral:** [impresión, sentimiento, etc.] Intenso, profundo y arraigado: odio visceral.

### 1.9. Elementos ajenos a las colocaciones.

Por la utilidad que puede suponer a la hora de implementar un reconocedor automático de colocaciones, se incluye una recopilación de elementos que se deben descartar por no formar parte de las mismas.

El verbo *ser* no puede formar colocaciones sustantivo + verbo.

Los adjetivos demostrativos, indefinidos y relativos.

Verbos auxiliares.

Los verbos derivados de sustantivos que son bases de colocaciones funcionales: *dar un consejo / aconsejar*.

Los verbos *querer, creer, pensar, indicar, afirmar, preguntar, verbos de decir, asegurar, advertir, vender, comprar, agradecer, obligar, ayudar, invitar, gustar, parecer, interesar, considerar, creer, hallar, poder, saber, querer, intentar, etc.*

Los sustantivos propios en las sustantivo + verbo.

Los numerales o indefinidos.

Los adjetivos situacionales.

Los adverbios funcionales que solo intensifican o cuantifican la acción denotada por el verbo: *absolutamente, enteramente, completamente, intensamente, tremendamente, terriblemente, etc.*

Adverbios de tiempo, de cantidad, de afirmación, de negación, de duda o de lugar.

## 2. Extracción automática de colocaciones utilizando criterios estadísticos.

### 2.1. Medidas de asociación

Las técnicas de extracción automática de colocaciones se basan en el criterio de la frecuencia de aparición conjunta de palabras en corpus textuales como manifestación cuantificable de la característica del uso preferente frente a las combinaciones libres. Los primeros estudios se deben a Firth (1957), quien utiliza por primera vez el término **collocation** para referirse a la coaparición de palabras. A partir de los trabajos de Firth surge la llamada escuela sistémica británica a la que pertenecen Halliday, Sinclair, Jones; para tal corriente, las colocaciones se reducen a combinaciones frecuentes de palabras; por utilizar el criterio de la preferencia como rasgo principal de las colocaciones, en sus trabajos se propone el uso de técnicas estadísticas para la exploración de textos en busca de combinaciones de palabras cuya frecuencia de aparición pueda considerarse superior a la de otras combinaciones.

El análisis de las combinaciones en esta línea forma parte de la disciplina conocida como Lingüística de Corpus, en la que se consideran esenciales los indicadores de los vínculos motivados entre los dos elementos de la colocación. Por lo general, determinan una puntuación para cada par de palabras que puede usarse para establecer un ranking, o bien seleccionar aquellos casos que se confirman como colocaciones por medio de valores de corte. Se basan en la frecuencia de aparición conjunta de palabras y en general intentan captar anomalías respecto a lo que se espera que suceda si la aparición de la combinación en el corpus se debiera a la mera casualidad: las variables aleatorias que representarían la aparición en el corpus de los elementos que constituyeran las combinaciones libres serían independientes. La independencia estadística ocurre cuando no existe relación alguna entre ambos fenómenos, ninguno de ellos incide en la ocurrencia del otro. En otras palabras, las probabilidades de ocurrencia conjunta,  $P(x, y)$  e individual de base y colocativos,  $P(x)$ ,  $P(y)$  respectivamente, verifican la relación:

$$P(x, y) = P(x) * P(y)$$

Desde los primeros trabajos en la materia se han ido proponiendo distintas modalidades de indicadores que van desde el más sencillo y directo, la **frecuencia relativa** a los test estadísticos como **z-score**, **t-score**, **fórmula de Dunning**, **test de Poisson**, o la **información mutua** adaptada de la teoría de la información. Todos ellos se calculan a partir de valores de frecuencias de aparición conjunta de dos palabras en el corpus:  $f(x, y)$  y de la frecuencia de aparición de forma individual de cada una de ellas:  $f(x)$ ,  $f(y)$ . En base a estos valores se puede

obtener como estimación de la probabilidad la frecuencia relativa:  $p(x) \cong \frac{f(x)}{N}$ , siendo  $N$  la cantidad de veces que se repite el experimento, en este caso, la cantidad total de palabras en el corpus. En virtud de la ley de los grandes números, la frecuencia relativa de un suceso se aproxima a su verdadera probabilidad cuanto mayor es el número de repeticiones de un suceso. Es decir, cuanto mayor sea el número de repeticiones del experimento, mayor será la fiabilidad de la aproximación.

### 2.1.1. Frecuencia Relativa

También denominada frecuencia de aparición de  $x$  con  $y$  (Koike, 2001); en realidad, es un indicador del porcentaje de veces que una palabra aparece con otra respecto al número total de veces que aparece.

$$frec(x|y) = \frac{frec(x,y)}{frec(x)} \times 100$$

Si este valor es elevado, quiere decir que si aparece  $x$  es muy posible que también aparezca  $y$ ; no tiene por qué presentar el mismo valor que la frecuencia de aparición de  $y$  con  $x$ . Para Koike es suficiente indicio de colocación que el valor sea superior al 20%.

### 2.1.2. Información Mutua

La medida de **información mutua** puntúa la coocurrencia de dos palabras mediante la expresión (Church, Hanks, 1990):

$$I(x,y) = \log_2 \frac{p(x,y)}{p(x) * p(y)}$$

Tal valor mide la cantidad de información (en bits) que la coocurrencia aporta sobre la información de las apariciones individuales de las dos palabras. De forma intuitiva, la información mutua mide en qué cantidad se reduce la incertidumbre de una variable  $X$  debido al conocimiento de la variable  $Y$ . En este caso se interpreta como la información que la aparición de una palabra aporta sobre la otra, cuando la información mutua es nula, estamos ante sucesos independientes. Las probabilidades se calculan directamente usando frecuencias relativas. Esta medida es inestable para cantidades pequeñas –se calcula solo si la frecuencia absoluta  $f(x,y) > 5$ .

La información mutua constituye un indicador más complejo y preciso que la frecuencia relativa utilizado para la extracción de colocaciones en Lingüística de Corpus.

### 2.1.3. z-score

El **z-score** mide si es significativa la frecuencia de aparición de  $x$  con  $y$  en el corpus, respecto a la frecuencia esperada, bajo la hipótesis de que las variables son independientes. Se basa en la aproximación normal de la variable aleatoria binomial que considera como éxito el que la palabra  $x$  ocurra con  $y$ :

$$z = \frac{f(x, y) - \bar{f}(x, y)}{\sqrt{f(x, y) * (1 - p(y))}}$$

$$\bar{f}(x, y) = p(y) * f(x) * |D|$$

$$p(y) = \frac{f(y)}{N - f(x)}$$

Siendo  $p(y)$  la probabilidad de que ocurra  $y$ , en una posición distinta en el corpus,  $f(x)$  la frecuencia de  $x$ ,  $N$  es el número de palabras en el corpus,  $|D|$  el número de posibilidades en las que puede aparecer  $y$  alrededor de  $x$  –coincide con dos veces la distancia colocacional que se esté usando en las líneas de concordancias (Pearce, 2002).

### 2.1.4. t-score

Este valor se utiliza ampliamente para la extracción de colocaciones y se inspira en la aplicación del *t-test de Student* para la media con varianzas desconocidas en distribuciones normales. Sin embargo, la distribución de las frecuencias de coocurrencia de palabras en el corpus no cumple con las condiciones de la teoría estadística en que se fundamenta el test original, por lo que su buen comportamiento le da rango de valor heurístico del z-score y no de un test de hipótesis (Evert, 2005).

$$t = \frac{f(x, y) - \frac{f(x) * f(y)}{N}}{\sqrt{f(x, y)}}$$

### 2.1.5. Test de Dunning

Se determina el ranking de colocaciones utilizando la razón de verosimilitud como test estadístico fiable con independencia del tamaño del corpus, ya que no se requiere la exigencia de normalidad en la distribución de la variable (Dunning, 1993). La razón de verosimilitudes se

calcula para la distribución binomial en la que bajo la hipótesis nula, la aparición de las dos palabras es independiente:

$$\log \lambda = \log L(f(x, y), f(x), p) + \log L(f(y) - f(x, y), N - f(x), p) \\ - \log L(f(x, y), f(x), p_1) - \log L(f(y) - f(x, y), N - f(x), p_2)$$

$$p = \frac{f(y)}{N}, p_1 = \frac{f(x, y)}{f(x)}, p_2 = \frac{f(y) - f(x, y)}{N - f(x)}$$

$$\log L(k, n, p) = k \log p + (n - k) \log(1 - p)$$

$$p = \frac{f(y)}{N}, p_1 = \frac{f(x, y)}{f(x)}, p_2 = \frac{f(y) - f(x, y)}{N - f(x)} \quad (\text{Manning, Schütze, 1 999})$$

### 2.1.6. Test de Poisson

Este test se construye bajo la suposición de que las coocurrencias siguen una distribución de Poisson, que corresponde a la extensión para grandes muestras de la distribución binomial. El valor se determina (Quasthoff, Wolff, 2002):

$$sig(x, y) = \frac{-\log p_k}{\log n} \approx \frac{\lambda - k \log \lambda + \log k!}{\log n}$$

$$\lambda = np_x p_y$$

$$p_k = (\text{más de } k \text{ ocurrencias de } x, y) = \sum_{l=k}^{\infty} \frac{1}{l!} \lambda^l e^{-\lambda}$$

## 2.2. Base de datos de información combinatoria.

Los valores de las distintas medidas de asociación de una combinación de palabras quedan determinados por las frecuencias de aparición conjunta de los elementos que la componen así como de las apariciones individuales de cada uno de ellos. Estos datos se deben obtener a partir de un corpus textual que será considerado como una muestra del lenguaje, por lo que es importante garantizar que sea representativo. Asimismo, su volumen debe ser tal que permita su tratamiento mediante computadores en un tiempo razonable. Este trabajo se centra en las colocaciones léxicas del español, por lo que se recopiló una gran cantidad de textos, de una amplia gama de géneros, albergando una colección de unos 11000 textos aproximadamente. Entre otros géneros este corpus reúne obras de literatura, tanto clásica como contemporánea,

española y universal, poesía y prosa, teatro, narrativa, ensayos, discursos y artículos periodísticos.

### 2.2.1. Características del Corpus

El corpus utilizado consiste en una colección de textos planos, sin ningún tipo de etiquetas que proporcionen información lingüística. Además de obras periodísticas se pueden encontrar obras literarias y obras no literarias. De las primeras, entre las que se han catalogado cuentos, relatos, poesía, novelas y otras encontramos 7758. En el segundo grupo se recopilaron 4108 tratados sobre arte, biografías, ciencias, cine, cristianismo, derecho, economía, educación, esoterismo, etnología, filosofía, geografía, historia, lingüística, literatura, política, psicología, religión, teatro y otros. Además, un grupo aparte lo conforman los periódicos. La cantidad de obras recogidas en cada uno de estos grupos se muestra en la Tabla 2-1 y la Tabla 2-2.

<b>Tipo</b>	
arte	21
biografías	85
ciencias	35
cine	220
cristianismo	44
derecho	513
economía	35
educación	40
esoterismo	22
etnología	21
filosofía	170
geografía	85
historia	265
lingüística	25
literatura	1068
política	346
psicología	70
religión	103
teatro	63
otros	877
<b>Total</b>	<b>4108</b>

Tabla 2-1 Distribución de las obras no literarias

<b>Tipo</b>	
novela	2053
teatro	1265
relato	2248
cuentos	632
otros	596
<b>Total</b>	<b>7758</b>

Tabla 2-2 Distribución de las obras literarias

En definitiva, una amplia muestra del español con un número total de palabras que está en torno a los 300 000 000, cuya distribución se resalta en el Gráfico 2-1 y Gráfico 2-2

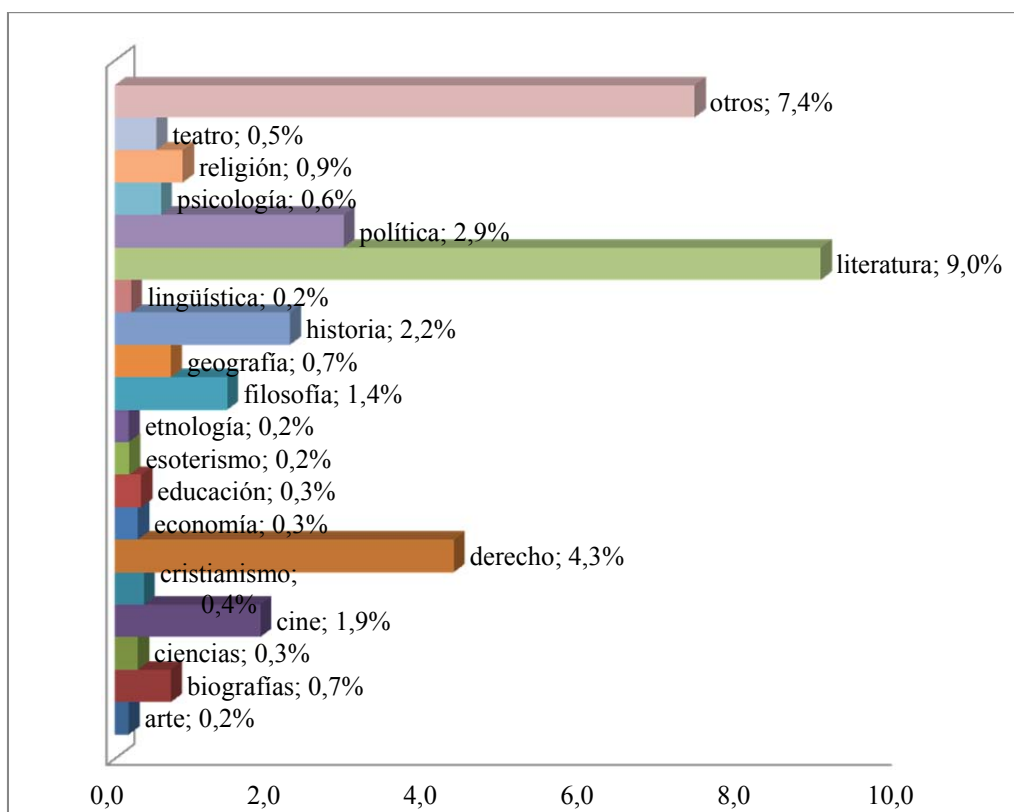


Gráfico 2-1 Distribución de las obras no literarias del corpus



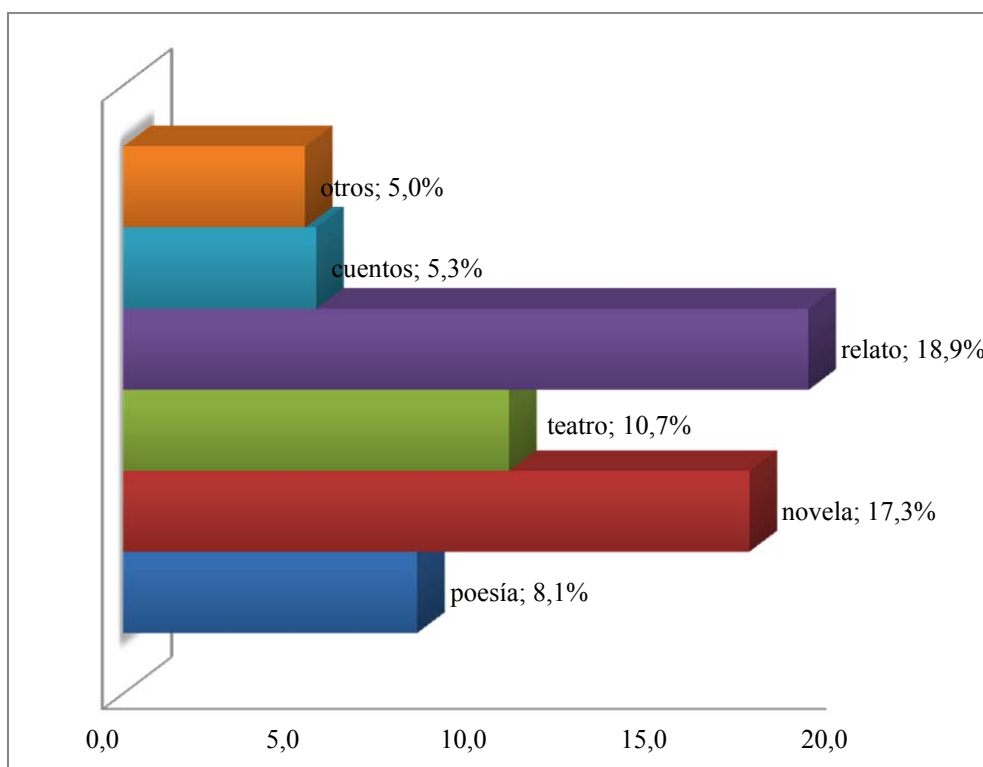


Gráfico 2-2 Distribución de las obras literarias del corpus

### 2.2.1.1. Tratamiento del corpus

El primer proceso llevado a cabo en esta investigación lo constituye la obtención de todas las combinaciones que se toman en consideración así como su frecuencia de aparición en el corpus. Los cálculos se realizan a partir de las llamadas *líneas de concordancia* extraídas de éste. Una línea de concordancia es un determinado fragmento de una longitud fija, en términos de cantidad de palabras que abarca, en el que aparece una palabra dada en un texto o corpus textual. La aparición conjunta de una determinada combinación en el ámbito de la línea de concordancia se considera una manifestación del uso conjunto de la misma o muestra de la combinación. En la Figura 2-1 se exponen líneas de concordancias obtenidas en el *corpus de referencia del español actual* (CREA) cuando se consulta sobre la palabra *arrojar*.

puesto muy cuesta arriba para su equipo, pero sin **arrojar** la toalla, porque "restan 90 minutos y no se

cuatro años. Pero al nuevo presidente le será difícil **arrojar** de sí la mala conciencia de tenerse de presta

el exilio no significa que el volcán deje de **arrojar** incertidumbre. Florentino Pérez vuelve a la carga, reclamando la

profesionales y características psicológicas de cada persona, así como **arrojar** un listado de arrojar un listado de los quince

Figura 2-1 Líneas de concordancia de la palabra arrojar en el CREA

Se salva en parte la falta de cualquier tipo de información lingüística en los textos gracias a las pautas marcadas por algunas características fundamentales de las colocaciones léxicas en español que se expusieron en el capítulo 1 y que se enumeran a continuación:

- a) Debido a la flexibilidad formal de las colocaciones se recuentan combinaciones de formas canónicas en lugar de combinaciones de palabras gráficas. Según esto se consideran muestras de la combinación *progreso-civilización*: *el progreso de las civilizaciones*, *la civilización progresó* o *los progresos de una civilización*, o bien de *reprobar-condena*: *reprobaron la condena*, *reprobó la condena*, *las condenas fueron reprobadas* o de *guerra-civil*: *guerra civil*, *las guerras civiles*, *la guerra entre civiles*. Para ello se ha utilizado el “*Flexionador y Lematizador de palabras del español*” del Grupo de Estructura de Datos y Lingüística Computacional de la Universidad de Las Palmas de Gran Canaria (Santana et al, 2007,1999,1997).
- b) Se considera un catálogo de *palabras vacías* como artículos, preposiciones, conjunciones, interjecciones, determinantes numerales, determinantes ordinales, adjetivos demostrativos, indefinidos y el relativo *cuanto*, los verbos *ser*, *estar* y *haber*, etc. descartándose del proceso todas las que se incluyeron en él. (Koike, 2001).
- c) Se consideran exclusivamente las estructuras colocacionales (Koike, 2001) que se pueden encontrar en el español, desde la perspectiva del Lematizador se pueden agrupar en los casos:

*sustantivo + verbo*

*sustantivo + adjetivo*

*sustantivo + de + sustantivo*

*verbo + adverbio*

*verbo + adjetivo*

*adverbio + adjetivo*

por lo que en el recuento se contabilizan solo aquellas combinaciones que se corresponden con alguno de estos patrones en orden indistinto.

El programa simula la obtención de líneas de concordancias para cada una de las formas canónicas en el corpus, generando un diccionario de formas canónicas y otro de combinaciones de formas canónicas a partir de ellas. Cada línea de concordancia está dentro del ámbito de una frase y se toma un ancho de ventana 10. Es decir, se recopilan los valores de las frecuencias conjuntas para distancias entre palabras que van desde 1 hasta 10, tanto a la izquierda como a la derecha. Se obtienen las combinaciones de formas canónicas que se corresponden con alguna de las estructuras de las colocaciones en español y su frecuencia en el corpus. Si el tamaño de frase es menor que la amplitud fijada para la concordancia, no se complementa con palabras en la

---

siguiente frase. Cada elemento en la línea de concordancia actual que no sea una palabra vacía se procesa con el lematizador que proporciona la correspondiente información morfológica.

Para una mayor eficiencia en el procesamiento, cada palabra en el corpus solo se lee una vez; se almacena temporalmente la información morfológica que se genera gracias a una estructura de cinta circular, con tamaño igual a la distancia colocacional. En cada iteración se determina el elemento de dicha cola donde se deben registrar las formas canónicas con las que tiene correspondencia la palabra actual, en los restantes estarán almacenadas las formas canónicas que se obtuvieron cuando fueron procesadas las palabras previas en el corpus. Cuando se lee la palabra siguiente, pasará a ser palabra actual. Las palabras previas permanecen en la cinta siempre que se esté en el ámbito de la línea de concordancia. Para lograr esto, en cada iteración antes de insertar la información de la palabra actual extraída del corpus, se agregan al diccionario las combinaciones compuestas por las formas canónicas en el registro de la cinta que se va a actualizar y las formas canónicas en todos los restantes elementos en ella, siempre que correspondan a alguna de las estructuras colocacionales del español. En cada caso se actualizará la frecuencia si fue registrado previamente o creando la nueva entrada en caso contrario. Por último, es necesario especificar que si la palabra actual es vacía, no se lematiza y se deja la posición correspondiente en la cola vacía. De esta forma se logra obtener las parejas de formas canónicas con palabras previas y con palabras posteriores en la línea de concordancia. A continuación se ilustra el procesamiento realizado al corpus sobre un fragmento de *Odisea 2001*, de Arthur C. Clarke, novela que forma parte del corpus, en concreto, el primer párrafo:

La sequía había durado ya diez millones de años, y el reinado de los terribles saurios tiempo a que había terminado. Aquí en el ecuador, en el continente que había de ser conocido un día como África, la batalla por la existencia había alcanzado un nuevo clímax de ferocidad, no avistándose aún al victorioso. En este terreno baldío y disecado solo podía medrar, o aun esperar sobrevivir, lo pequeño, lo raudo o lo feroz. (*Odisea 2001*, Arthur C. Clarke).

La tabla se utiliza para representar la cola circular en la que se irá insertando temporalmente la información morfológica de cada palabra que se va leyendo en el fichero. Se muestra dicha estructura después de haber leído las 5 primeras palabras, las casillas sin información se corresponden a palabras vacías (Figura 2-2).

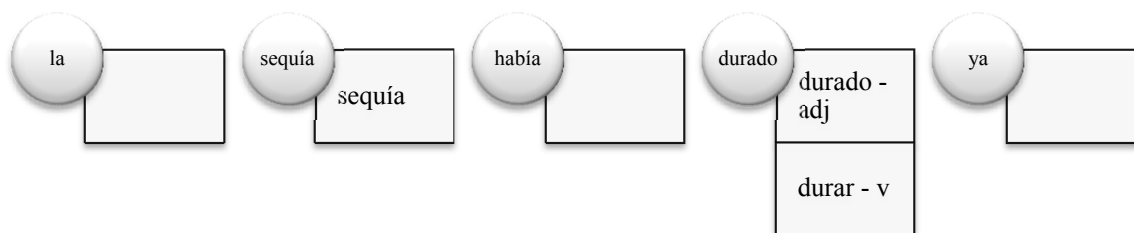


Figura 2-2 Inicialización de la cinta para la extracción de combinaciones

En la siguiente iteración la ventana se desplaza hacia la palabra *diez*, que será la siguiente que se lea en el corpus. La palabra actual es *diez* y se inserta en la cinta, no es necesario agregar información a los diccionarios de combinaciones, porque en la posición en la cola circular no hay elementos que correspondan previamente a una palabra vacía. Por el contrario, el diccionario de formas canónicas procesadas será actualizado con la información que produce la lematización de *diez* (Figura 2-3).

La sequía ~~había~~ durado ya diez millones ~~de~~ años, y ~~el~~ reinado ~~de~~ los terribles saurios tiempo a ~~que~~ había terminado. Aquí ~~en el~~ Ecuador, ~~en el~~ continente ~~que~~ había ~~de~~ ser conocido un día ~~como~~ África, ~~la~~ batalla ~~por~~ la existencia ~~había~~ alcanzado un nuevo clímax ~~de~~ ferocidad, ~~no~~ avistándose aún ~~al~~ victorioso. En este terreno baldío y disecado ~~se~~ podía medrar, ~~o~~ ~~aun~~ esperar sobrevivir, ~~lo~~ pequeño, ~~lo~~ raudo ~~o~~ ~~lo~~ feroz.

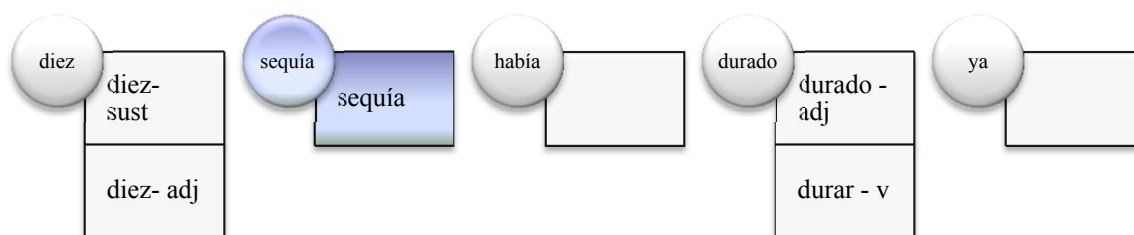


Figura 2-3 Actualización de la cinta para la extracción de combinaciones.

A continuación, se inserta la información morfológica relativa a *millones*; debe actualizarse las combinaciones de formas canónicas, es decir, la forma canónica *sequía* / *sustantivo* y sus combinaciones factibles con todas las formas canónicas en la cinta en ese instante: *sequía* + *durado*, *sequía* + *durar*, *sequía* + *diez*, *sequía* + *diez*

La sequía ~~había~~ durado ya diez millones ~~de~~ años, y ~~el~~ reinado ~~de~~ los terribles saurios tiempo a ~~que~~ había terminado. Aquí ~~en el~~ Ecuador, ~~en el~~ continente ~~que~~ había ~~de~~ ser conocido un día ~~como~~ África, ~~la~~ batalla ~~por~~ la existencia ~~había~~ alcanzado un nuevo clímax ~~de~~ ferocidad, ~~no~~ avistándose aún ~~al~~ victorioso. En este terreno baldío y disecado ~~se~~ podía medrar, ~~o~~ ~~aun~~ esperar sobrevivir, ~~lo~~ pequeño, ~~lo~~ raudo ~~o~~ ~~lo~~ feroz.

Una vez actualizado el diccionario de combinaciones, se elimina seqúa de la cinta y se agrega en su lugar millones, quedando con la configuración en la Figura 2-4.

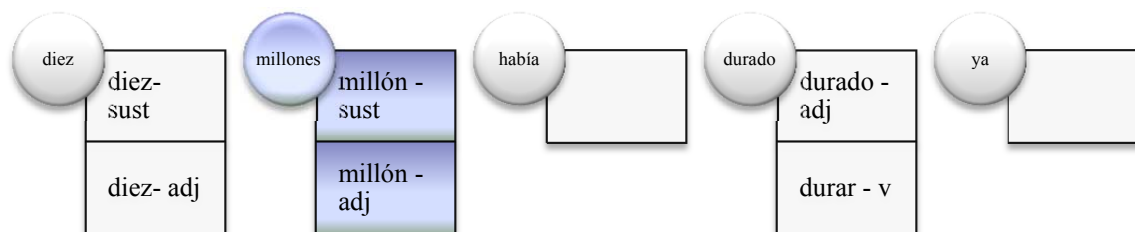



Figura 2-4 Actualización de una posición no vacía de la cinta de combinaciones.

Para mantener ambos diccionarios, el de formas canónicas y el de combinaciones de éstas se utiliza tablas hash cuyo tipo base permite registrar tanto el elemento como la frecuencia de aparición del mismo. De esta forma la ejecución permite generar la información de formas canónicas correspondientes a palabras no vacías y sus frecuencias en el corpus.

El lematizador del GEDLC permite identificar las formas canónicas cuya flexión o derivación produce la palabra por la que se consulta. Sin sumar la ampliación inherente a los prefijos y a los pronombres enclíticos, el universo está compuesto por alrededor de cuatro millones y medio de formas flexionadas y derivadas obtenidas a partir de las más de ciento veinticuatro mil formas canónicas con que opera el sistema. En los verbos, trata la conjugación simple y compuesta, los pronombres enclíticos, la flexión del participio como adjetivo verbal (género, número, grado superlativo y adverbialización) y el diminutivo del gerundio. Con las formas no verbales, considera: género y número en los sustantivos, adjetivos, pronombres y artículos; heteronimia por cambio de sexo en los sustantivos; grado superlativo en los adjetivos y adverbios; adverbialización y adverbialización del superlativo en los adjetivos; derivación apreciativa en los sustantivos, adjetivos y adverbios; variantes gráficas en todas las categorías gramaticales; formas invariantes tales como preposiciones, conjunciones, exclamaciones, palabras de otros idiomas y locuciones o frases. Se contempla la prefijación. (Santana et al, 2007). Por motivos de rendimiento, debido a la extensión del corpus, se ha utilizado la versión local, si bien está disponible como servicio web. En la Figura 2-6 se muestra la información que se obtiene al lematizar una palabra, se ha utilizado para generar el ejemplo la aplicación web disponible para consultas.

  
Grupo de Estructuras de Datos  
y Lingüística Computacional

## Flexionador y lematizador de palabras del español

( etiquetador morfológico - analizador morfológico )  
( lemmatizer - morphologic tagger - morphologic analyzer )

Introduzca una forma para lematizar:

[English](#)

- La aplicación informática que se presenta lematiza cualquier palabra del español al identificar su forma canónica, categoría gramatical y la flexión o derivación que la produce, y obtiene las formas correspondientes a partir de una forma canónica y de la flexión o derivación solicitada; tanto el reconocimiento como la generación operan sobre una misma estructura de datos, recorrerla en sentidos contrarios implica que la herramienta funcione en una u otra modalidad.
- En los verbos, trata la conjugación simple y compuesta, los pronombres enclíticos, la flexión del participio como adjetivo verbal (género, número) y el diminutivo del gerundio. Con las formas no verbales, considera: género y número en los sustantivos, adjetivos, pronombres y artículos; heteronimia por cambio de sexo en los sustantivos; grado superlativo en los adjetivos y adverbios; adverbialización del superlativo en los adjetivos; derivación apreciativa en los sustantivos, adjetivos y adverbios; variantes gráficas en todas las categorías gramaticales; formas invariantes tales como preposiciones, conjunciones, exclamaciones, palabras de otros idiomas y locuciones o frases. Tanto en la lematización como en la generación se contempla la prefijación.
- Permite el reconocimiento, la generación y la manipulación de las relaciones morfológicas a partir de cualquier palabra, incluye la recuperación de toda su información lexicogenética hasta llegar a una primitiva, la gestión y control de los afixos en el tratamiento de sus relaciones, así como la regularidad en la relación establecida. Proporciona una visión global del comportamiento y productividad de las palabras del español en los principales procesos de formación (sufijación, prefijación, parasíntesis, supresión, regresión, modificación-cero, apócope, metátesis y otros no clasificables que generan grafías alternativas).
- A partir de 151103 formas canónicas (incluye 14859 nombres de personas y apellidos), se obtienen algo más de 4900000 formas flexionadas y derivadas (sin contar la ampliación inherente a los prefijos y a los pronombres enclíticos) y se establecen unas 90000 relaciones morfológicas. El sistema incluye todas las entradas del Diccionario de la Lengua Española de la Real Academia, del Diccionario General de la Lengua Española, del Diccionario de Uso del Español de María Moliner, del Gran Diccionario de la Lengua Española de Larousse Planeta, del Diccionario de Uso del Español Actual Clave SM<sup>(1)</sup>, del Diccionario de voces de uso restringido por Manuel Alvar Ezquerro<sup>(2)</sup>, del Gran Diccionario de Sinónimos y Antónimos de Espasa-Calpe y del Diccionario Ideológico de la Lengua Española de Julio Casares.

Figura 2-5 Lematizador del GEDLC de la ULPGC

<http://www.gedlc.ulpgc.es/investigacion/scogeme02/lematiza.htm>

  
Grupo de Estructuras de Datos  
y Lingüística Computacional

## Resultados de la lematización

Resultado del reconocimiento de millones

- Forma canónica: millón
- Categoría: sustantivo masculino
- Flexión: masculino plural

- Forma canónica: millo
- Categoría: sustantivo masculino
- Flexión: augmentativo ( masculino plural )

[\[ Flexión verbal \]](#) | [\[ Flexión sustantivo \]](#) | [\[ Flexión adjetiva \]](#) | [\[ Flexión otras formas \]](#) | [\[ Lematización \]](#) | [\[ Relaciones morfológicas \]](#)  
[\[ Desambiguador local \]](#) | [\[ Desambiguador morfosintáctico \]](#)

Enviar sugerencias a [José R. Pérez Aguir](#)  
(C) Grupo de Estructuras de Datos-ULPGC

Figura 2-6 Respuesta a una consulta realizada al lematizador

Se hace necesario hacer notar que al utilizar exclusivamente información morfológica es inviable determinar qué forma canónica es la que produce flexión o palabra derivada que se haya leído del corpus. Por este motivo para cada palabra se contemplan todas las formas canónicas posibles. Desde el punto de vista de las colocaciones, en una gran parte de los casos se producirán combinaciones de formas canónicas con distintas categorías correspondientes a la misma colocación, como manifestación de la característica de la flexibilidad formal.

La ejecución del programa sobre el corpus de trabajo produce un total de 14475136 combinaciones distintas con frecuencia de aparición mayores o iguales que 3 en alguna de las

distancias. Las combinaciones con frecuencia 1 ó 2 se desprecian debido a la explosión de combinaciones irrelevantes que producen.

### 2.3. La Base de Datos

Con el objetivo de manipular la cantidad ingente de información generada a partir del tratamiento del corpus de la forma más eficiente posible se almacenaron los resultados en una BDD (Base de Datos) en la que inicialmente se recogen formas canónicas y su frecuencia de uso así como combinaciones de formas canónicas y su frecuencia de uso. Se muestra en el Gráfico 2-3 la distribución de las combinaciones registradas en la BDD según las categorías gramaticales de los dos elementos de la combinación, como grupos más numerosos se catalogan 6 551 979 en el tipo *sustantivo + verbo* y 3 743 476 en el tipo *sustantivo + adjetivo*.

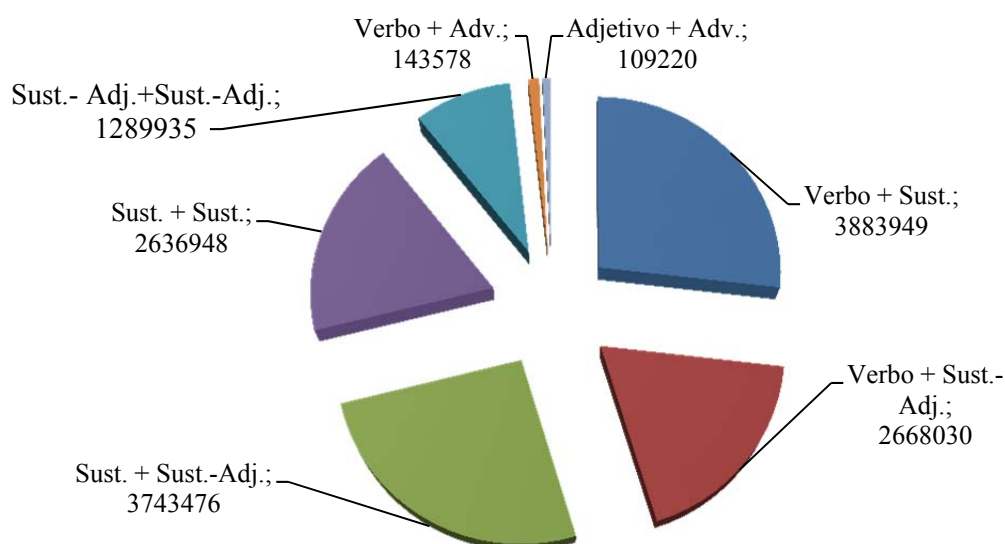


Gráfico 2-3 Clasificación de las combinaciones por categorías gramaticales de los elementos

A partir de la información inicial de frecuencias extraídas mediante el tratamiento del corpus, y utilizando el sistema gestor se amplía el conjunto de datos con los valores de los indicadores: **frecuencia relativa**, **z-score**, **t-score** y **test de Dunning** evaluados sobre todos los datos del corpus. Se incorpora también el catálogo de colocaciones recopiladas en los artículos que versan sobre colocaciones léxicas en español consultados para el desarrollo del trabajo (Alonso, 1994-1995), (Alonso, 2002), (Bargalló, 1997-1998), (Blasco, 2002), (Bosque, 2001), (Castillo, 1997-1998), (Castillo, 1998), (Castillo, 2001), (Corpas, 1996), (Corpas, 2001), (Corpas, 2003), (García 2002), (García-Page, 2001), (Koike, 2001), (Moreno, 1998), (Penadés, 2001), (Varela, 1994), (Zuluaga, 2002) y en el DCECR. Tanto unas como otras se usan para contrastar

la capacidad de las distintas técnicas basadas en la frecuencia a la hora de detectar colocaciones. El primer grupo se denominará a partir de ahora *combinaciones recopiladas*, donde se encontraron 2356 combinaciones, y en el DCECR se registran 55 052.

En la Figura 2-7 se muestra el grafo relacional que refleja cómo está organizada la información.

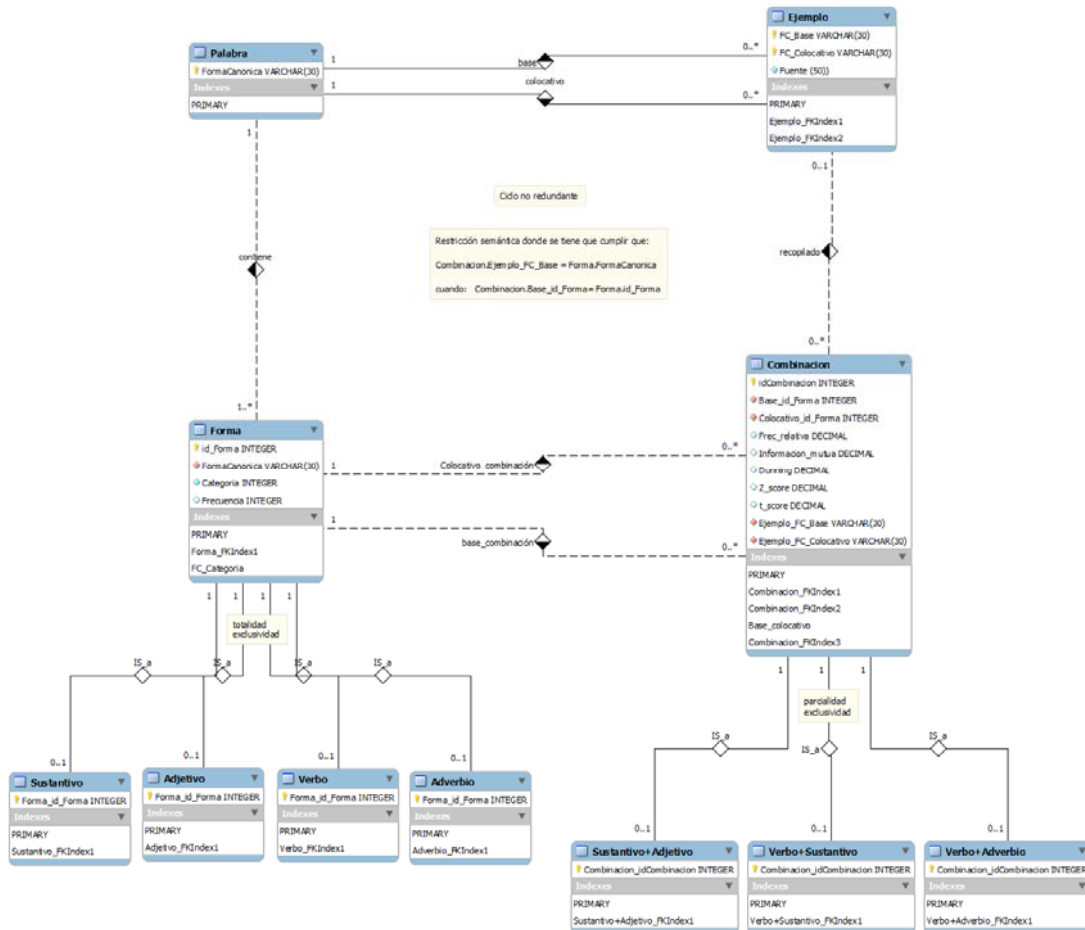


Figura 2-7 Grafo relacional de la Base de Datos de Combinaciones

A continuación se describen cada una de las entidades y sus atributos.

<b>Entidad</b>
Ejemplo
<b>Descripción</b>
Entidad que contiene la información de las combinaciones recopiladas y del DCECR.
<b>Atributos</b>
FC_Base: Forma Canónica que es base de la colocación.
FC_Colocativo Forma Canónica que es colocativo de la colocación.
Fuente: Puede ser Recopilada o DCECR.



<b>Entidad</b>
Palabra
<b>Descripción</b>
Entidad que contiene la información de la palabra que se toma para la forma canónica, a diferencia de ésta que abarca también la categoría, la frecuencia entre otros.
<b>Atributos</b>
FormaCanonica: Palabra que corresponde a una forma canónica.

<b>Entidad</b>
Forma
<b>Descripción</b>
Entidad que contiene la información de una Forma Canónica.
<b>Atributos</b>
id_Forma: identificador de la Forma Canónica
FormaCanonica: Forma Canónica
Categoría: Categoría de la Forma Canónica
Frecuencia: Frecuencia de aparición de la forma canónica en el corpus.

<b>Entidad</b>
pg
<b>Descripción</b>
Entidad que contiene para cada grupo del Diccionario Ideológico Vox su probabilidad a priori. Pueden ser grupos de <i>verbos, sustantivos, adjetivos o adverbios</i> .
<b>Atributos</b>
idGrupo: identificador del grupo
pg: probabilidad a priori del grupo idGrupo

<b>Entidad</b>
Combinación
<b>Descripción</b>
Entidad que contiene la información de una combinación que se ha extraído del corpus. Puede ser de alguno de los tipos: <i>verbo + sustantivo</i> , <i>verbo + adverbio</i> , <i>sustantivo + adjetivo</i> .
<b>Atributos</b>
id_Combinación: identificador de la combinación.
Base_id_Forma: Indentificador de la base
Colocativo_id_Forma: Identificador del colocativo
Frec_absoluta: Frecuencia absoluta de la combinación
Frec_relativa_base: Valor de la frecuencia relativa de la combinación respecto a la base
Frec_relativa_colocativo: Valor de la frecuencia relativa de la combinación respecto al colocativo
Informacion_mutua: Valor de la información mutua de la combinación
Dunning: Valor del test de Dunning de la combinación
z_score: Valor del z-score de la combinación
t_score: Valor del t-score de la combinación
Ejemplo_FC_Base:
Ejemplo_FC_Colocativo:

<b>Entidad</b>
IdeologicoVox
<b>Descripción</b>
Entidad que contiene la clasificación semántica de palabras del Diccionario Ideológico Vox (DIV)
<b>Atributos</b>
NumeroCabecera: Identificador de la Cabecera a la que pertenece el Grupo
NumeroGrupo: Identificador del Grupo al que pertenece la palabra
NumeroPalabra: Identificador de la palabra
Palabra: Palabra en la clasificación
Categoria: Categoría de la palabra

<b>Entidad</b>
Ideologico_VOX_Cabeceras
<b>Descripción</b>
Entidad que contiene la información de las cabeceras del DIV.
<b>Atributos</b>
Numero: Identificador de la cabecera
Cabecera: Palabra que representa a la cabecera.
Cuadro: Cuadro al que pertenece la cabecera, representa un nivel superior en la clasificación.
CardSust: Cardinal del conjunto de sustantivos en la cabecera.
CardVerbos: Cardinal del conjunto de verbos en la cabecera.
CardAdj: Cardinal del conjunto de adjetivos en la cabecera.
CardAdv: Cardinal del conjunto de adverbios en la cabecera

<b>Entidad</b>
Cuartiles
<b>Descripción</b>
Entidad que contiene los cuartiles y rango intercuartílico de las frecuencias relativas de cada forma del corpus para poder calcular estadísticos de valores atípicos. Pueden ser cuartiles de combinaciones <i>verbo + sustantivo</i> , <i>verbo + adverbio</i> o <i>sustantivo + adjetivo</i> .
<b>Atributos</b>
id_Forma: identificador de la Forma Canónica
Q1: Cuartil 1 de las frecuencias relativas de las combinaciones en las que aparece la forma canónica.
Q2: Cuartil 2 o mediana de las frecuencias relativas de las combinaciones en las que aparece la forma canónica.
Q3: Cuartil 3 de las frecuencias relativas de las combinaciones en las que aparece la forma canónica.
RI: Rango intercuartílico de las frecuencias relativas de las combinaciones en las que aparece la forma canónica.

## 2.4. Análisis de los resultados

Se utilizan los valores obtenidos de los indicadores para las combinaciones recopiladas y las del DCECR que aparecen en el corpus explorado para analizar la eficacia de las medidas de colocabilidad en la literatura de referencia, junto con su comparación con los datos extraídos de

la totalidad del corpus. Hay que destacar que a pesar de la extensión de éste, de las 2 356 combinaciones recopiladas se encuentran muestras de 1832 de ellas y en cuanto a las 55 052 que provienen del DCECR se obtiene información de 28403. Por tanto, si bien los datos registrados se consideran una amplia muestra de las capacidades de combinatoria de las palabras del español, el catálogo se considera totalmente abierto.

En los primeros epígrafes de este apartado se analizan los resultados obtenidos en general, para posteriormente realizar el estudio según la estructura de la combinación corresponda a *sustantivo + verbo*, *sustantivo + adjetivo* o *verbo + adverbio*.

#### 2.4.1. Frecuencia Relativa e Información Mutua

Inicialmente se valora la posibilidad de establecer como puntos de corte para delimitar la frontera entre las colocaciones y las combinaciones libres una puntuación de 0,2 en el caso de la frecuencia relativa (Koike, 2001) y de 3 en el caso de la información mutua (Kenneth, 1990). En la Tabla 2-3 se presenta un resumen de los valores de ambos indicadores en la muestra tratada. A partir de ella se aprecia que el valor de corte en el caso de la frecuencia relativa es excesivamente restrictivo. Por otra parte, los valores máximos no necesariamente lo alcanzan las colocaciones. En este sentido, se detecta que valores de frecuencia relativa 1 corresponden con combinaciones cuya frecuencia absoluta es baja (Tabla 2-6), siendo raras las combinaciones que aparecen en el corpus con una frecuencia absoluta superior a 5 dentro de este grupo.

Fuente	Casos con Frec. Rel. $\geq$ 0,2	Casos con Inf. Mutua $\geq$ 3	Rango de la Frec. Rel.	Rango de la Inf. Mutua
Recopiladas	7%	92,79%	[0,00018196 0,8666667]	[-0,72 15,45]
DCECR.	2%	83%	[9,3194E-05 0,8833333]	[-1,25 16,50]
Corpus	0%	25,74%	[1,536E-06 2,7]	[-3,66 26,57]

Tabla 2-3. Resumen de valores de los indicadores obtenidos en el corpus

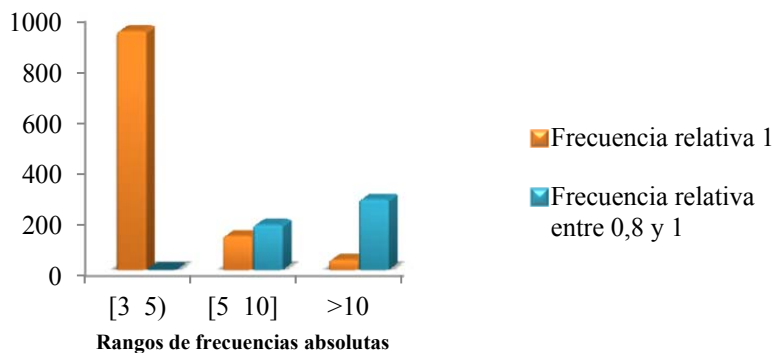


Gráfico 2-4 Distribución de las combinaciones según frecuencias totales

Los rangos encontrados para los valores de la información mutua indican que cuando se utiliza esta medida de asociación se hace necesario establecer además un límite superior. Una inspección exhaustiva de las combinaciones que alcanzan dichos valores máximos revela que en bastantes de ellas, ambas formas canónicas pertenecen al mismo campo semántico: *cambaleo-gangarilla*, *anabolismo-catabolismo*, *obrepción-subrepción* son ejemplos de ello, como se observa en las definiciones del DRAE. Este análisis manual pone de manifiesto que el límite inferior de la información mutua también debe ser modificado, ya que en caso contrario se seleccionarían como colocaciones una gran cantidad de combinaciones que no lo son.

**cambaleo**

**m.** Compañía antigua de la legua, compuesta ordinariamente de cinco hombres y una mujer que cantaba.

**gangarilla**

**f.** Compañía antigua de cómicos o representantes, compuesta de tres o cuatro hombres y un muchacho que hacía de dama.

**anabolismo**

**m. Biol.** Conjunto de procesos metabólicos de síntesis de moléculas complejas a partir de otras más sencillas.

**catabolismo**

**m. Biol.** Conjunto de procesos metabólicos de degradación de sustancias para obtener otras más simples.

**obrepción**

**f. Der.** Falsa narración de un hecho, que se hace al superior para sacar o conseguir de él un rescripto, empleo o dignidad, de modo que oculta el impedimento que haya para su logro.

**subrepción**

**f.** Acción oculta y a escondidas.

**2. f. Der.** Ocultación de un hecho para obtener lo que de otro modo no se conseguiría.

Cabe destacar que los valores de frecuencias relativas mayores que 1 se producen en los casos en el que en la misma línea de concordancia se repiten formas canónicas, se tienen solo 55 combinaciones de este tipo en el corpus, por lo que se opta por desecharlas del estudio. Uno de los ejemplos que mejor explica estos casos es la frase:

*“Matarile rile rile, matarile rile rón”.*

La lematización produce los resultados:

Forma	<i>matarile</i>	<i>rilar</i>	<i>rilar</i>	<i>matarile</i>	<i>rilar</i>	<i>rón</i>
<b>Canónica</b>						
<b>Categoría</b>	Sustantivo	Verbo	Verbo	Sustantivo	Verbo	No reconocida

Tabla 2-4 Frecuencias relativas mayores que 1

Los valores de frecuencia que se recogen para las combinaciones *verbo + sustantivo* registradas en el texto del ejemplo se muestran en la Tabla 2-5.

	Frecuencia	Frec. Relativa
<i>matarile + rilar</i>	6	
<i>A = Matarile</i>	2	$6/2 > 1$
<i>B = Rilar</i>	3	$6/3 > 1$

Tabla 2-5 Frecuencias calculadas en el ejemplo

Se muestra en la Tabla 2-6 el ranking de las 20 combinaciones con mejores puntuaciones de la frecuencia relativa, en la Tabla 2-7 el de las 20 mejores cuando se usa la información mutua. Estos resultados dejan patente que las combinaciones libres se entremezclan en zonas de las que deberían quedar excluidas para hacer fiable la extracción automática por medio de estos indicadores. A la luz de estos resultados se ve la necesidad de fijar la atención también en los valores de la frecuencia absoluta de la combinación. Esto se debe a que gran parte de las combinaciones con este dato bajo y frecuencia relativa o información mutua superando los umbrales establecidos no son colocaciones. A partir del ensayo con distintas posibilidades se fija la exigencia de al menos 10 muestras de la combinación para considerar fiable la puntuación alcanzada por una combinación. En el Gráfico 2-5 se contrasta el valor mínimo de muestras requerido para considerar fiables los indicadores sobre las colocaciones recopiladas y las del DCECR.

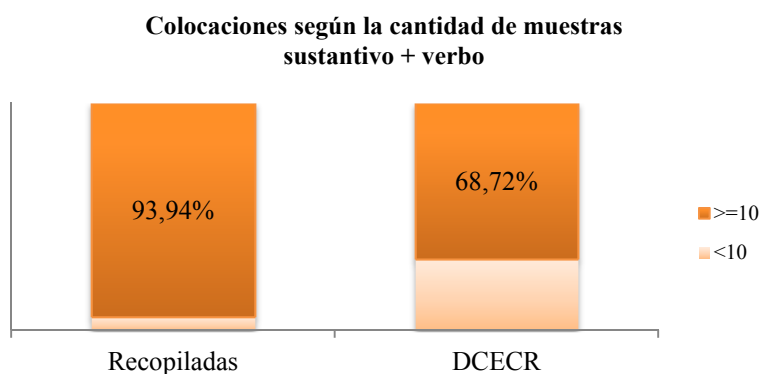


Gráfico 2-5 Colocaciones de los conjuntos de prueba que cumplen con la restricción del mínimo número de muestras

<b>Palabra x</b>	<b>Palabra y</b>	<b>Frec. Relat. x</b>	<b>Frec. Relat. y</b>
<i>recostadero</i>	<i>tapiz</i>	1	0,0008219178
<i>acristalamiento</i>	<i>doblar</i>	1	6,739301E-05
<i>antuvador</i>	<i>disoluta</i>	1	0,003115265
<i>antimoniuro</i>	<i>india</i>	1	0,0001105899
<i>guadianés</i>	<i>palado</i>	1	0,00175644
<i>desmelenadura</i>	<i>telaraña</i>	1	0,001273885
<i>ajedrista</i>	<i>juan</i>	1	1,958608E-05
<i>candonguero</i>	<i>confirmar</i>	1	0,0001316713
<i>caratulero</i>	<i>morar</i>	1	7,478127E-05
<i>bereque</i>	<i>decir</i>	1	1,536072E-06
<i>antuvador</i>	<i>poner</i>	1	5,858757E-06
<i>microfilmadora</i>	<i>atar</i>	1	0,000108452
<i>caciquesco</i>	<i>hazañar</i>	1	0,0002938871
<i>gomosidad</i>	<i>húmedo</i>	1	0,0001965795
<i>bereque</i>	<i>curro</i>	1	0,00123406
<i>caratulero</i>	<i>infinito</i>	1	9,875568E-05
<i>desalumbradamente</i>	<i>decir</i>	1	1,536072E-06
<i>galonista</i>	<i>novatada</i>	1	0,06818182
<i>impugnativo</i>	<i>desahogar</i>	1	0,0008186655
<i>breguero</i>	<i>valeroso</i>	1	0,0003867474
<i>desalinizadora</i>	<i>planta</i>	1	0,0001832988

Tabla 2-6 Combinaciones en el corpus con valores máximos de frecuencia relativa

Porcentajes de combinaciones según la distancia

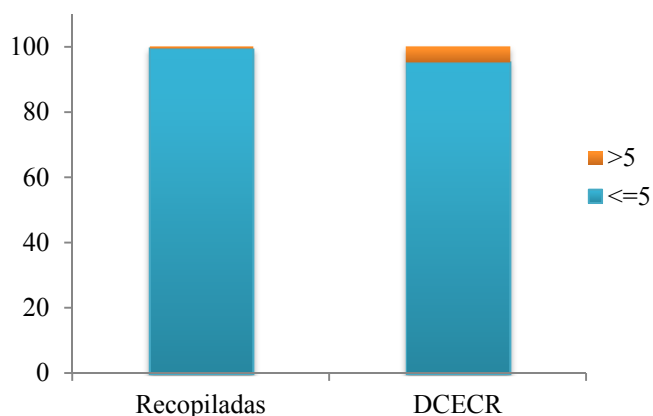


Gráfico 2-6 Distancia entre los colocados en los conjuntos de prueba

Palabra x	Palabra y	Inf. Mutua
<i>candonguero</i>	<i>candongear</i>	26,5754247590989
<i>carroñoso</i>	<i>lapidífico</i>	25,8384591649327
<i>amentáceo</i>	<i>balanífero</i>	25,8384591649327
<i>formiato</i>	<i>propilo</i>	25,8384591649327
<i>antuviador</i>	<i>piltraca</i>	25,8384591649327
<i>lairén</i>	<i>pedrojiménez</i>	25,5754247590989
<i>hiperlipidemia</i>	<i>lipemia</i>	25,5165310700453
<i>espigadilla</i>	<i>zollipar</i>	25,3530323762343
<i>amarguero</i>	<i>morrilla</i>	25,3530323762343
<i>helenización</i>	<i>panhelenismo</i>	25,1603872598201
<i>cochinita</i>	<i>pibil</i>	25,1603872598201
<i>adarguero</i>	<i>algareador</i>	25,1603872598201
<i>mandilejo</i>	<i>oracionero</i>	24,768069909176
<i>cespitoso</i>	<i>flexuoso</i>	24,7009556796546
<i>abstergente</i>	<i>restringente</i>	24,6160667596262
<i>abstergente</i>	<i>cefalálgico</i>	24,6160667596262
<i>epigénesis</i>	<i>preformismo</i>	24,5754247590989
<i>norato</i>	<i>redel</i>	24,5754247590989
<i>rupicapra</i>	<i>rupicapra</i>	24,5754247590989
<i>ruado</i>	<i>ruado</i>	24,5456775383339

Tabla 2-7 Combinaciones en el corpus con valores máximos de Información mutua



Asimismo, la revisión de los casos en que la combinación solo se da a distancias grandes revela que éstos no aportan información sobre la combinatoria de los lexemas, por lo que se desechan todas aquellas muestras de la combinación a distancias mayores que 5 (Gráfico 2-6)

En síntesis, se puede concluir que la metodología de trabajo exige fijar un número mínimo de muestras y una distancia máxima entre los elementos de la colocación. Empíricamente, se determina que si no se llega al menos a 0,0001 en la frecuencia relativa o de 1,5 en la información mutua no se deben tomar en consideración. Por otra parte, parecía suficiente indicio de colocación el disponer de al menos 10 muestras de la combinación y frecuencia relativa que supere el valor de 0,05 o la información mutua mayor que 6. Fruto de los hallazgos realizados surge la catalogación de las combinaciones extraídas en: combinaciones libres, casos dudosos y colocaciones aplicando los criterios expuestos anteriormente y que se resumen en la Tabla 2-8.

	<b>Frec. absoluta</b>	<b>Frec. Relativa</b>	<b>Inf. mutua</b>
<b>Combinaciones libres</b>		$\leq 0,0001$	O $\leq 1,5$
<b>Casos dudosos</b>	< 10	O $> 0,0001$ Y $< 0,05$	Y $> 1,5$ Y $< 6$
<b>Colocaciones</b>	$\geq 10$	Y $\geq 0,05$	O $\geq 6$

Tabla 2-8. Catalogación de las combinaciones del corpus.

El grupo de combinaciones que ni se aceptan como colocaciones ni como combinaciones libres será identificado como combinaciones sin catalogar. Entre ellas se pueden encontrar combinaciones libres,

*intimidar-ciclo*

*cobrar-madera*

y colocaciones,

*enternecer-caricia*

*acunar-sueño*

*estirar-cabello*

Si se aplican estas reglas sobre las colocaciones con estructura *sustantivo + verbo* se observa la gran cantidad de combinaciones rechazadas como colocaciones (Gráfico 2-7) y que el conjunto de combinaciones que se identifican por ambos indicadores tiene un alto grado de coincidencia (Gráfico 2-8).

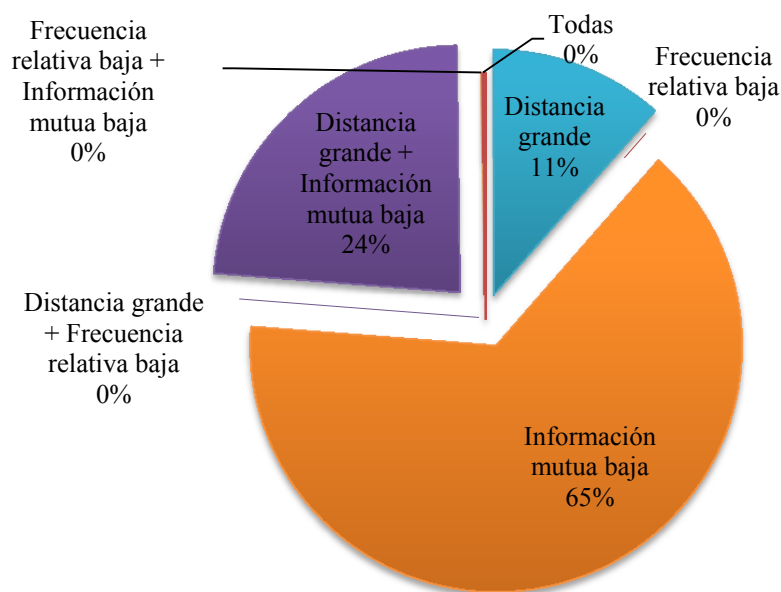


Gráfico 2-7 Combinaciones sustantivo + verbo que no cumplen ninguna de las reglas establecidas

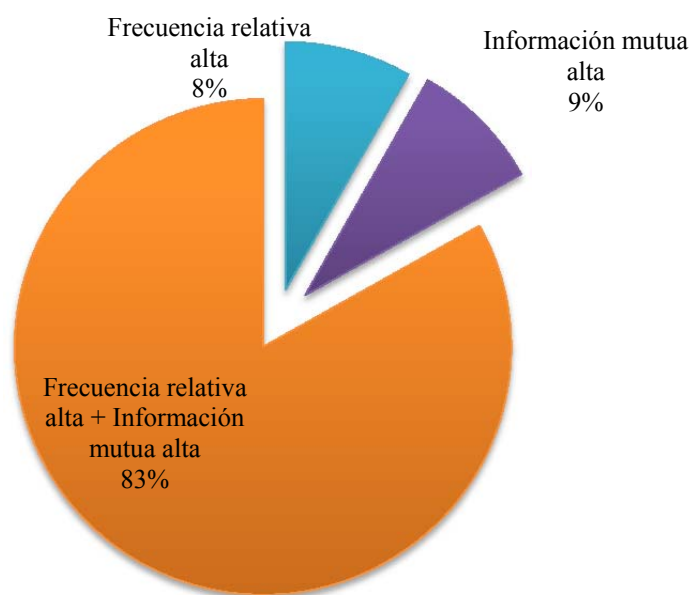


Gráfico 2-8 Combinaciones sustantivo + verbo aceptadas según las reglas establecidas.

#### 2.4.2. Test estadísticos

Se analizan los datos obtenidos para el z-score, t-score y test de Dunning, los rangos en que varían tanto el z-score como el t-score no presentan gran diferencia entre las colocaciones recopiladas o las del DCECR. Este resultado parece útil para automatizar la catalogación de combinaciones como colocaciones si pertenecen a dicho intervalo respetando un cierto margen (Tabla 2-9).

	<b>z-score</b>	<b>t-score</b>	<b>Dunning</b>
<b>Recopiladas</b>	[-55,41 695,99]	[-1,57 127,55]	[-2293,62 1904484,55]
<b>D. CC .RR.</b>	[-72,64 695,99]	[-2,38 197,37]	[357,84 19664670]
<b>CORPUS</b>	[-172,32 5539,22 ]	[-28,66 304,99]	[-36533,65 19664695,48]

Tabla 2-9 Rangos de variación de los estadísticos z-score, t-score y test de Dunning

Son pocas las combinaciones cuyo valor del z-score superan el umbral marcado por los conjuntos de ensayo, en total 567. En este grupo se recogen anomalías como *matarile-rilar* o bien palabras del mismo campo semántico como *fulano-mengano*, *iza-rabiza*, *antineutrón-antiprotón* que también presentan como característica la baja frecuencia de los elementos que las componen. En este rango del z-score, al ceñirse a los casos en que las formas canónicas constituyentes tienen frecuencias altas, aparece una alta densidad de colocaciones (Tabla 2-10).

<b>z-score</b>	<b>FC x</b>	<b>Frec x</b>	<b>FC y</b>	<b>Frec y</b>
876,01	<i>administrativo</i>	8842	<i>contencioso</i>	1505
1189,47	<i>ceja</i>	10904	<i>enarcar</i>	1620
1160,71	<i>ceño</i>	11392	<i>fruncido</i>	3040
2388,16	<i>ceño</i>	11392	<i>fruncir</i>	14661
700,07	<i>distrito</i>	12837	<i>federal</i>	10248
796,14	<i>barro</i>	15357	<i>recreativo</i>	1639
719,03	<i>archivo</i>	15515	<i>protocolo</i>	3227
795,52	<i>columna</i>	19662	<i>vertebral</i>	965
1673,44	<i>estado</i>	40175	<i>unido</i>	42708
1226,58	<i>estado</i>	40175	<i>unir</i>	79751
711,36	<i>presidenta</i>	47988	<i>república</i>	42959
711,36	<i>presidente</i>	47988	<i>república</i>	42959
1407,71	<i>hombro</i>	61358	<i>encoger</i>	20376
770,38	<i>puerta</i>	203887	<i>abrir</i>	207390

Tabla 2-10 Combinaciones cuyos elementos aparecen con elevada frecuencia y z-score alto.

La representación de la frecuencia de la combinación frente al valor alcanzado por el z-score muestra que valores altos del indicador se dan en todos los espectros de frecuencia, tanto si ésta es alta como si no (Gráfico 2-9). Otro aspecto a destacar en el comportamiento del z-score es que los valores negativos corresponden a colocaciones funcionales (Tabla 2-11).

Recopiladas			DCECR.		
z-score	FC A	FC B	z-score	FC A	FC B
-27,29	<i>embargo</i>	<i>levantar</i>	-36,75	<i>profundo</i>	<i>saber</i>
-27,63	<i>tristeza</i>	<i>tener</i>	-38,17	<i>clase</i>	<i>dar</i>
-28,05	<i>locura</i>	<i>tener</i>	-38,13	<i>clases</i>	<i>dar</i>
-28,40	<i>reflejo</i>	<i>tener</i>	-38,754	<i>sitio</i>	<i>dar</i>
-28,43	<i>función</i>	<i>tener</i>	-39,15	<i>aire</i>	<i>dar</i>
-29,73	<i>angustia</i>	<i>tener</i>	-39,21	<i>amor</i>	<i>dar</i>
-30,38	<i>solución</i>	<i>hacer</i>	-40,21	<i>fuerza</i>	<i>dar</i>
-30,81	<i>término</i>	<i>dar</i>	-40,25	<i>imagen</i>	<i>dar</i>
-32,49	<i>sueño</i>	<i>tener</i>	-40,77	<i>vivo</i>	<i>dar</i>
-32,78	<i>idea</i>	<i>llevar</i>	-41,42	<i>colores</i>	<i>dar</i>
-32,89	<i>belleza</i>	<i>tener</i>	-41,59	<i>color</i>	<i>dar</i>
-35,00	<i>pasión</i>	<i>tener</i>	-42,92	<i>voz</i>	<i>dar</i>
-35,68	<i>terror</i>	<i>tener</i>	-43,00	<i>espacio</i>	<i>dar</i>
-35,76	<i>fin</i>	<i>dar</i>	-43,29	<i>atenciones</i>	<i>dar</i>
-37,36	<i>metro</i>	<i>ir</i>	-43,29	<i>atención</i>	<i>dar</i>
-41,35	<i>problema</i>	<i>señor</i>	-49,53	<i>pie</i>	<i>dar</i>
-41,35	<i>problema</i>	<i>señor</i>	-51,97	<i>guerra</i>	<i>dar</i>
-42,97	<i>relación</i>	<i>hacer</i>	-53,65	<i>forma</i>	<i>dar</i>
-46,09	<i>bajo</i>	<i>hombre</i>	-58,67	<i>cuerpo</i>	<i>dar</i>
-52,17	<i>causa</i>	<i>tener</i>	-64,27	<i>alto</i>	<i>dar</i>
-55,41	<i>amor</i>	<i>tener</i>	-68,23	<i>ver</i>	<i>dar</i>
			-72,62	<i>tiempo</i>	<i>dar</i>

Tabla 2-11 Combinaciones con z-score negativo

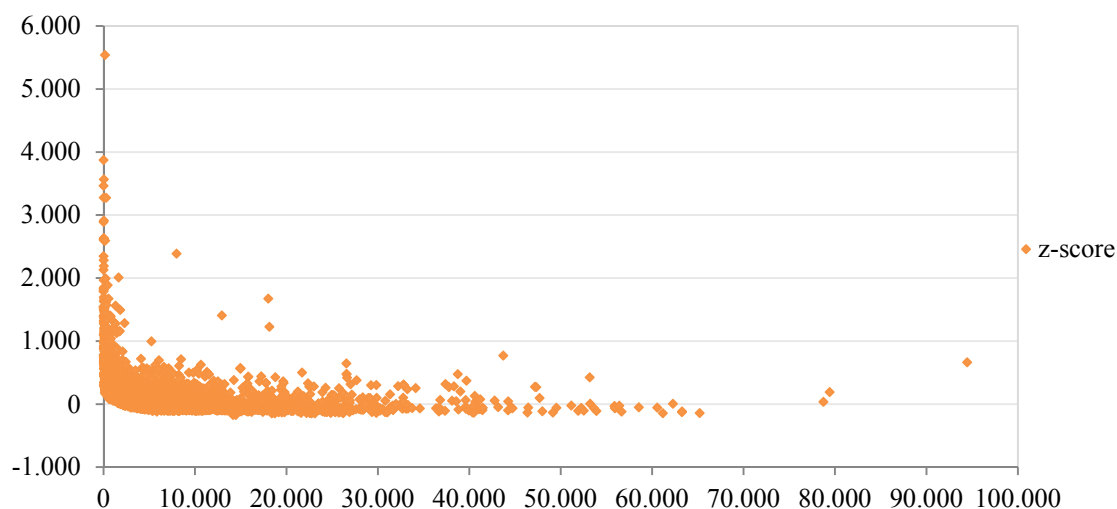


Gráfico 2-9 Valores del z-score según la frecuencia de la combinación

En contraste con los resultados observados para el *z-score* ocurre que los valores negativos del *t-score* se pueden considerar prácticamente inexistentes entre las colocaciones recopiladas y las del DCECR: 2 y 93 respectivamente. En el total del corpus se observa que estos casos se obtienen en combinaciones de frecuencias mínimas. También destaca la dependencia directa entre el valor del *t-score* y la cantidad de veces que aparece la combinación, a mayor frecuencia mayor *t-score*. Por tanto, no se recomienda establecer un ranking en el que se entremezclen palabras poco usadas con aquellas de frecuencia de aparición alta. En la Tabla 2-12 se resalta, a modo de ejemplo, una combinación libre que queda igual o mejor catalogada que colocaciones cuya frecuencia de aparición es igual o menor a la de ella. Este error está inducido por utilizar un único ranking entre todas las combinaciones en un texto o corpus dado, basado en este indicador.

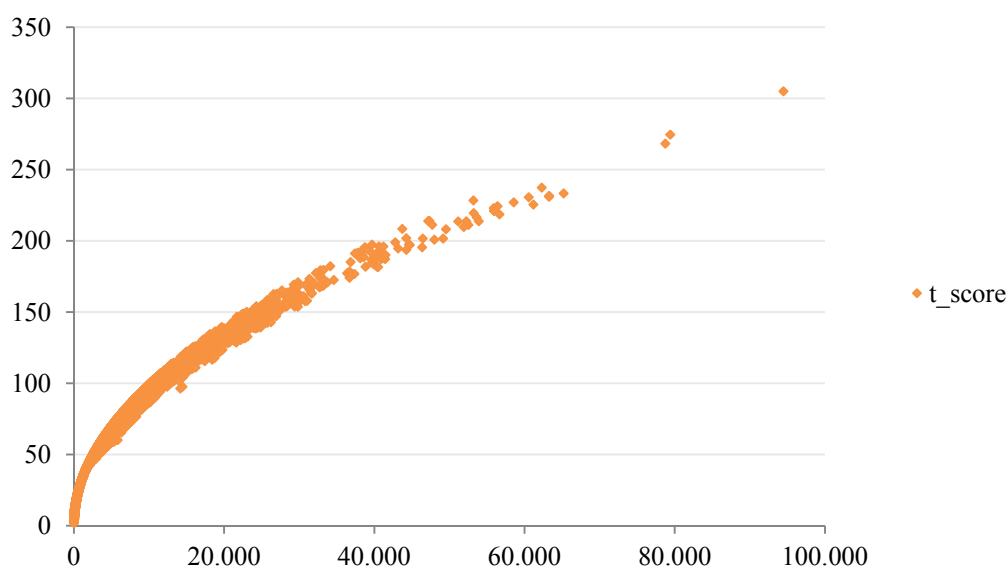


Gráfico 2-10 Valores máximos de t-score entre las combinaciones agrupadas por frecuencia de aparición.

Pal x	Pal y	Frec (x,y)	t-score
dar	discípulo	681	22,8
fuerte	emoción	426	20,29
dar	emoción	681	23,86
controlar	emoción	165	12,73
dominar	emoción	390	19,6
prestar	colaboración	143	11,86

Tabla 2-12 Influencia de la frecuencia de la combinación en el t-score

Respecto al *test de Dunnig*, se atenúa la influencia de la cantidad de muestras de una combinación en la puntuación que se obtiene. Este resultado es de esperar, ya que está diseñado para hacerlo más estable ante las variaciones en la cantidad de muestras. En la Tabla 2-13 se observan las colocaciones recopiladas para el verbo *ir*, destacando el amplio rango de variación de las frecuencias computadas para cada una de ellas, y la baja variabilidad del *test de Dunning*.

FC x	FC y	Frec(x, y)	t-Dunning
coche	ir	3088	11394028,69
metro	ir	669	11432604,93
automóvil	ir	250	11441127,50
automóvil	ir	250	11441127,50
taxi	ir	249	11441428,53
autobús	ir	223	11441910,29
detrimiento	ir	60	11444571,56
descrédito	ir	7	11445482,55
desdoro	ir	7	11445487,50

Tabla 2-13 Test de Dunning evaluado en colocaciones recopiladas para el verbo ir

Para finalizar, se aprecia inviable la evaluación del *test de Poisson* en un corpus tan extenso. Se recuerda que la expresión incluye valores de  $k!$ , siendo  $k$  la frecuencia de la combinación.

### 2.4.3. Análisis por tipos de combinación según la categoría gramatical

Un segundo análisis de los resultados obtenidos se llevó a cabo según el tipo de combinación. Los grupos se han establecido según la categoría gramatical de las formas canónicas de los elementos que intervienen en las mismas, teniéndose en cuenta los casos:

*sustantivo + verbo*

*sustantivo + adjetivo*

*verbo + adverbio.*

### 2.4.3.1. Combinaciones sustantivo + verbo

De ellas se reconocen aproximadamente un 55% de las colocaciones recopiladas que se encuentran en el corpus y un 23 % de las que se encuentran en el DCECR. De las colocaciones recopiladas ninguna infringe el requisito de la frecuencia relativa, sin embargo 11 solo aparecen a una distancia de más de 5 palabras, 18 no superan el umbral de la información mutua y 153 no cumplen el requisito de la frecuencia absoluta mínima. En cuanto a las combinaciones del DCECR., 4 se descartarían por no llegar a la frecuencia relativa mínima, 1065 por no superar el valor impuesto de 1,5 a la información mutua, 1332 casos solo aparecen a distancia alta y 9995 no cumplen el criterio de contar con al menos 10 muestras de las mismas en el corpus (Tabla 2-14).

	Distancia > 5	Frec. Relativa < 0,05	Inf. Mutua < 1,5	Frec. Total ≤ 10
<b>Recopiladas</b>	11	-	18	153
<b>DCECR.</b>	1 332	4	1 065	9 995

Tabla 2-14 Colocaciones que infringen las restricciones determinadas en los conjuntos de prueba

Se comprueban los rangos de variación de *z-score*, *t-score* y *test de Dunning* evaluados para cada uno de los grupos de ensayo; en todos los casos, de uno u otro modo, se observa la influencia de la frecuencia de uso de los elementos de la combinación con la puntuación obtenida: valores mínimos de *z-score* para verbos como *hacer*, *tener*, *formar*, *decir*, etc. o valores máximos de *t-score* cuando ambos elementos de la combinación tienen una frecuencia alta de aparición (Tabla 2-15), o el caso del *test de Dunning*, donde ocurre que los puestos más altos en el ranking corresponden a las llamadas colocaciones funcionales producidas por verbos de elevada frecuencia de uso:

*hacer, dar, tener, ir, poner, haber, ... + sustantivo*

	<i>z-score</i>	<i>t-score</i>	<i>Dunning</i>
<b>Recopiladas</b>	[-55,41 695,99]	[-0,66 127,55]	[83,53 1904484]
<b>D. CC .RR.</b>	[-72,64 695,99]	[-2,38 197,37]	[1712,65 19664670]
<b>Corpus</b>	[-172,32 5539,22 ]	[-20,31 268,26]	[83,53 19044484,55]

Tabla 2-15 rangos de variación de los estadísticos en los conjuntos de prueba y en el corpus

Se puede considerar que los criterios establecidos en el caso de las frecuencias relativas y la información mutua para descartar combinaciones son efectivos; sin embargo, la cobertura del método en cuanto a la catalogación de las colocaciones resulta insuficiente. Se concluye además que cuando se usa cualquiera de los indicadores estudiados no parece muy útil la comparación mediante rankings, o valores de corte sobre el total del corpus, ya que su comportamiento queda determinado por la cantidad de muestras de la combinación y por la frecuencia de aparición de los elementos que la componen, es decir, el uso que se le da a una palabra. No se puede comparar el valor de un indicador basado en la frecuencia entre los verbos *dar* (1118012), *desempeñar* (18903), o *traspasar* (100).

<i>z-score</i> , valores máximos		<i>z-score</i>	
Frecuencia Combinación $\geq 10$		valores mínimos	
<i>aranzada majolar</i>	<i>nobel premiar</i>	<i>capitán poder</i>	<i>antiguo poder</i>
<i>ruga rugar</i>	<i>monóxido carbonar</i>	<i>vuelto poder</i>	<i>amor hacer</i>
<i>puerta cerrar</i>	<i>pistilo estambrar</i>	<i>negro hacer</i>	<i>ante tener</i>
<i>cloruro clorurar</i>	<i>rodilla hincar</i>	<i>bienes seguir</i>	<i>llama poder</i>
<i>vito vitar</i>	<i>dactilar hollar</i>	<i>bien seguir</i>	<i>siglo tener</i>
<i>cigarro fumar</i>	<i>viento soplar</i>	<i>alrededores decir</i>	<i>hombre formar</i>
<i>i her</i>	<i>café tazar</i>	<i>alrededor decir</i>	<i>año parecer</i>
<i>mamante plantar</i>	<i>salce pitar</i>	<i>alto querer</i>	<i>propio dar</i>
<i>holístico agenciar</i>	<i>diamela engraciar</i>	<i>artículo saber</i>	<i>negro poder</i>
<i>hidrato carbonar</i>	<i>atenciones prestar</i>	<i>mesa decir</i>	<i>junto saber</i>
<i>sant sacramentar</i>	<i>atención prestar</i>	<i>gobierno ver</i>	<i>voz ir</i>
<i>vistazo echar</i>	<i>iris arcar</i>	<i>español hacer</i>	<i>colores poder</i>
		<i>vida llegar</i>	<i>capitana poder</i>

Tabla 2-16 combinaciones del corpus con valores extremos del *z-score*

#### 2.4.4. Combinaciones sustantivo + adjetivo

En este grupo se incluyó el total de combinaciones del tipo *sustantivo – adjetivo* o *sustantivo – sustantivo (que puede usarse como adjetivo en el corpus)*. El resumen de los datos registrados revela ligeras diferencias con las correspondientes a la estructura *sustantivo + verbo*. Se destaca que la distancia colocacional no debe restringirse a 5, las frecuencias relativas alcanzan valores mayores cuando se evalúan respecto al adjetivo en lugar del sustantivo (Gráfico 2-11). Entre las coincidencias en el comportamiento se observa la necesidad de exigir un número mínimo de



muestras, el que las combinaciones con frecuencia relativa 1 se deben descartar, el valor del *t-score* está directamente relacionado con la cantidad de veces que aparece la combinación en el corpus, y que el rango de variación de la información mutua es similar entre las colocaciones recopiladas y las del DCECR.

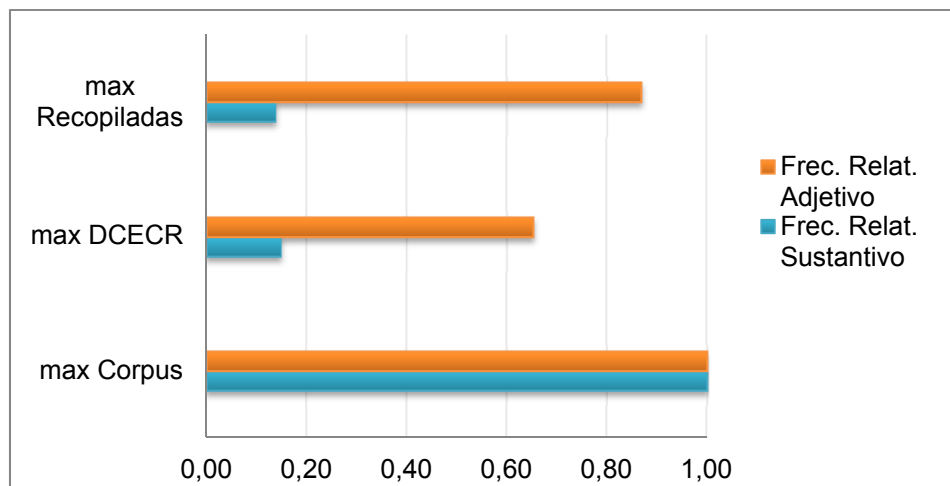


Gráfico 2-11 Frecuencias relativas en las combinaciones Sustantivo + Adjetivo

	Corpus	Recopiladas	DCECR.
<b>Frec. Total <math>\geq 10</math></b>	41,34%	87,56%	57,58%
<b>Frec. Relat. Sust.</b>	[0,000014 1]	[0,000012 0,13826]	[0,0000051 0,14894]
<b>Frec. Relat. Adj.</b>	[0,000005 1]	[0,00015 0,86667]	[0,000014 0,65165]
<b>Inf. Mutua</b>	[-3,41759 25,83846]	[-0,71624 15,45569]	[-0,73025 16,50326]
<b>z-score</b>	[-104,74 2634,92]	[-46,1 607,86]	[-17,32 335,7]
<b>t-score</b>	[-16,78 191,71]	[1,57 100,58]	[-1,61 40,41]
<b>Dunning</b>	[-19846,94 971272,33]	[-2293,62 6740597,91]	[357,85 5320378,84]

Tabla 2-17 Valores obtenidos por las combinaciones Sustantivo + Adjetivo.

En el caso del *t-score* se entremezclan numerosas combinaciones libres que obtienen buenas puntuaciones por ser palabras de uso frecuente en lengua española:

*madre padre, hijo parte, don señor, aire bueno*

### 2.4.5. Combinaciones verbo + adverbio

No se encuentra ningún caso del conjunto de colocaciones recopiladas en este grupo, por lo que el contraste se realiza entre las 2 148 extraídas del DCECR., y las 143 578 encontradas en el corpus. Los rangos de variación de los indicadores presentan comportamientos similares al caso *sustantivo + adjetivo*. Se observa que las frecuencias relativas son mayores cuando se evalúan respecto al adverbio en lugar del verbo, se debe exigir un número de muestras mínimo y se aconseja utilizar los rangos obtenidos en el conjunto de pruebas añadiéndoles un cierto margen. Destaca el caso del test de Dunning, donde se aprecia la influencia de la frecuencia de aparición del verbo en la puntuación obtenida: las combinaciones mejor posicionadas corresponden al verbo *decir*, las siguientes al verbo *hacer*, etc. Este resultado nuevamente revela cierta dependencia de los valores calculados con la cantidad de muestras de la forma canónica en el corpus.

	Corpus	D. CC.RR.
<b>Frec. Total <math>\geq 10</math></b>	37,78%	55,87%
<b>Frec. Relat. Adv.</b>	[5,907373E-05 1]	[0,00018 0,78588]
<b>Frec. Relat. Verbo</b>	[1,536072E-06 0,3571429]	[2,987902E-06 0,08926]
<b>Inf. Mutua</b>	[-2,00833 21,83846]	[0,39011 13,18894]
<b>z-score</b>	[-48,86878 749,99601]	[-17,71749 265,77742]
<b>t-score</b>	[-5,23625 96,45938]	[0,47386 42,63691]
<b>Dunning</b>	[228,8713 19664695,48157]	[8348,81265 19664515,42659]

Tabla 2-18 Valores obtenidos por las combinaciones verbo + adverbio

<b>Frec. Rel</b>	<i>mirar fijamente, oponer diametralmente, cerrar herméticamente, mirar fijamente, enamorar perdidamente, mirar interrogativamente, poder lícitamente, vivir maritalmente, abrir desmesuradamente, vivir disolutamente, consumir improductivamente, recibir hospitalariamente, girar preceptivamente, realizar preceptivamente, disponer preceptivamente, asignar preceptivamente, llorar desconsoladamente, pecar venialmente, llorar inconsolablemente</i>
<b>Inf. Mutua</b>	<i>enjaezar vistosamente, legalizar solidariamente, enjaezar ricamente, amueblar vistosamente, guarnir ricamente, emanar patrióticamente, asignar preceptivamente, absolver sacramentalmente, departir amigablemente, encuadernar lujosamente, enjaezar lujosamente, infringir civilmente, lisar llanamente, amueblar suntuosamente, reclinar muellemente</i>
<b>z-score</b>	<i>mirar fijamente, lisar llanamente, guarnir ricamente, oponer diametralmente, enamorar perdidamente, enjaezar ricamente, amueblar lujosamente, cerrar herméticamente, iluminar brillantemente, participar activamente, ataviar ricamente, emanar patrióticamente, ligar íntimamente, sudar copiosamente</i>
<b>t-score</b>	<i>mirar fijamente, poder fácilmente, poder solamente, poder realmente, ver claramente, decir realmente, decir solamente, tener solamente, tener realmente, hacer solamente, hacer realmente, decir propiamente, decir exactamente, saber exactamente, hacer exactamente, decir finalmente, saber perfectamente,</i>
<b>Dunning</b>	<i>decir macarrónicamente, decir rasamente, decir chuscamente, decir aficionadamente, ...</i>

Tabla 2-19 Combinaciones Verbo + Adverbio mejor puntuadas por cada indicador

## 2.5. Influencia del tamaño del corpus

A la vista de la dependencia que muestran los resultados con respecto al uso que se hace de las palabras en el corpus, es necesario comprobar en qué medida la estrategia queda determinada por el tamaño del corpus. Para ello, se evalúan los datos utilizando un corpus sensiblemente menos extenso que compila una colección de novelas de D. Benito Pérez Galdós, con un total de 2299920 palabras.

Los rangos de variación de todos los indicadores se modifican sustancialmente, comprobándose que los valores umbral se deben alterar en cuanto se cambia el tamaño del corpus. Se manifiesta de nuevo la influencia de la cantidad de muestras que es posible recolectar de una combinación sobre los valores que marcan el límite entre lo que es o no colocación. Sería necesario revisar los puntos de corte, incluido el del número mínimo de muestras que se considere fiable. Una forma de resolverlo consiste en la revisión de los datos que se obtengan y establecerlo de forma empírica, al igual que se hizo en el epígrafe anterior. Otro método es recurrir a un conjunto de entrenamiento, en tal caso las colocaciones recopiladas o el DCECR,

con el inconveniente de la dependencia de disponer de ellas y de que el corpus que se procese albergue suficientes muestras de las mismas (Gráfico 2-12).

	<i>Frec. Relativa</i>	<i>Inf. Mutua</i>	<i>z-score</i>	<i>t-score</i>	<i>Dunning</i>
<b>Recopiladas</b>	[0,0049 0,7846]	[0,28 13,13]	[-57,07 13,96]	[-46,72 5,93]	[-26,71 46634]
<b>D. CC .RR.</b>	[0,0028 1]	[0,5 13,26]	[-61,11 10,98]	[-39,79 14,54]	
<b>Corpus</b>	[0,00092 1]	[-0,99 18,19]	[-136 76,48]	[-116,19 14,54]	

Tabla 2-20 Rangos de variación de los indicadores en el corpus de Galdós

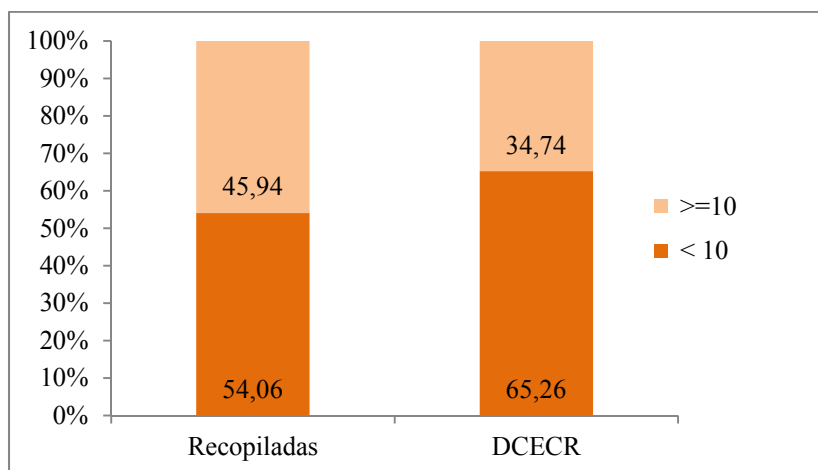


Gráfico 2-12 Muestras de las colocaciones en los grupos de ensayo en el corpus de Galdós

## 2.6. Las colocaciones como combinaciones atípicas

A la vista de los resultados obtenidos se hace obvia la necesidad de modificar la estrategia basada únicamente en umbrales de corte para delimitar la frontera entre colocaciones y combinaciones libres. El objetivo último es obtener una herramienta que permita automatizar en lo posible la extracción de las colocaciones presentes en un corpus con independencia de su tamaño. Para ello, se tendrá en cuenta la característica de preferencia presente en el concepto de colocación, según la cual parece lógico que, fijado un colocado, el análisis de sus capacidades combinatorias se realice entre las muestras obtenidas de él, y no sobre el total de combinaciones extraídas del corpus. De esta forma se evita la comparación entre palabras distintas desde el punto de vista de sus frecuencias de uso, cuyas diferencias aumentan a medida que el corpus sea más extenso.

Bajo estas hipótesis, se pretende construir un indicador que detecte aquellas combinaciones cuyo uso destaque con respecto a lo que es habitual, en contraposición a la idea subyacente en los estadísticos analizados en las secciones previas, en las que la decisión se toma en términos de rechazar la independencia en el uso de los colocados. En este sentido, se considera que las frecuencias relativas serán fundamentales, ya que aportan información del uso de una determinada combinación con respecto a lo que ha sido utilizada una determinada

---

palabra en el corpus sobre el que se trabaje. Así las cosas, la hipótesis se centra en qué valores atípicos de ésta, en relación a lo que es común en su rango de variación, proporcionarán indicios del fenómeno de la preferencia. El uso de este concepto permite contrastar los datos en términos relativos al uso que se da en el corpus a una palabra concreta, con mayor independencia respecto a cuánto haya sido usada en el corpus de referencia.

Siguiendo esta metodología se asume que cada palabra tiene su propia muestra, y se considera qué casos son atípicos dentro de ella, frente a considerar el corpus como una única muestra en la que se detecta dónde se rechaza la independencia de los colocados. Esto, además, permite alejarse del problema de establecer valores de corte dependientes del tamaño del corpus o de establecer rankings en los que se comparan datos que dependen directamente de lo frecuente que es el uso de las palabras implicadas en la combinación.

Este epígrafe se centra a partir de ahora en la experimentación con técnicas que tradicionalmente se usan para la identificación de valores atípicos en muestras. En estadística se utiliza el término *outlier* para referirse a los datos de una muestra que parecen ser inconsistentes con el resto del conjunto, es decir, son valores que parecen demasiado grandes o demasiado pequeños en comparación con las demás observaciones.

Un outlier es una observación que se desvía tanto de las restantes como para sospechar que fue generada por un mecanismo diferente. (Aggarwal, 2013).

La hipótesis que se plantea consiste en considerar las colocaciones como anomalías en el mecanismo de generación de combinaciones, se presupone que el fenómeno de generación origina combinaciones libres y que la característica de la preferencia determina la diferencia en la producción de las colocaciones.

Se analizan 3 estrategias para la detección de estos elementos singulares ampliamente utilizadas en otros ámbitos: método basado en la *desigualdad de Chebyshev*, método de los *cuartiles* y por último el método basado en la *MEDA*. En el primero se determinan los outliers a partir de la media muestral de las frecuencias mientras que en los otros dos se hace a partir de los valores de los cuartiles de la muestra. Cabe destacar que todos ellos proporcionan un criterio objetivo para establecer un punto de corte respecto a las combinaciones que se deben considerar, frente a las técnicas propuestas anteriormente, en la que la frontera se establece de forma aproximada y es totalmente dependiente del tamaño del corpus.

### 2.6.1. Método basado en la desigualdad de Chebyshev

La desigualdad de Chebyshev garantiza que los datos de una muestra verifican:

$$fr(i: |x_i - \bar{x}| > ks) < \frac{1}{k^2}$$

Siendo  $\bar{x}$  la media muestral y  $s$  la desviación típica de la muestra. Aplicado a distribuciones normales, significa que se cumplen las relaciones de la Tabla 2-21

% de datos que verifican	
68%	$fr(i:  x_i - \bar{x}  < s)$
95%	$fr(i:  x_i - \bar{x}  < 2s)$
99%	$fr(i:  x_i - \bar{x}  < 3s)$

Tabla 2-21 Aplicación de la desigualdad de Chebyshev

Sobre la base de estas propiedades se justifica considerar atípicas las observaciones que verifican:  $|Z| > 3$ , siendo:

$$Z = \frac{x_i - \bar{x}}{s}$$

De este modo se estarían limitando las anomalía al 0,13% de los casos (Leys *et al.*, 2013), algunos investigadores adoptan un criterio menos restrictivos, considerando una desviación de la media de 2 o 2,5 veces la desviación de la media.

En este trabajo: fijada una palabra, los valores  $x_i$  corresponden a las frecuencias relativas con las restante con las que combina y  $\bar{x}$ ,  $s$  a la media muestral y desviación típica de dichos datos respectivamente, por lo que el estadístico que se usa:

$$Z_{Chebyshev}(X_i, Palabra) = \left| \frac{FrecRel(X_i, Palabra) - \overline{FrecRel(Palabra)}}{S_{FrecRel(Palabra)}} \right|$$

Las combinaciones  $(X_i, Palabra)$  tales que  $Z_{Chebyshev}(X_i, Palabra) > 3$  serán aceptados como colocaciones.

La Tabla 2-22 y la Tabla 2-23 resumen los resultados de aplicar esta estrategia sobre las combinaciones en el corpus completo y el corpus de Galdós, mostrando que el requisito que se exigió a las combinaciones para considerarlas colocaciones fue excesivamente restrictivo. No obstante, se aprecia que se filtran notablemente la gran cantidad de combinaciones libres que se introducen erróneamente en los catálogos generados por los restantes métodos. Se muestran,

como ejemplos, combinaciones que superan la restricción de Chebyshev en los 3 conjuntos de ensayo. Se presentan tanto las de mejor como las de peor puntuación (Tabla 2-24).

	<b>Rango</b>	<b>Combinaciones con <i>ZChebyshev</i> ≥ 3</b>
<b>Recopiladas</b>	[0 22,289966]	933
<b>DCECR</b>	[0 78,607647]	2735
<b>Corpus</b>	[0 78,607647]	455738

Tabla 2-22 Resultados de *ZChebyshev* en el corpus completo

	<b>Rango</b>	<b>Combinaciones con <i>ZChebyshev</i> ≥ 3</b>
<b>Recopiladas</b>	[0 4,35341]	45
<b>DCECR</b>	[0 11,650217]	83
<b>Corpus</b>	[0 16,537422]	2922

Tabla 2-23 Resultados de *ZChebyshev* en el corpus de Galdós

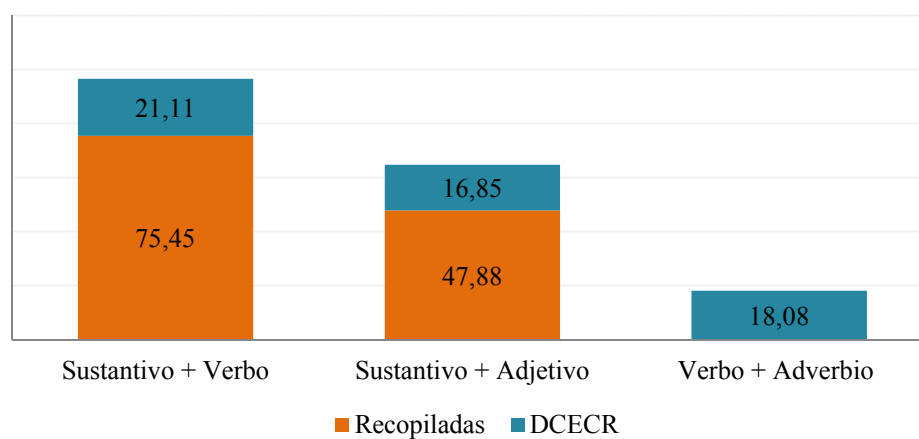


Gráfico 2-13 Combinaciones de los conjuntos de ensayo que superan el criterio de *ZChebyshev* ( $Z\text{-Chebyshev} \geq 3$ )

	Recopiladas		DCECR		Corpus	
<b>ZChebyshev <math>\geq 3</math></b>	<i>dar</i>	<i>golpe</i>	<i>dar</i>	<i>vuelta</i>	<i>dar</i>	<i>vuelta</i>
<b>Max</b>	<i>hacer</i>	<i>pregunta</i>	<i>dar</i>	<i>paso</i>	<i>don</i>	<i>doña</i>
	<i>poner</i>	<i>marcha</i>	<i>dar</i>	<i>gracia</i>	<i>estado</i>	<i>unido</i>
	<i>tener</i>	<i>intención</i>	<i>tomar</i>	<i>decision</i>	<i>tener</i>	<i>razón</i>
	<i>tomar</i>	<i>decision</i>	<i>dar</i>	<i>salto</i>	<i>hacer</i>	<i>falta</i>
	<i>dar</i>	<i>salto</i>	<i>dar</i>	<i>orden</i>	<i>tener</i>	<i>miedo</i>
	<i>latir</i>	<i>corazón</i>	<i>soplar</i>	<i>viento</i>	<i>decir</i>	<i>verdad</i>
<b>ZChebyshev <math>\geq 3</math></b>	<i>alisar</i>	<i>cabello</i>	<i>sumir</i>	<i>cuidado</i>	<i>cavado</i>	<i>fosa</i>
<b>Min</b>	<i>alisar</i>	<i>pelo</i>	<i>reponer</i>	<i>fuerza</i>	<i>alternativo</i>	<i>formar</i>
	<i>salir</i>	<i>crisis</i>	<i>camino</i>	<i>accidentado</i>	<i>llamada</i>	<i>recibir</i>
	<i>viajar</i>	<i>coche</i>	<i>perder</i>	<i>control</i>	<i>dar</i>	<i>evento</i>
	<i>dar</i>	<i>taconazo</i>	<i>poner</i>	<i>objection</i>	<i>día</i>	<i>evento</i>
	<i>haber</i>	<i>relación</i>	<i>dar</i>	<i>ubicación</i>	<i>torre</i>	<i>elevado</i>
	<i>poner</i>	<i>crisis</i>	<i>reinar</i>	<i>tranquilidad</i>	<i>enemigo</i>	<i>padre</i>

Tabla 2-24 Ejemplos de combinaciones que superan el criterio de Chebyshev

La Tabla 2-25 y el Gráfico 2-13 recogen el resumen de los datos según la estructura de la colocación, se muestran los rangos alcanzados y los porcentajes de combinaciones que serían reconocidas según este criterio. Se aprecia escasa cobertura entre las combinaciones del DCECR. Sin embargo, los listados que se obtienen se consideran mejores desde el punto de vista de la precisión, logrando filtrar la gran cantidad de ruido en el proceso automatizado.

	<b>ZChebyshev</b>
<b>Verbo + Sustantivo</b>	[0 78,6076]
<b>Sustantivo + Adjetivo</b>	[0 68,8782]
<b>Verbo + Adverbio</b>	[0 41,06655]

Tabla 2-25 Rango de ZChebyshev según la estructura de la combinación

### 2.6.2. Métodos robustos

Una fuente de error en los métodos de identificación de outliers la constituye precisamente el uso de los datos anómalos en el cálculo de la media aritmética y de la desviación típica usados para evaluar los estadísticos. El valor de la media y la desviación típica muestral están fuertemente influenciados por la presencia de outliers en la muestra, además la desigualdad de Chebyshev se cumple para distribuciones normales. El siguiente ejemplo ilustra estos problemas. Sea la muestra en (1), en la que el valor 1000 claramente es atípico, no sería



detectado mediante el indicador  $Z_{Chebyshev}$  como tal debido a la distorsión que él mismo provoca en los estimadores muestrales que se usan en su cálculo. Este ejemplo (Leys *et al.*, 2013) ilustra los problemas que se originan al considerar este indicador, aunque su uso es habitual.

$$\begin{aligned} &\{1, 3, 3, 6, 8, 10, 10, 1000\} \\ &\bar{x} = 130,13 \\ &s = 328,8 \quad (1) \\ &\bar{x} + 3s = 1116,52 \\ &\bar{x} - 3s = -856,27 \end{aligned}$$

La mediana, sin embargo, es un estimador para la medida de tendencia central que es invariante a la presencia de outliers, lo que sugiere que se utilicen indicadores basados en este estadístico para resolver el problema que estamos tratando. En este epígrafe se presenta la adaptación de dos métodos de detección de outliers a la extracción automática de colocaciones, ambos basados en los cuartiles de la muestra.

### 2.6.3. Método basado en el diagrama de cajas y bigotes

En el análisis exploratorio de datos estadísticos se utilizan cálculos sencillos para resumir los datos. El **diagrama de cajas y bigotes de Tukey**, o **Box-Plot** es una herramienta para la visualización de 5 valores que aportan información sobre la tendencia, la simetría y la variación de la distribución del conjunto. Además permite localizar la existencia de outliers. En un diagrama de cajas y bigotes se representan  $Q_1, Q_2$  y  $Q_3$ , cuartiles 1, 2 y 3 de la distribución y se utilizan para delimitar la caja. Los bigotes se marcan en los puntos  $Q_1 - 1,5 \cdot RI$  y  $Q_3 + 1,5 \cdot RI$ . Estos últimos fijan la zona a partir de la que se considerarán *valores atípicos moderados* en la distribución. Finalmente, destacan los *valores atípicos extremos*, que rebasan las fronteras del intervalo  $Q_1 - 3 \cdot RI$  y  $Q_3 + 3 \cdot RI$ . Los valores en esta zona se marcan con un símbolo que los identifica como tales, un ejemplo puede verse en la Figura 2-8. En el interior de la caja se representa el valor de la mediana, si éste no aparece centrado la distribución de la variable es asimétrica.

En la Figura 2-9 se muestra el diagrama de cajas de las frecuencias relativas del verbo *palpitar* con los sustantivos en su contexto en el corpus, la mayor parte de combinaciones se concentran en la zona de frecuencia relativa inferior a 0,002 apreciándose un conjunto de valores atípicos extremos. El perfil de la representación en un diagrama de cajas y bigotes del resumen de la variable “*frecuencia relativa fijada una palabra*” ( $FrecRel(Palabra)$ ) por lo general será similar al que se muestra en la Figura 2-9. Se observa que corresponde a una variable totalmente asimétrica, y en la que aparecen valores atípicos extremos, por lo que se justifica la propuesta de utilizar indicadores más robustos frente a outliers que  $Z_{Chebyshev}$ .

media	desv. típica	Q1	Q2	Q3	n
0.0112348	0.02690609	0.003634361	0.005947137	0.009581498	167

Tabla 2-26 Estadísticos resumen de las frecuencias relativas cuando se fija el verbo palpar.

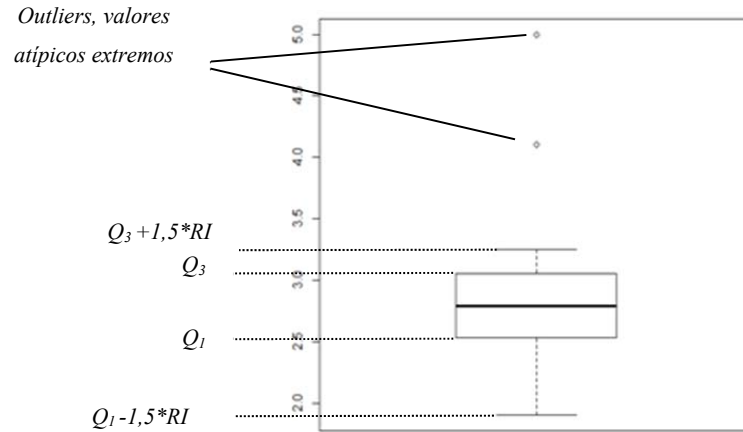


Figura 2-8 Explicación del diagrama de cajas de Tukey

Se define a continuación otro indicador basado en estas consideraciones, que identificaremos como *Box-plot(palabra,  $x_i$ )*. Este método también se encuentra en la literatura con el nombre *método de los cuartiles*. Para calcularlo, se debe fijar de antemano sobre qué palabra se calculará, y habrá que determinar el primer y tercer cuartil ( $Q1$ ,  $Q3$ ) de las frecuencias relativas respecto a la palabra de interés así como el rango intercuartílico ( $RI$ ).

$$Box - plot(palabra, x_i) = \frac{FrecRelativa(x_i, palabra) - Q3}{RI}$$

$FrecRelativa(x_i, palabra)$  es la frecuencia relativa respecto a *palabra* de  $x_i$  con *palabra*

En la Tabla 2-27, Tabla 2-28 y la Tabla 2-29 se presentan las combinaciones *sustantivo + verbo*, *sustantivo + adjetivo* y *verbo + adverbio* que alcanzan los valores máximos que se obtienen al calcular este indicador fijando el colocativo en cada uno de esos 3 casos respectivamente. La cobertura que se alcanza se recoge en el Gráfico 2-14, donde se puede apreciar cómo suben los resultados frente al método de Chebyshev.

<i>palabra</i>	<i>x<sub>i</sub></i>	<i>FrecRelativa(palabra, x<sub>i</sub>)</i>	<i>Box-plot (palabra, x<sub>i</sub>)</i>
<i>encoger</i>	<i>hombro</i>	0,6348646	993,846208698718
<i>trizar</i>	<i>hecho</i>	0,3452381	688,999867633008
<i>enguantar</i>	<i>mano</i>	0,7609428	673,999826155181
<i>abrir</i>	<i>puerta</i>	0,210825	671,584602763095
<i>enguantar</i>	<i>manos</i>	0,7485971	662,999835490524
<i>cerrar</i>	<i>puerta</i>	0,2387841	576,282633476278
<i>fruncir</i>	<i>ceño</i>	0,5443012	530,800016942488
<i>menear</i>	<i>cabeza</i>	0,617914	481,083363372164
<i>atracar</i>	<i>puerta</i>	0,6423659	462,999946849891
<i>decir</i>	<i>señora</i>	0,04031934	456,767419796828
<i>decir</i>	<i>señor</i>	0,04031934	456,767419796828
<i>pronunciar</i>	<i>palabra</i>	0,3318395	434,450978686954
<i>vezar</i>	<i>mucho</i>	0,1426878	431,531244894402
<i>cerrar</i>	<i>ojo</i>	0,1766763	426,108704142493
<i>latir</i>	<i>corazón</i>	0,3690867	412,812470856838
<i>unir</i>	<i>estado</i>	0,2274454	393,239122510575
<i>mirar</i>	<i>ojo</i>	0,07979666	382,259992494583
<i>decir</i>	<i>hombre</i>	0,03339114	378,098820930561
<i>pablar</i>	<i>san</i>	0,3634326	374,421032936214
<i>roncar</i>	<i>voz</i>	0,4186728	362,363633128052

Tabla 2-27 Combinaciones Sustantivo + Verbo con máximo valor de Box-plot

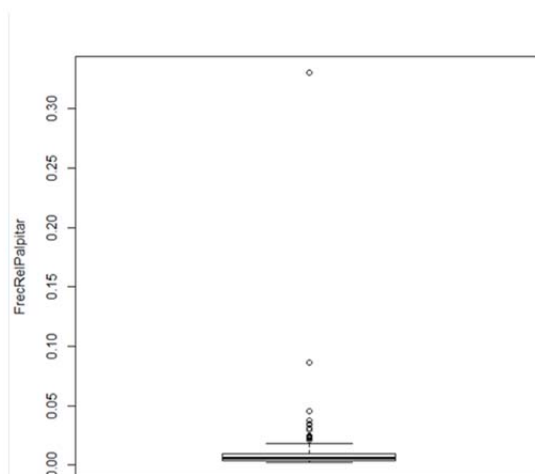


Figura 2-9 Diagrama de cajas para la distribución de las frecuencias relativas del verbo palpitar

<i>palabra</i>	$x_i$	<i>FrecRelativa(palabra, <math>x_i</math>)</i>	<i>Box-plot(palabra, <math>x_i</math>)</i>
<i>pecado</i>	<i>venial</i>	0,6516464	185,499984013151
<i>prueba</i>	<i>fehaciente</i>	0,278481	149,999992789701
<i>fuerza</i>	<i>centrifugo</i>	0,5914894	137,000007441268
<i>nota</i>	<i>discordante</i>	0,1366623	100,000013820363
<i>camino</i>	<i>trillado</i>	0,3438228	96,3333324454725
<i>testigo</i>	<i>ocular</i>	0,3409549	90,6363666121319
<i>mirada</i>	<i>penetrante</i>	0,1941656	87,214282613629
<i>precio</i>	<i>irrisorio</i>	0,094	86,999982304874
<i>necesidad</i>	<i>imperioso</i>	0,2272065	85,3333331557094
<i>voz</i>	<i>estridente</i>	0,2045455	84,7142806504476
<i>ojo</i>	<i>penetrante</i>	0,1874216	84,1428532289289
<i>edad</i>	<i>avanzado</i>	0,1275232	82,1249971597993
<i>edad</i>	<i>avanzada</i>	0,1275232	82,1249971597993
<i>movimiento</i>	<i>brusco</i>	0,1216772	79,6315753950049
<i>documento</i>	<i>fehaciente</i>	0,1500904	78,9999953648077
<i>pregunta</i>	<i>formular</i>	0,2516515	78,6999982346979
<i>discusión</i>	<i>acalorado</i>	0,1558185	76,9999971669168
<i>libertad</i>	<i>condicional</i>	0,1933497	76,0000066170476
<i>luz</i>	<i>cegador</i>	0,3292156	74,7142820524585
<i>luces</i>	<i>cegador</i>	0,3292156	74,7142820524585

Tabla 2-28 Combinaciones Sustantivo + Adjetivo con máximo valor de Box-plot

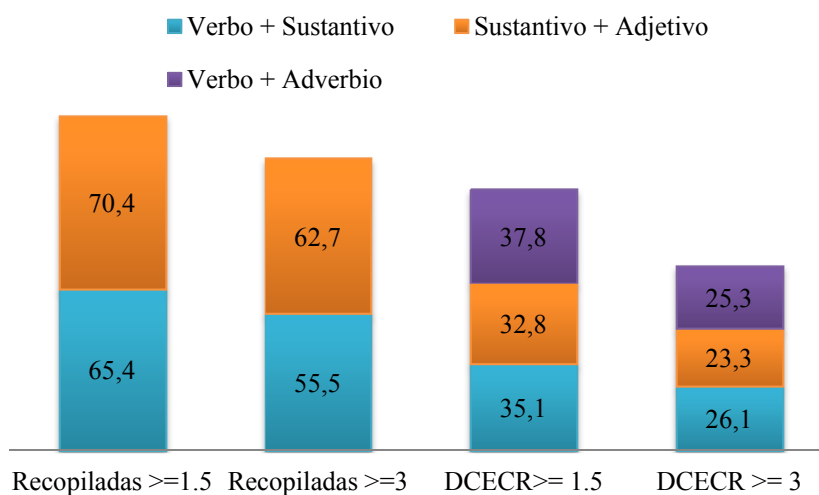


Gráfico 2-14 Cobertura en el caso de aplicar Box-Plot

<i>palabra</i>	$x_i$	<i>FrecRelativa(palabra, <math>x_i</math>)</i>	<i>Box-plot(palabra, <math>x_i</math>)</i>
<i>herméticamente</i>	<i>cerrar</i>	0,7858796	153,805255159383
<i>comercialmente</i>	<i>estrenar</i>	0,08759124	144,595113843274
<i>olímpicamente</i>	<i>ignorar</i>	0,3606557	139,279349376819
<i>alfabéticamente</i>	<i>ordenar</i>	0,508772	118,27913930758
<i>celosamente</i>	<i>guardar</i>	0,3695652	106,073907808064
<i>rematadamente</i>	<i>locar</i>	0,2102564	105,617925682556
<i>activamente</i>	<i>participar</i>	0,2290419	102,565217422258
<i>frontalmente</i>	<i>chocar</i>	0,2113821	100,519609091312
<i>herméticamente</i>	<i>sellar</i>	0,05555556	93,2356047672532
<i>gratamente</i>	<i>sorprender</i>	0,3082192	82,7287324380578
<i>encarecidamente</i>	<i>rogar</i>	0,2931243	81,9901037563911
<i>ardientemente</i>	<i>desear</i>	0,4692623	78,9741632248684
<i>descaradamente</i>	<i>mentir</i>	0,08625731	73,420234354581
<i>terminantemente</i>	<i>prohibir</i>	0,25459	65,3472818468441
<i>rotundamente</i>	<i>negar</i>	0,3542977	63,152643313143
<i>acaloradamente</i>	<i>discutir</i>	0,2864321	62,5256543244676
<i>fervientemente</i>	<i>desear</i>	0,3698225	61,957777978053
<i>atentamente</i>	<i>escuchar</i>	0,2001596	60,3643561695563
<i>copiosamente</i>	<i>sudar</i>	0,1366384	60,3583868324945
<i>cordialmente</i>	<i>detestar</i>	0,03237743	59,4035125720244

Tabla 2-29 Combinaciones sustantivo + verbo con máximo valor de Box-plot

### 2.6.4. Método de Hampel

El estadístico  $MAD$  (2), basado en la  $MEDA$  de los datos, es robusto frente a observaciones atípicas, por lo que métodos más efectivos para la identificación de outliers se basan en él, si bien los cálculos son más complejos (Leys y otros, 2013). La  $MEDA$  (3) es la mediana de las desviaciones a la mediana de los datos en valores absolutos (3) y  $b$  (4) en caso de distribuciones normales debe ser: 1,4826.

$$(2) \mathbf{MAD} = \mathbf{b} * \mathbf{MEDA}$$

$$(3) \mathbf{MEDA} = \mathbf{mediana}\{|x_1 - \mathbf{mediana}(X)|, |x_2 - \mathbf{mediana}(X)|, \dots, |x_n - \mathbf{mediana}(X)|\}$$

$$(4) \mathbf{b} = \frac{\mathbf{1}}{\mathbf{Q}_3}$$

$X = \{x_1, x_2, \dots, x_n\}$  es la muestra de la variable.

$MAD$  es un estadístico más adecuado para el problema por las siguientes propiedades:

- El punto de ruptura de  $MAD$  es de 0,5, superior al del rango intercuartílico que es de 0,25, o 0 en el caso de la media. Este valor indica la proporción máxima de observaciones que pueden ser infinito (estar contaminadas) sin que el estimador se resienta, es decir, que de un valor falso.
- $MAD$  no tiene restricciones respecto al tamaño muestral.
- Permite establecer un criterio objetivo fácilmente automatizable para discriminar entre las combinaciones que se consideran outliers y las que no. Por extensión proporciona un mecanismo para identificar automáticamente los candidatos a colocaciones.

El método de Hampel acepta como outlier los valores  $x_i$  que cumplen la condición  $|MAD(x_i)| \geq 4,5$  (5).

$$(5) \mathbf{MAD}(x_i) = \frac{|x_i - \mathbf{mediana}(X)|}{\mathbf{MEDA}}$$

La adaptación del método a este trabajo lleva a considerar:

$$\mathbf{MADF}(X_i, \mathbf{Palabra}) = \frac{|\mathbf{FrecRel}(X_i, \mathbf{Palabra}) - \mathbf{mediana}(\mathbf{Palabra})|}{\mathbf{MEDAF}(\mathbf{Palabra})}$$

Siendo:

$$\mathbf{mediana}(\mathbf{Palabra}) = \mathbf{mediana}\left\{|\mathbf{FrecRel}(x_j, \mathbf{Palabra}) - \mathbf{mediana}(\mathbf{Palabra})|_{j=1 \dots n}\right\}$$

y

$$MEDAF(Palabra) = mediana \left\{ |FrecRel(x_j, Palabra) - mediana(Palabra)|_{j=1 \dots n} \right\}$$

Las Tabla 2-31 y Tabla 2-32 muestran los rangos obtenidos en ambos corpus cuando se intentan detectar los outliers a partir de las desviaciones respecto a la mediana de las frecuencias relativas. Hay que tener en cuenta que se pueden realizar los cálculos en los casos en que la *MEDA* sea distinta de cero y que es necesario calcular los valores del indicador para cada uno de los colocados, ya que si el uso respecto a la base no es preferente si lo puede ser respecto al colocativo y viceversa. En este caso, la cobertura aumenta significativamente, si bien en los ejemplos extraídos en torno al punto de corte teórico se aprecian falsas colocaciones (Tabla 2-30).

	Recopiladas		DCECR.		Corpus	
<b>MAD<math>\geq</math>4,5</b>	<i>hacer</i>	<i>alarde</i>	<i>dar</i>	<i>vuelta</i>	<i>encoger</i>	<i>hombro</i>
<b>Max</b>	<i>echar</i>	<i>vistazo</i>	<i>echar</i>	<i>vistazo</i>	<i>dar</i>	<i>vuelta</i>
	<i>dar</i>	<i>palmada</i>	<i>dar</i>	<i>gracia</i>	<i>don</i>	<i>doña</i>
	<i>hacer</i>	<i>pregunta</i>	<i>dar</i>	<i>paso</i>	<i>abrir</i>	<i>puerta</i>
	<i>ladrar</i>	<i>perro</i>	<i>soplar</i>	<i>viento</i>	<i>decir</i>	<i>señor</i>
	<i>latir</i>	<i>corazón</i>	<i>dar</i>	<i>prisa</i>	<i>don</i>	<i>quijo</i>
	<i>poder</i>	<i>ejecutivo</i>	<i>dar</i>	<i>salto</i>	<i>menear</i>	<i>cabeza</i>
<b>MAD<math>\geq</math>4,5</b>	<i>andar</i>	<i>barco</i>	<i>dorar</i>	<i>fama</i>	<i>juir</i>	<i>cerrado</i>
<b>Min</b>	<i>incurrir</i>	<i>ira</i>	<i>arranque</i>	<i>patriotismo</i>	<i>servido</i>	<i>regentar</i>
	<i>cifra</i>	<i>grande</i>	<i>fórmula</i>	<i>predicción</i>	<i>botón</i>	<i>suicidio</i>
	<i>memoria</i>	<i>frágil</i>	<i>sofocar</i>	<i>levantamiento</i>	<i>adquirir</i>	<i>círculo</i>
	<i>fuelle</i>	<i>fidedigno</i>	<i>influenciar</i>	<i>dominante</i>	<i>sensación</i>	<i>huir</i>
	<i>acarrear</i>	<i>molestia</i>	<i>distorsionar</i>	<i>versión</i>	<i>maldito</i>	<i>amenazar</i>
	<i>abrir</i>	<i>crisis</i>	<i>estrechar</i>	<i>acuerdo</i>	<i>pluma</i>	<i>relación</i>

Tabla 2-30 Ejemplos de combinaciones que superan el criterio MAD $\geq$ 4,5

	Rango	Combinaciones con <i>MAD</i> $\geq$ 4,5
<b>Recopiladas</b>	[0 2087,0032]	1532
<b>DCECR</b>	[0 4304,111]	15039
<b>Corpus</b>	[0 4309,9998]	4318191

Tabla 2-31 Resultados obtenidos en el corpus para MAD

	Rango	Combinaciones con <i>MAD</i> ≥ 4,5
<b>Recopiladas</b>	[0 68,00004]	933
<b>DCECR</b>	[0 344,9999]	2735
<b>Corpus</b>	[0 344,9999]	455738

Tabla 2-32 Resultados obtenidos en el corpus de Galdós para MAD

Los ejemplos extraídos reflejan cierta tendencia a que los verbos ocupen las posiciones más altas del ranking, y entre éstos, los funcionales como *dar*, *echar* o *hacer*.

Los rangos alcanzados en cada uno de los grupos, según la estructura de la combinación, se muestran en la Tabla 2-33. Se aprecian intervalos similares en los grupos *sustantivo + adjetivo* y *verbo + dverbio*, y destaca la amplitud mucho mayor en el caso *verbo + sustantivo*. Respecto a la cobertura en los conjuntos de ensayo, en todas las estructuras sube drásticamente respecto a los resultados obtenidos con el criterio de *ZChebyshev*.

	<i>MAD</i>
<b>Verbo + Sustantivo</b>	[0 4309,9999]
<b>Sustantivo + Adjetivo</b>	[0 2531,476]
<b>Verbo + Adverbio</b>	[0 2332,5601]

Tabla 2-33 Rango de Mad según la estructura de la combinación.

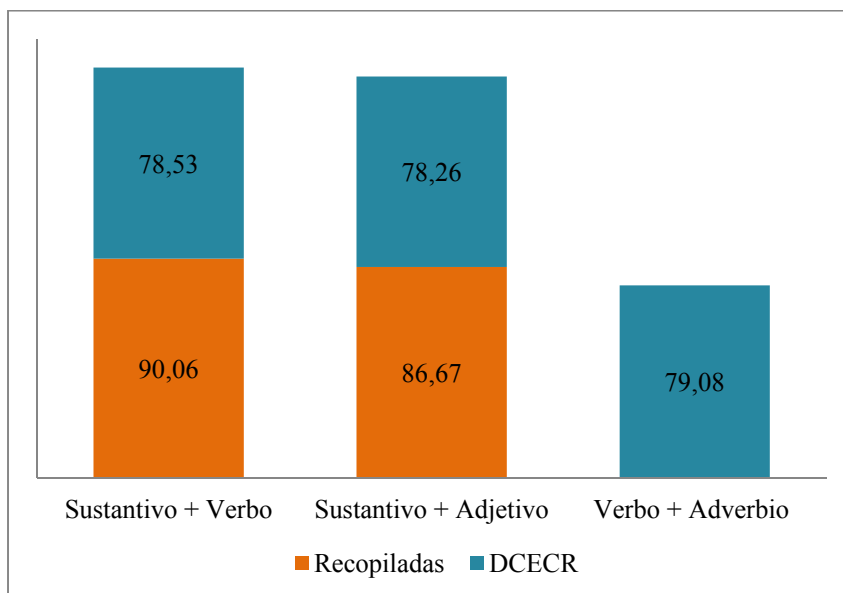


Gráfico 2-15 Combinaciones de los conjuntos de ensayo que superan el criterio de MAD ≥ 4,5

Para finalizar, se contrasta la cobertura de los distintos criterios barajados sobre los conjuntos de ensayo, teniendo en cuenta también la estructura de la colocación. Los resultados se resumen en



el Gráfico 2-16; en él se aprecia la gran diferencia de cobertura entre *MAD* y el resto. Se observa que en orden de más restrictivos a menos se tiene: frecuencia relativa ( $\geq 0,05$ ), información mutua ( $\geq 6$ ), *Zchebyshev* ( $\geq 3$ ), *Box-Plot* y *MAD* ( $\geq 4,5$ ) respectivamente. Sin embargo, ya se ha advertido que la frontera teórica marcada por *MAD* introduce ruido en los resultados.

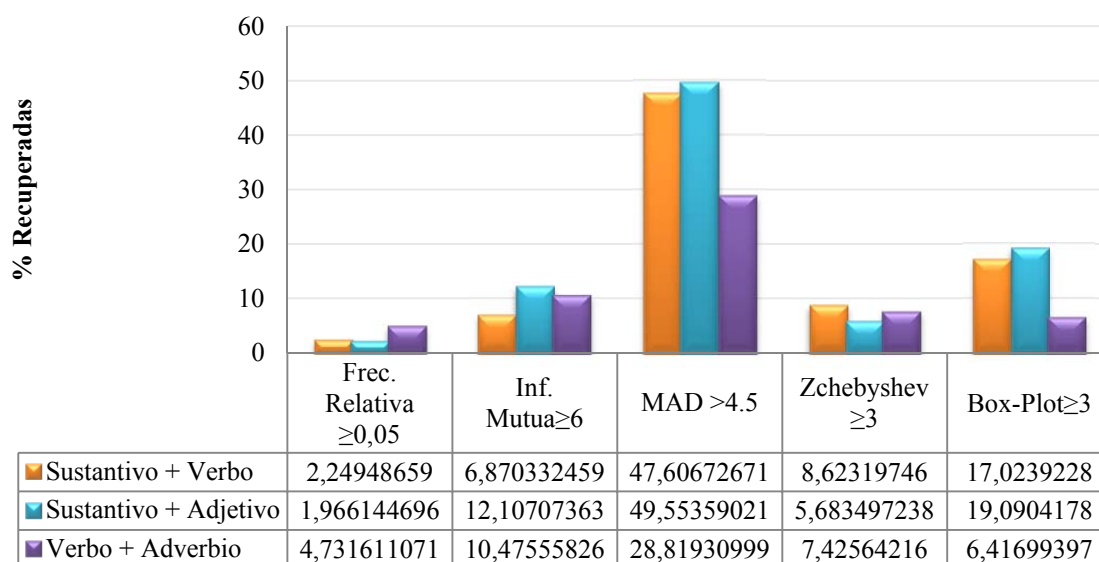


Gráfico 2-16 Comparativa de los criterios de extracción automática

## 2.7. Estrategia para la extracción de información colocacional generada por un corpus extenso

Fruto del análisis de los resultados registrados en la BD de colocaciones surgen las líneas que marcan su explotación en el desarrollo de un módulo de software que dé soporte a la presentación de las características combinatorias de las palabras en español: sugieren consideraciones a tener en cuenta a la hora de explotar corpus textuales extensos cuando se utilizan indicadores basados en las frecuencias de uso para medir las relaciones de asociación entre palabras. Los resultados de los epígrafes previos llevan a la necesidad de considerar la composición del corpus como una muestra en la que el uso que se da en él a las palabras influye en la posición en cualquier ranking que se intente establecer, y no solo el fenómeno fraseológico. Una muestra muy extensa del lenguaje agranda las diferencias de comportamiento según factores tales como el tipo de elementos que conformen la combinación, la estructura de la colocación o si intervienen palabras funcionales, de gran frecuencia de uso. Por otra parte, el límite entre lo que es colocación cuando los indicadores se basan en la independencia se debe variar en función del tamaño del corpus; sin embargo, los indicadores que pretenden captar el uso preferente permiten establecer un criterio objetivo de forma independiente a este factor. Sobre la base de estas consideraciones se propone seguir una estrategia en la que dada una palabra se muestren sus capacidades combinatorias, de forma que se presente al usuario un

resumen de los datos de frecuencias y se filtran aquellas que se comportan como colocaciones en su rango de variación de frecuencias relativas; es decir, considerando la muestra de su comportamiento en el corpus y no todo el corpus como una muestra de su comportamiento. Siguiendo esta estrategia, se consulta sobre una determinada palabra y se representa el ranking según el o los criterios seleccionados para cada una de las estructuras colocacionales en las que podría formar parte. En la Figura 2-10 se expone un esquema de este proceso que se describe a continuación:

- i. Lematizar la palabra para obtener las formas canónicas que la representan en la B. DD.
- ii. Dada una de las formas canónicas, se extraen de la BDD las listas ordenadas de combinaciones según alguno de las medidas de asociación: *frecuencia relativa*, *información mutua*, *z-score*, *ZChebyshev*, *Box-Plot* o *MAD*. Este proceso se lleva a cabo para cada una de las posibles estructuras: *verbo + sustantivo*, *sustantivo + adjetivo* o *verbo + adverbio* en las que pueda intervenir la forma canónica.

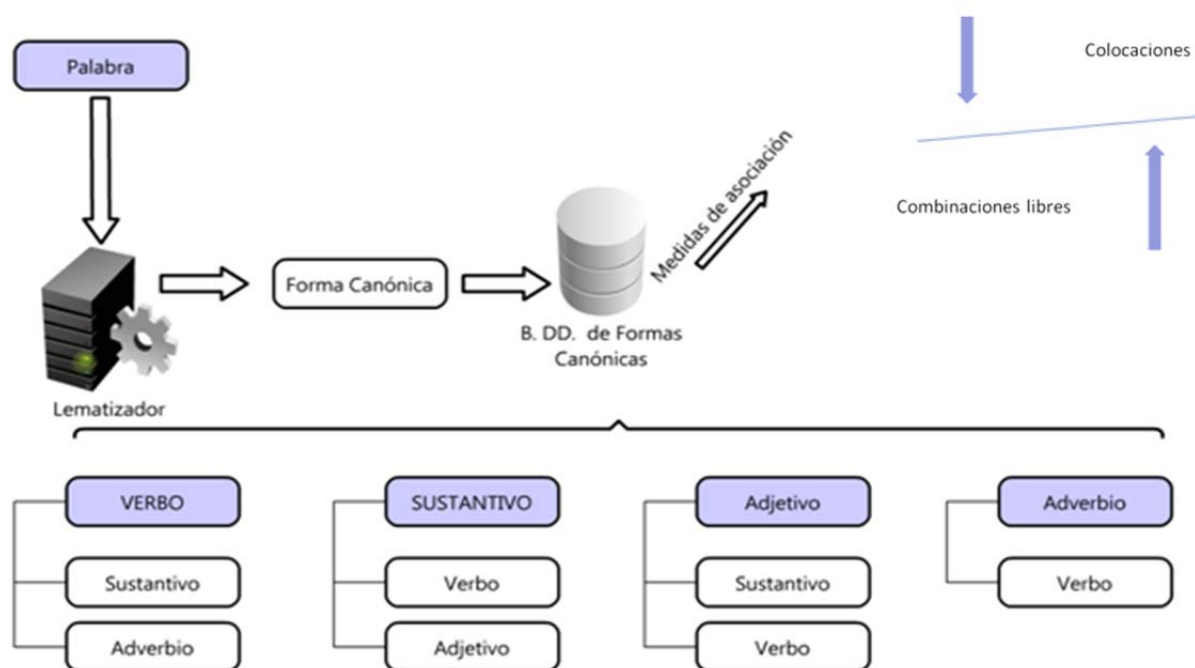
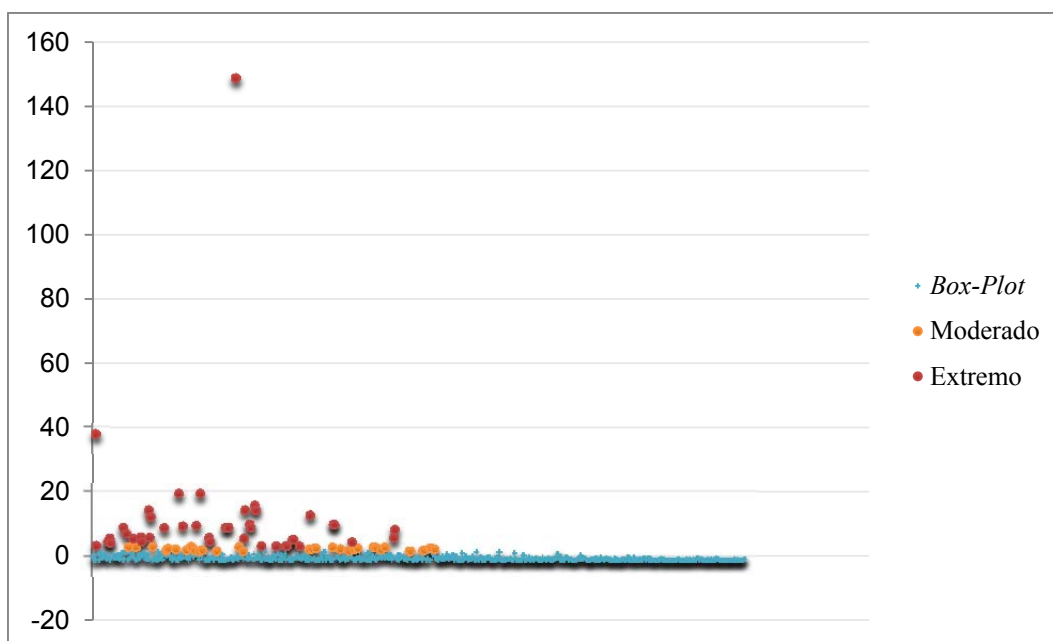


Figura 2-10 Extracción de colocaciones en corpus textuales extenso.

Un ejemplo de uso de esta metodología cuando la palabra de referencia es el verbo *palpitar* y se aplica el indicador *Box-Plot* evaluado sobre los sustantivos con los que aparece en el corpus, arroja los resultados mostrados en el Gráfico 2-17. Las 20 combinaciones que se obtienen los valores máximos de *Box-Plot* se listan en la Tabla 2-34.

Gráfico 2-17 *Box-Plot* para las combinaciones palpitar + Sustantivo en el corpus

	$x_i$	$FrecRelativa(x_i, palpitar)$	$Box-Plot(x_i, palpitar)$
<b>palpitar</b>	<i>corazón</i>	0,330837	148,899993456854
	<i>pecho</i>	0,08634361	37,8999958880714
	<i>vida</i>	0,04537445	19,2999986998272
	<i>vido</i>	0,04537445	19,2999986998272
	<i>ojo</i>	0,03744493	15,6999990803656
	<i>baja</i>	0,03436123	14,2999983827119
	<i>bajo</i>	0,03436123	14,2999983827119
	<i>amor</i>	0,03414097	14,1999996617437
	<i>fuerza</i>	0,03083701	12,6999993974809
	<i>sangre</i>	0,02951542	12,099999460904
	<i>mano</i>	0,0246696	9,89999969345519
	<i>seno</i>	0,0246696	9,89999969345519
	<i>manos</i>	0,02444934	9,79999928120528
	<i>cabeza</i>	0,02378855	9,49999973573724
	<i>vez</i>	0,02334802	9,29999975687826
	<i>lucos</i>	0,02246696	8,89999895351946
<i>luz</i>	0,02246696	8,89999895351946	

Tabla 2-34 Valores máximos de *Box-Plot* para el verbo palpitar + Sustantivo

### 2.7.1. Resultados obtenidos sobre el corpus completo

En las tablas Tabla 2-35 a la Tabla 2-58 se muestran ejemplos de los listados obtenidos siguiendo esta estrategia cuando se consultan formas canónicas (FC) que aparecen en distintos rangos de frecuencia en el corpus. Se han elegido los verbos *dar* y *propinar* (Tabla 2-35, a Tabla 2-40), los sustantivos *tierra* y *contrato* (Tabla 2-41 a Tabla 2-46), los adjetivos *público* y *omiso* (Tabla 2-47 a Tabla 2-52), y los adverbios *libremente* y *frugalmente* (Tabla 2-53). En todos los casos se utiliza *ZChebyshev*, *Mad* y *Box-Plot*. Se muestran los 10 mejores, los 10 más cercanos al umbral utilizado para discriminar y 10 de los que obtienen valor mínimo, siempre que sea posible. Hay que señalar que para el adverbio *frugalmente* se muestran todos los casos de combinaciones que aparecen en el corpus.

Max <i>ZChebyshev</i>		Punto de corte <i>ZChebyshev</i>		Min <i>ZChebyshev</i>	
<i>vuelto</i>	78,60	<i>bodega</i>	3,00	<i>pipado</i>	0
<i>vuelta</i>	74,39	<i>alusión</i>	3,00	<i>pretura</i>	0
<i>golpe</i>	52,28	<i>pe</i>	3,00	<i>elogiado</i>	0
<i>pasa</i>	51,76	<i>tarifa</i>	3,00	<i>pecunia</i>	0
<i>paso</i>	51,76	<i>ubicación</i>	3,00	<i>sene</i>	0
<i>gracia</i>	48,06	<i>desconcertado</i>	3,00	<i>anclote</i>	0
<i>muestra</i>	46,03	<i>marchado</i>	3,00	<i>molinosismo</i>	0
<i>salto</i>	45,83	<i>brote</i>	3,00	<i>canguelo</i>	0
<i>orden</i>	43,30	<i>inoportuno</i>	3,00	<i>remado</i>	0
<i>órdenes</i>	43,30	<i>evento</i>	3,00	<i>robla</i>	0

Tabla 2-35 Resultados *ZChebyshev*, verbo dar + Sustantivo

Max MAD		Punto de corte Mad		Min Mad	
<i>vuelta</i>	4304,11	<i>cáliz</i>	4,51	<i>jengibre</i>	0
<i>vuelto</i>	3969,09	<i>presunción</i>	4,51	<i>yunque</i>	0
<i>palmado</i>	1963,00	<i>aprobado</i>	4,51	<i>palabrita</i>	0
<i>palmada</i>	1962,49	<i>editorial</i>	4,51	<i>naire</i>	0
<i>gracia</i>	1674,36	<i>lentitud</i>	4,51	<i>mosquetería</i>	0
<i>pasa</i>	1471,76	<i>manual</i>	4,51	<i>molusco</i>	0
<i>paso</i>	1471,76	<i>montañoso</i>	4,50	<i>manillar</i>	0
<i>priso</i>	1382,75	<i>suplicación</i>	4,50	<i>judicatura</i>	0
<i>prisa</i>	1378,25	<i>candelaria</i>	4,50	<i>landa</i>	0
<i>golpe</i>	1329,11	<i>tensó</i>	4,50	<i>latamente</i>	0

Tabla 2-36 Resultados MAD, verbo dar + Sustantivo

<b>Max Box-Plot</b>		<b>Punto de corte Box-Plot</b>		<b>Min Box-Plot</b>	
<i>vuelto</i>	265,42	<i>intelectual</i>	3,03	<i>sociable</i>	-1,03
<i>vuelta</i>	259,01	<i>docena</i>	3,03	<i>socioeconómico</i>	-1,03
<i>hombre</i>	207,23	<i>comentarios</i>	3,03	<i>subjefe</i>	-1,03
<i>señor</i>	179,59	<i>comentario</i>	3,03	<i>substituto</i>	-1,03
<i>señora</i>	179,59	<i>acero</i>	3,02	<i>tecnicismo</i>	-1,03
<i>diosa</i>	178,23	<i>altar</i>	3,01	<i>teleología</i>	-1,03
<i>dios</i>	178,23	<i>cuartel</i>	3,01	<i>teologal</i>	-1,03
<i>mujer</i>	173,11	<i>tradición</i>	3,01	<i>sufijo</i>	-1,03
<i>bienes</i>	169,78	<i>traducción</i>	3,00	<i>susurrante</i>	-1,03
<i>bien</i>	169,78	<i>contacto</i>	3,00	<i>tetilla</i>	-1,03

Tabla 2-37 Resultados *Box-Plot*, verbo dar + Sustantivo

<b>Min ZChebyshev</b>		<b>Min ZChebyshev</b>	
<i>golpe</i>	13,69	<i>fin</i>	0
<i>paliza</i>	5,588	<i>muchacha</i>	0
<i>bueno</i>	5,075	<i>muchacho</i>	0
<i>patada</i>	4,936	<i>sonoro</i>	0
<i>buena</i>	4,424	<i>trasero</i>	0
<i>puntapié</i>	4,377	<i>joven</i>	0
<i>puñetazo</i>	4,191	<i>mordisco</i>	0
<i>puñete</i>	4,191	<i>sonora</i>	0
<i>fuerte</i>	3,632	<i>bofetón</i>	0
<i>vez</i>	3,25	<i>trasera</i>	0

Tabla 2-38 Resultados ZChebyshev, verbo propinar + Sustantivo

<b>Max MAD</b>		<b>Punto de corte MAD</b>		<b>Min MAD</b>	
<i>golpe</i>	301,00	<i>amor</i>	5,00	<i>vero</i>	0
<i>paliza</i>	127,00	<i>momento</i>	5,00	<i>vera</i>	0
<i>paliza</i>	126,99	<i>mandíbula</i>	5,00	<i>obra</i>	0
<i>bueno</i>	116,00	<i>hermana</i>	5,00	<i>problemas</i>	0
<i>patada</i>	113,00	<i>hermano</i>	5,00	<i>descubierta</i>	0
<i>buena</i>	102,00	<i>gracia</i>	5,00	<i>actitud</i>	0
<i>puntapié</i>	101,00	<i>oro</i>	5,00	<i>soldado</i>	0
<i>puñete</i>	97,000	<i>bebida</i>	5,00	<i>complacencia</i>	0
<i>puñetazo</i>	97,000	<i>sonrisa</i>	5,00	<i>complacencia</i>	0
<i>fuerte</i>	85,000	<i>caso</i>	5,00	<i>problema</i>	0

Tabla 2-39 Resultados MAD, verbo propinar +Sustantivo

<b>Max Box-Plot</b>		<b>Punto de corte Box-Plot</b>		<b>Min Box-Plot</b>	
<i>golpe</i>	42,14	<i>dólar</i>	6,14	<i>lejos</i>	-1
<i>paliza</i>	17,29	<i>generoso</i>	4,86	<i>manera</i>	-1
<i>bueno</i>	15,71	<i>padre</i>	4,57	<i>manero</i>	-1
<i>patada</i>	15,29	<i>codo</i>	3,71	<i>nalga</i>	-1
<i>buena</i>	13,71	<i>madre</i>	3,57	<i>esclava</i>	-1
<i>puntapié</i>	13,57	<i>tiempo</i>	3,43	<i>esclavo</i>	-1
<i>puñetazo</i>	13,00	<i>señor</i>	3,29	<i>estilo</i>	-1
<i>puñete</i>	13,00	<i>señora</i>	3,29	<i>vicioso</i>	-1
<i>fuerte</i>	11,29	<i>ante</i>	3,00	<i>felpa</i>	-1
<i>vez</i>	10,14	<i>antes</i>	3,00	<i>gratificado</i>	-1

Tabla 2-40 Resultados Box-Plot, verbo propinar

Se observa que los valores máximos de *ZChebyshev* y *MAD* arrojan buenos resultados en las regiones de máximos, produciendo además listados similares. Los valores con puntuaciones más altas de *Box-Plot* son peores para el caso del verbo *dar*, en el que las combinaciones que se detectan como anomalías corresponden a palabras de uso frecuente, más que de uso preferente con este verbo. Asimismo, los valores mínimos de los tres se producen entre combinaciones libres. La diferencia principal surge en las inmediaciones del umbral que

marca la zona de lo que se propondría como colocación. Tanto en *MAD* como *Box-Plot* corresponde a una región difusa donde se mezclan combinaciones libres: *propinar hermana* con otras con cierta fijación como *propinar amor* o *propinar sonrisa*. *ZChebyshev* establece un punto de corte más drástico, lo que se traduce en que la zona de aceptación, si se ciñe a palabras de uso más restringido, prácticamente coincide con los máximos seleccionados. Un comportamiento similar está presente en los restantes ejemplos barajados (Tablas Tabla 2-41 a Tabla 2-49).

<b>Max ZChebyshev</b>		<b>Punto de corte ZChebyshev</b>		<b>Min ZChebyshev</b>	
<i>firme</i>	33,83	<i>corona</i>	3,00	<i>nulidad</i>	0
<i>cielo</i>	31,00	<i>monte</i>	3,00	<i>mutación</i>	0
<i>firmar</i>	30,71	<i>venida</i>	3,00	<i>perturbación</i>	0
<i>adentrar</i>	30,60	<i>bordo</i>	3,00	<i>domeñar</i>	0
<i>hacer</i>	30,08	<i>demasiar</i>	3,00	<i>lotófago</i>	0
<i>mar</i>	29,74	<i>contornos</i>	3,00	<i>amainar</i>	0
<i>cultivar</i>	28,06	<i>contorno</i>	3,00	<i>ufanía</i>	0
<i>tener</i>	28,06	<i>muerte</i>	3,00	<i>impeditivo</i>	0
<i>poder</i>	27,49	<i>amarrar</i>	3,00	<i>desjugar</i>	0
<i>decir</i>	24,82	<i>carga</i>	3,00	<i>remado</i>	0

Tabla 2-41 Resultados ZChebyshev, sustantivo tierra + Adjetivo

<b>Max MAD</b>		<b>Punto de corte MAD</b>		<b>Min MAD</b>	
<i>hacer</i>	826,72	<i>fluvial</i>	4,66	<i>carrete</i>	0
<i>cielo</i>	784,81	<i>sinuoso</i>	4,66	<i>carbonero</i>	0
<i>adentrar</i>	780,66	<i>enemistad</i>	4,66	<i>caminata</i>	0
<i>tener</i>	771,54	<i>carbono</i>	4,66	<i>chachapoyo</i>	0
<i>poder</i>	756,09	<i>confinado</i>	4,62	<i>consciencia</i>	0
<i>decir</i>	683,36	<i>usura</i>	4,59	<i>cárdeno</i>	0
<i>firme</i>	682,84	<i>acorde</i>	4,59	<i>cincelado</i>	0
<i>mar</i>	680,78	<i>recitar</i>	4,54	<i>cigarrillo</i>	0
<i>ver</i>	654,00	<i>tenso</i>	4,50	<i>deliberación</i>	0
<i>cielo</i>	570,50	<i>medieval</i>	4,5	<i>desacato</i>	0

Tabla 2-42 Resultados MAD, sustantivo tierra + Adjetivo

<b>Max Box-Plot</b>		<b>Punto de corte Box-Plot</b>		<b>Min Box-Plot</b>	
<i>hacer</i>	182,07	<i>tronco</i>	3,05	<i>desasir</i>	-0,03
<i>tener</i>	169,93	<i>social</i>	3,04	<i>yate</i>	-0,03
<i>poder</i>	166,53	<i>recibido</i>	3,04	<i>languidecer</i>	-0,03
<i>decir</i>	150,53	<i>morado</i>	3,04	<i>empedernir</i>	-0,03
<i>ver</i>	144,07	<i>morada</i>	3,04	<i>amordazar</i>	-0,03
<i>cielo</i>	125,70	<i>rumbo</i>	3,03	<i>erar</i>	-0,03
<i>dar</i>	123,11	<i>cruel</i>	3,02	<i>reportar</i>	-0,03
<i>hombre</i>	109,40	<i>cortés</i>	3,02	<i>violar</i>	-0,03
<i>ir</i>	101,36	<i>dilatar</i>	3,02	<i>tortuoso</i>	-0,03
<i>mar</i>	95,42	<i>acostumbrar</i>	3,02	<i>traicionero</i>	-0,03

Tabla 2-43 Resultados *Box-Plot*, sustantivo tierra + Adjetivo

<b>Max ZChebyshev</b>		<b>Punto de corte ZChebyshev</b>		<b>Min ZChebyshev</b>	
<i>poder</i>	15,93	<i>firmado</i>	3,54	<i>fundamento</i>	0
<i>hacer</i>	15,79	<i>ver</i>	3,44	<i>idea</i>	0
<i>firmar</i>	14,88	<i>contrato</i>	3,41	<i>ideo</i>	0
<i>tener</i>	13,55	<i>contrato</i>	3,40	<i>número</i>	0
<i>celebrar</i>	11,16	<i>hecho</i>	3,38	<i>recordar</i>	0
<i>decir</i>	9,83	<i>cancelar</i>	3,25	<i>reinar</i>	0
<i>matrimonial</i>	7,97	<i>gobierno</i>	3,23	<i>solar</i>	0
<i>firmar</i>	7,10	<i>haber</i>	3,21	<i>ir</i>	0
<i>social</i>	6,48	<i>contrato</i>	3,16	<i>mercantil</i>	0
<i>dar</i>	6,27	<i>nuevo</i>	3,06	<i>contrato</i>	0

Tabla 2-44 Resultados ZChebyshev, sustantivo contrato + Adjetivo



<b>Max MAD</b>		<b>Punto de corte MAD</b>		<b>Min MAD</b>	
<i>poder</i>	192,33	<i>fecho</i>	3,22	<i>condicional</i>	0
<i>hacer</i>	190,66	<i>contrato</i>	4,62	<i>bastir</i>	0
<i>firmar</i>	179,83	<i>promesa</i>	4,62	<i>destruir</i>	0
<i>matrimonial</i>	174,00	<i>millo</i>	4,59	<i>demasiar</i>	0
<i>tener</i>	164,16	<i>millón</i>	4,57	<i>encargar</i>	0
<i>celebrar</i>	135,83	<i>nombre</i>	4,57	<i>dotar</i>	0
<i>firmar</i>	134,62	<i>pueblo</i>	4,52	<i>extinguir</i>	0
<i>decir</i>	119,99	<i>mercantil</i>	4,52	<i>manifestar</i>	0
<i>celebrar</i>	116,28	<i>arrendatario</i>	4,50	<i>salariar</i>	0
<i>rescindir</i>	86,999	<i>usar</i>	4,50	<i>cincuenta</i>	0

Tabla 2-45 Resultados MAD, sustantivo contrato + Adjetivo

<b>Max Box-Plot</b>		<b>Punto de corte Box-Plot</b>		<b>Min Box-Plot</b>	
<i>poder</i>	49,39	<i>vido</i>	3,17	<i>auxiliador</i>	-1,04
<i>hacer</i>	48,96	<i>parecer</i>	3,17	<i>contracto</i>	-1,04
<i>firmar</i>	46,13	<i>vida</i>	3,17	<i>contracto</i>	-1,04
<i>tener</i>	42,04	<i>padre</i>	3,17	<i>comodato</i>	-1,04
<i>celebrar</i>	34,65	<i>conocer</i>	3,09	<i>ayuno</i>	-1,04
<i>decir</i>	30,52	<i>pago</i>	3,09	<i>acordado</i>	-1,04
<i>social</i>	20,17	<i>volver</i>	3,09	<i>contacto</i>	-1,04
<i>dar</i>	19,52	<i>pueblo</i>	3,04	<i>abuelo</i>	-1,04
<i>partir</i>	17,87	<i>conseguir</i>	3,04	<i>ayunas</i>	-1,04
<i>deber</i>	17,78	<i>nombre</i>	3,04	<i>condesa</i>	-1,04

Tabla 2-46 Resultados Box-Plot, sustantivo contrato + Adjetivo

<b>Max ZChebyshev</b>		<b>Punto de corte ZChebyshev</b>		<b>Min ZChebyshev</b>	
<i>hacer</i>	29,05	<i>judicial</i>	3,03	<i>animal</i>	0
<i>poder</i>	27,70	<i>económico</i>	3,03	<i>terma</i>	0
<i>privado</i>	25,57	<i>acogida</i>	3,03	<i>distraído</i>	0
<i>funcionario</i>	24,70	<i>gasto</i>	3,02	<i>intrigante</i>	0
<i>privar</i>	24,27	<i>caso</i>	3,02	<i>imprescriptible</i>	0
<i>instrucciones</i>	23,51	<i>universidad</i>	3,02	<i>precaución</i>	0
<i>instrucción</i>	23,51	<i>dirigir</i>	3,01	<i>latas</i>	0
<i>ministerio</i>	20,97	<i>ilustrar</i>	3,01	<i>gana</i>	0
<i>administración</i>	20,54	<i>tener</i>	3,01	<i>arriesgar</i>	0
<i>tener</i>	20,47	<i>militar</i>	3,01	<i>trimestralmente</i>	0

Tabla 2-47 Resultados ZChebyshev, adjetivo público

<b>Max MAD</b>		<b>Punto de corte, MAD</b>		<b>Min MAD</b>	
<i>hacer</i>	636,21	<i>mentalidad</i>	4,50	<i>suficiencia</i>	0
<i>poder</i>	606,85	<i>adaptar</i>	4,50	<i>inválido</i>	0
<i>subastar</i>	591,00	<i>estampa</i>	4,50	<i>público</i>	0
<i>privar</i>	481,14	<i>suelo</i>	4,50	<i>quieto</i>	0
<i>privado</i>	469,14	<i>ventana</i>	4,50	<i>público</i>	0
<i>tener</i>	449,99	<i>transformación</i>	4,50	<i>conminar</i>	0
<i>privada</i>	427,00	<i>gubernamental</i>	4,5	<i>etiquetar</i>	0
<i>instrucción</i>	416,75	<i>bofetada</i>	4,5	<i>mortificar</i>	0
<i>instrucciones</i>	416,75	<i>sobrevenir</i>	4,5	<i>semillar</i>	0
<i>funcionario</i>	370,00	<i>agencia</i>	4,5	<i>público</i>	0

Tabla 2-48 Resultados MAD, adjetivo público + Adjetivo

<b>Max Box-Plot</b>		<b>Punto de corte, Box-Plot</b>		<b>Min Box-Plot</b>	
<i>bien</i>	54,41	<i>género</i>	3,09	<i>prosa</i>	-0,93
<i>bienes</i>	54,41	<i>tarde</i>	3,07	<i>rienda</i>	-0,93
<i>leyes</i>	51,77	<i>consideración</i>	3,07	<i>romería</i>	-0,93
<i>ley</i>	51,77	<i>decreto</i>	3,06	<i>acabo</i>	-0,93
<i>obra</i>	51,03	<i>garantía</i>	3,04	<i>basura</i>	-0,93
<i>servicio</i>	48,64	<i>voluntad</i>	3,03	<i>basuras</i>	-0,93
<i>instrucción</i>	47,41	<i>gestión</i>	3,03	<i>calificación</i>	-0,93
<i>instrucciones</i>	47,41	<i>comunidad</i>	3,03	<i>cuestor</i>	-0,93
<i>opinión</i>	46,43	<i>presencia</i>	3,01	<i>aristocracia</i>	-0,93
<i>orden</i>	46,13	<i>sentimiento</i>	3,01	<i>cuánto</i>	-0,93

Tabla 2-49 Resultados *Box-Plot*, adjetivo público + Adjetivo

Se destaca el comportamiento de formas canónicas fácilmente decodificables por tener una única acepción como el adjetivo *omiso*: en la zona de transición cuando se utiliza *MAD*, las combinaciones que aparecen tienen una relación estrecha que induce a catalogarlas como colocaciones: *remitente omiso*, *insulto omiso*, *parte omiso*, *protesta omiso*. Se observan también los efectos excesivamente restrictivos de *ZChebyshev* en palabras con estas características. En este caso, solo las formas canónicas que provienen de *hacer caso omiso*, superan la restricción impuesta por este indicador (Tabla 2-50). También destaca el buen comportamiento de *Box-Plot* que parece mejorar los resultados, introduciendo entre los casos anómalos palabras con menor frecuencia de uso.

<b>Max ZChebyshev</b>		
<i>omiso</i>	<i>hacer</i>	10,04
<i>omiso</i>	<i>casar</i>	10,01
<i>omiso</i>	<i>caso</i>	9,889
<i>omiso</i>	<i>hacendar</i>	2,980
<i>omiso</i>	<i>omiso</i>	2,814
<i>omiso</i>	<i>hacendar</i>	1,871
<i>omiso</i>	<i>casar</i>	1,855
<i>omiso</i>	<i>hacer</i>	0,563
<i>omiso</i>	<i>vanguardia</i>	0,492
<i>omiso</i>	<i>convención</i>	0,479

Tabla 2-50 Resultados *ZChebyshev*, adjetivo *omiso* + Adjetivo

<b>Max MAD</b>		<b>Punto de corte, MAD</b>		<b>Min MAD</b>	
<i>hacer</i>	1495,00	<i>remitente</i>	6,00	<i>culpar</i>	0
<i>casar</i>	1491,00	<i>señalar</i>	6,00	<i>insultar</i>	0
<i>caso</i>	1472,00	<i>soler</i>	6,00	<i>observar</i>	0
<i>hacendar</i>	457,000	<i>insulto</i>	6,00	<i>padre</i>	0
<i>caso</i>	91,0625	<i>llamar</i>	6,00	<i>procurar</i>	0
<i>poder</i>	76,0000	<i>sentir</i>	6,00	<i>recomendación</i>	0
<i>casar</i>	73,6000	<i>parte</i>	5,00	<i>regla</i>	0
<i>decir</i>	69,0000	<i>partes</i>	5,00	<i>sugerencia</i>	0
<i>tener</i>	58,0000	<i>ser</i>	5,00	<i>teléfono</i>	0
<i>hecho</i>	51,0000	<i>protesta</i>	5,00	<i>texto</i>	0

Tabla 2-51 Resultados MAD, adjetivo omiso + Adjetivo

<b>Max Box-Plot</b>		<b>Min Box-Plot</b>	
<i>caso</i>	367,25	<i>órdenes</i>	2,25
<i>hombre</i>	9,50	<i>mano</i>	2,25
<i>pregunta</i>	8,50	<i>manos</i>	2,00
<i>palabra</i>	7,00	<i>advertencia</i>	2,00
<i>mujer</i>	6,25	<i>presencia</i>	1,75
<i>comentario</i>	5,25	<i>interrupción</i>	1,50
<i>comentarios</i>	5,25	<i>bien</i>	1,50
<i>dolor</i>	4,50	<i>bienes</i>	1,50
<i>dolores</i>	4,50	<i>voz</i>	1,00
<i>ley</i>	3,50	<i>cámara</i>	0,75
<i>leyes</i>	3,5	<i>insulto</i>	0,75
<i>protesta</i>	3,5		

Tabla 2-52 Resultados Box-Plot, omiso + Sustantivo

Nótese como el comportamiento se reproduce también en los adverbios seleccionados para ilustrar los resultados, destacando el efecto restrictivo de *ZChebyshev* que excluye todas las combinaciones halladas del conjunto de colocaciones.

<b>Max ZChebyshev</b>		<b>Punto de corte, ZChebyshev</b>		<b>Min ZChebyshev</b>	
<i>hacer</i>	29,05	<i>relacionado</i>	4,24	<i>cierto</i>	0,01
<i>poder</i>	27,70	<i>vincular</i>	4,07	<i>entrelazar</i>	0,01
<i>privado</i>	25,57	<i>poder</i>	4,03	<i>limitado</i>	0,01
<i>funcionario</i>	24,70	<i>relacionado</i>	3,90	<i>social</i>	0,01
<i>privar</i>	24,27	<i>emparentado</i>	3,69	<i>llamar</i>	0,01
<i>instrucciones</i>	23,51	<i>abrazado</i>	3,53	<i>pasar</i>	0,01
<i>instrucción</i>	23,51	<i>vinculado</i>	3,30	<i>permanecer</i>	0,01
<i>ministerio</i>	20,97	<i>tener</i>	3,22	<i>tratar</i>	0,01
<i>administración</i>	20,54	<i>emparentar</i>	3,21	<i>humanar</i>	0,01
<i>tener</i>	20,47	<i>abrazar</i>	3,19	<i>entrelazado</i>	0,00

Tabla 2-53 Resultados ZChebyshev, adverbio estrechamente + Verbo

<b>Max MAD</b>		<b>Punto de corte MAD</b>		<b>Min MAD</b>	
<i>unir</i>	165,66	<i>bajar</i>	5,00	<i>suficiencia</i>	0
<i>abrazar</i>	120,99	<i>guardar</i>	5,00	<i>inválido</i>	0
<i>unido</i>	110,99	<i>ideo</i>	5,00	<i>público</i>	0
<i>ligar</i>	107,99	<i>general</i>	5,00	<i>quieto</i>	0
<i>relacionar</i>	103,33	<i>señor</i>	5,00	<i>público</i>	0
<i>ligado</i>	92,666	<i>conectar</i>	4,74	<i>conminar</i>	0
<i>ligado</i>	92,666	<i>puntar</i>	4,66	<i>etiquetar</i>	0
<i>abrazar</i>	90,500	<i>sentido</i>	4,66	<i>mortificar</i>	0
<i>ligar</i>	80,749	<i>poblar</i>	4,66	<i>semillar</i>	0
<i>relacionado</i>	71,999	<i>apretar</i>	4,60	<i>público</i>	0

Tabla 2-54 Resultados MAD, adverbio estrechamente + Verbo

<b>Max Box-Plot</b>		<b>Min Box-Plot</b>	
<i>unir</i>	39,97	<i>preguntar</i>	-1,27
<i>relacionar</i>	33,93	<i>resultar</i>	-1,27
<i>ligar</i>	33,45	<i>esperar</i>	-1,27
<i>vincular</i>	24,17	<i>dirigir</i>	-1,29
<i>abrazar</i>	21,14	<i>pensar</i>	-1,30
<i>enlazar</i>	9,08	<i>sentar</i>	-1,30
<i>vigilar</i>	8,05	<i>escribir</i>	-1,32
<i>asociar</i>	6,69	<i>salir</i>	-1,36
<i>aliar</i>	5,55	<i>vezar</i>	-1,39
<i>ceñir</i>	4,44	<i>llevar</i>	-1,45

Tabla 2-55 Resultados *Box-Plot*, adverbio estrechamente + Verbo

		<b>ZChebyshev</b>
<i>frugalmente</i>	<i>comer</i>	1,76
<i>frugalmente</i>	<i>desayunar</i>	0,72
<i>frugalmente</i>	<i>almorzar</i>	0,38
<i>frugalmente</i>	<i>vivir</i>	0,38
<i>frugalmente</i>	<i>desayunar</i>	0,33
<i>frugalmente</i>	<i>cenar</i>	0,27
<i>frugalmente</i>	<i>cenar</i>	0,24
<i>frugalmente</i>	<i>almorzar</i>	0,24
<i>frugalmente</i>	<i>vivir</i>	0,24
<i>frugalmente</i>	<i>comer</i>	0,20

Tabla 2-56 Resultados ZChebyshev, adverbio frugalmente + Verbo

		<b>Box-Plot</b>
<i>frugalmente</i>	<i>comer</i>	72,30
<i>frugalmente</i>	<i>cenar</i>	68,09
<i>frugalmente</i>	<i>desayunar</i>	45,83
<i>frugalmente</i>	<i>almorzar</i>	21,45
<i>frugalmente</i>	<i>vivir</i>	7,56

Tabla 2-57 Resultados Box-Plot, adverbio frugalmente + Verbo

		<b>MAD</b>
<i>frugalmente</i>	<i>comer</i>	18,99
<i>frugalmente</i>	<i>desayunar</i>	3
<i>frugalmente</i>	<i>comer</i>	1,399
<i>frugalmente</i>	<i>cenar</i>	1
<i>frugalmente</i>	<i>almorzar</i>	0,666
<i>frugalmente</i>	<i>vivir</i>	0,642
<i>frugalmente</i>	<i>desayunar</i>	0,500
<i>frugalmente</i>	<i>cenar</i>	0,199
<i>frugalmente</i>	<i>vivir</i>	0
<i>frugalmente</i>	<i>almorzar</i>	0

Tabla 2-58 Resultados MAD, adverbio frugalmente + Verbo

Para finalizar, cabe resaltar en este apartado la necesidad del tratamiento diferencial de los datos según la semántica de la palabra, ya que incide directamente en su uso. Véase en los ejemplos que palabras con valor funcional como el verbo *dar*, o el adverbio *estrechamente* abarcan un espectro más amplio de combinaciones, frente a *propinar*, *omiso* o *frugalmente*, que no solo son palabras de menor frecuencia de uso, sino que en el *DRAE* solo disponen de una acepción (Tabla 2-59).

	Acepciones	Frecuencia
<i>dar</i>	53	1.118.012
<i>propinar</i>	3	2.970
<i>contrato</i>	1	12.244
<i>tierra</i>	10	286.387
<i>público</i>	8	109.763
<i>omiso</i>	1	1.581
<i>estrechamente</i>	4	3.188
<i>frugalmente</i>	1	73

Tabla 2-59 Acepciones y frecuencias en el corpus de los casos analizados.

Cuando se contrastan las combinaciones mejor posicionadas según cada indicador analizado se manifiestan nuevamente las diferencias entre palabras cuya frecuencia de uso es dispar. Las tablas Tabla 2-60 y Tabla 2-61 muestran que en general los casos mejor posicionados coinciden cuando se trata de una palabra cuya frecuencia no es alta y presenta pocas acepciones. Por el contrario, la información mutua arroja resultados similares a las mayores frecuencias relativas evaluadas respecto al otro elemento de la combinación. El resto de los indicadores aportan resultados similares.

	FrecRel	FrecRel	InfMutua	ZScore	t-Score	ZChebyshev	MAD	Box-Plot
<i>dar</i>	<i>adestrador</i>	<i>vuelto</i>	<i>adestrador</i>	<i>vuelta</i>	<i>vuelto</i>	<i>vuelto</i>	<i>vuelta</i>	<i>vuelto</i>
	<i>moravedí</i>	<i>vuelta</i>	<i>moravedí</i>	<i>vuelto</i>	<i>vuelta</i>	<i>vuelta</i>	<i>vuelto</i>	<i>vuelta</i>
	<i>carpetazo</i>	<i>hombre</i>	<i>carpetazo</i>	<i>palmada</i>	<i>hombre</i>	<i>golpe</i>	<i>palmado</i>	<i>hombre</i>
	<i>traspíe</i>	<i>señora</i>	<i>traspíe</i>	<i>palmado</i>	<i>dios</i>	<i>pasa</i>	<i>palmada</i>	<i>señor</i>
	<i>palmada</i>	<i>señor</i>	<i>palmada</i>	<i>salto</i>	<i>diosa</i>	<i>paso</i>	<i>gracia</i>	<i>señora</i>
	<i>palmado</i>	<i>dios</i>	<i>palmado</i>	<i>golpe</i>	<i>pasa</i>	<i>gracia</i>	<i>pasa</i>	<i>diosa</i>
	<i>respingo</i>	<i>diosa</i>	<i>respingo</i>	<i>gracia</i>	<i>paso</i>	<i>muestra</i>	<i>paso</i>	<i>dios</i>
	<i>golpecito</i>	<i>mujer</i>	<i>golpecito</i>	<i>golpecito</i>	<i>señora</i>	<i>salto</i>	<i>priso</i>	<i>mujer</i>

Tabla 2-60 Colocaciones según el indicador, verbo dar + Sustantivo



	<b>Frec. Rel.</b>	<b>Frec. Rel. propinar</b>	<b>InfMutua</b>	<b>ZScore</b>	<b>t-Score</b>	<b>ZChebyshev</b>	<b>MAD</b>	<b>Box-Plot</b>
<b>propinar</b>	<i>paliza</i>	<i>golpe</i>	<i>paliza</i>	<i>paliza</i>	<i>golpe</i>	<i>golpe</i>	<i>golpe</i>	<i>golpe</i>
	<i>puntapié</i>	<i>paliza</i>	<i>puntapié</i>	<i>puntapié</i>	<i>paliza</i>	<i>paliza</i>	<i>paliza</i>	<i>paliza</i>
	<i>puñetazo</i>	<i>bueno</i>	<i>puñetazo</i>	<i>puñetazo</i>	<i>patada</i>	<i>bueno</i>	<i>paliza</i>	<i>bueno</i>
	<i>puñete</i>	<i>patada</i>	<i>puñete</i>	<i>puñete</i>	<i>bueno</i>	<i>patada</i>	<i>bueno</i>	<i>patada</i>
	<i>patada</i>	<i>buena</i>	<i>patada</i>	<i>patada</i>	<i>puntapié</i>	<i>buena</i>	<i>patada</i>	<i>buena</i>
	<i>codazo</i>	<i>puntapié</i>	<i>codazo</i>	<i>golpe</i>	<i>buena</i>	<i>puntapié</i>	<i>buena</i>	<i>puntapié</i>
	<i>bofetada</i>	<i>puñete</i>	<i>bofetada</i>	<i>bofetada</i>	<i>puñetazo</i>	<i>puñetazo</i>	<i>puntapié</i>	<i>puñetazo</i>
	<i>bofetón</i>	<i>puñetazo</i>	<i>bofetón</i>	<i>codazo</i>	<i>puñete</i>	<i>puñete</i>	<i>puñete</i>	<i>puñete</i>

Tabla 2-61 Colocaciones según el indicador, verbo propinar + Sustantivo

	<b>Frec. Rel<sub>Adj</sub></b>	<b>Frec. Rel<sub>tierra</sub></b>	<b>Inf. Mutua</b>	<b>ZScore</b>	<b>t-Score</b>	<b>ZChebyshev</b>	<b>Box Plot</b>	<b>MAD</b>
<b>tierra</b>	<i>temporejado</i>	<i>tanto</i>	<i>temporejado</i>	<i>firme</i>	<i>baja</i>	<i>firme</i>	<i>baja</i>	<i>baja</i>
	<i>lejas</i>	<i>baja</i>	<i>lejas</i>	<i>fértil</i>	<i>bajo</i>	<i>cielo</i>	<i>bajo</i>	<i>bajo</i>
	<i>estercolado</i>	<i>bajo</i>	<i>estercolado</i>	<i>apisonado</i>	<i>mucho</i>	<i>firmar</i>	<i>nuevo</i>	<i>firme</i>
	<i>labrantío</i>	<i>mucho</i>	<i>labrantío</i>	<i>natal</i>	<i>grande</i>	<i>adentrar</i>	<i>bien</i>	<i>medio</i>
	<i>apisonado</i>	<i>grande</i>	<i>apisonado</i>	<i>baldío</i>	<i>nuevo</i>	<i>hacer</i>	<i>padre</i>	<i>vido</i>
	<i>lantánido</i>	<i>nuevo</i>	<i>lantánido</i>	<i>labrantío</i>	<i>tanto</i>	<i>mar</i>	<i>firme</i>	<i>indio</i>
	<i>cultivable</i>	<i>señor</i>	<i>cultivable</i>	<i>lejas</i>	<i>firme</i>	<i>cultivar</i>	<i>medio</i>	<i>alto</i>
	<i>alquiladizo</i>	<i>señora</i>	<i>alquiladizo</i>	<i>poblado</i>	<i>señora</i>	<i>tener</i>	<i>vido</i>	<i>visto</i>

Tabla 2-62 Colocaciones según el indicador, sustantivo tierra + Adjetivo

	<b>Frec.</b>	<b>Frec.</b>	<b>Inf</b>	<b>ZScore</b>	<b>t-Score</b>	<b>ZChebyshev</b>	<b>Box</b>	<b>MAD</b>
	<b>Rel<sub>Adj</sub></b>	<b>Rel</b>	<b>Mutua</b>				<b>Plot</b>	
<b>contrato</b>	<i>rescindido</i>	<i>poder</i>	<i>rescindido</i>	<i>matrimonial</i>	<i>social</i>	<i>social</i>	<i>social</i>	<i>social</i>
	<i>sinalagmático</i>	<i>hacer</i>	<i>sinalagmático</i>	<i>firmado</i>	<i>derecha</i>	<i>derecha</i>	<i>derecha</i>	<i>derecha</i>
	<i>prenupcial</i>	<i>firmar</i>	<i>prenupcial</i>	<i>rescindido</i>	<i>derecho</i>	<i>derecho</i>	<i>derecho</i>	<i>derecho</i>
	<i>revisable</i>	<i>tener</i>	<i>revisable</i>	<i>celebrado</i>	<i>público</i>	<i>público</i>	<i>público</i>	<i>público</i>
	<i>consensual</i>	<i>celebrar</i>	<i>consensual</i>	<i>usurario</i>	<i>firmado</i>	<i>firmado</i>	<i>firmado</i>	<i>firmado</i>
	<i>usurario</i>	<i>decir</i>	<i>usurario</i>	<i>social</i>	<i>hecho</i>	<i>hecho</i>	<i>hecho</i>	<i>hecho</i>
	<i>enfitéutico</i>	<i>social</i>	<i>enfitéutico</i>	<i>consensual</i>	<i>celebrado</i>	<i>nuevo</i>	<i>nuevo</i>	<i>nuevo</i>
	<i>canjeado</i>	<i>dar</i>	<i>canjeado</i>	<i>sinalagmático</i>	<i>nuevo</i>	<i>celebrado</i>	<i>celebrado</i>	<i>celebrado</i>

Tabla 2-63 Colocaciones según el indicador, sustantivo contrato + Adjetivo

### 2.7.2. Resultados obtenidos sobre el corpus de Galdós

Se centra la discusión en las mismas formas utilizadas para ilustrar los resultados del epígrafe anterior, se incluye en la Tabla 2-64 sus frecuencias calculadas sobre el Corpus de Galdós.

	<b>Frecuencia</b>	<b>Combinaciones</b>
<i>dar</i>	9628	972
<i>propinar</i>	15	0
<i>tierra</i>	689	88
<i>contrato</i>	20	1
<i>público</i>	461	66
<i>omiso</i>	3	3
<i>estrechamente</i>	9	1
<i>frugalmente</i>	1	0

Tabla 2-64 Frecuencias en el corpus de Galdós de los ejemplos

Si se utiliza la información mutua o el valor de *Box-plot* en los valores máximos se encuentran las mismas colocaciones mejor posicionadas. Se observa que este último es más eficaz usándose en corpus menos extensos, ya que no surgen las distorsiones originadas por las palabras de uso muy frecuente frente a otras menos habituales. La palabra *carpetazo* aparece 3 veces en este corpus, todas ellas con el verbo *dar* (Tabla 2-65).

	<i>ZCh</i>	<i>MAD</i>	<i>Box- Plot</i>	<i>InfM</i>	<i>Frec R</i>					
<i>dar</i>	<i>vuelta</i>	10,07	<i>vuelta</i>	85,50	<i>carpetazo</i>	23,71	<i>carpetazo</i>	2,91	<i>vuelta</i>	0,03
	<i>vuelto</i>	10,07	<i>vuelto</i>	85,50	<i>atracque</i>	23,71	<i>atracque</i>	2,69	<i>vuelto</i>	0,03
	<i>pasa</i>	8,09	<i>paso</i>	69,25	<i>atraco</i>	19,07	<i>atración</i>	2,69	<i>pasa</i>	0,03
	<i>paso</i>	8,09	<i>pasa</i>	69,25	<i>atración</i>	19,07	<i>atraco</i>	2,69	<i>paso</i>	0,03
	<i>gana</i>	5,41	<i>gana</i>	47,25	<i>morrada</i>	12,79	<i>morrada</i>	2,59	<i>dar</i>	0,02
	<i>diosa</i>	5,01	<i>diosa</i>	44,00	<i>grima</i>	11,86	<i>grima</i>	2,59	<i>gana</i>	0,02
	<i>dios</i>	5,01	<i>dios</i>	44,00	<i>palmada</i>	11,07	<i>palnado</i>	2,56	<i>diosa</i>	0,02
	<i>mano</i>	4,67	<i>mano</i>	41,25	<i>palnado</i>	11,07	<i>palmada</i>	2,56	<i>dios</i>	0,02
	<i>manos</i>	4,64	<i>manos</i>	41,00	<i>pábulo</i>	11,00	<i>pábulo</i>	2,55	<i>mano</i>	0,02
	<i>cuenta</i>	4,34	<i>cuenta</i>	38,50	<i>traste</i>	10,29	<i>traste</i>	2,53	<i>manos</i>	0,02

Tabla 2-65 Colocaciones en el corpus de Galdós según el indicador, verbo dar

	<i>ZCh</i>	<i>MAD</i>	<i>Box- Plot</i>	<i>InfM</i>	<i>Frec.</i>					
<i>tierra</i>	<i>cielo</i>	5,67	<i>cielo</i>	38,00	<i>espuerta</i>	52,61	<i>espuerta</i>	4,49	<i>cielo</i>	0,060
	<i>echar</i>	3,12	<i>echar</i>	22,00	<i>entraña</i>	37,37	<i>entraña</i>	3,20	<i>echar</i>	0,037
	<i>ver</i>	3,12	<i>ver</i>	22,00	<i>entraño</i>	37,37	<i>entraño</i>	3,20	<i>ver</i>	0,037
			<i>hacer</i>	13,00	<i>expediente</i>	37,05	<i>expediente</i>	3,17	<i>hacer</i>	0,024
			<i>dar</i>	13,00	<i>expediente</i>	37,05	<i>expediente</i>	3,17	<i>dar</i>	0,024
			<i>venir</i>	12,00	<i>cielo</i>	29,61	<i>expedientar</i>	3,17	<i>haber</i>	0,023
			<i>haber</i>	12,00	<i>fecundo</i>	29,06	<i>entrañar</i>	3,05	<i>venir</i>	0,023
			<i>bien</i>	10,00	<i>montón</i>	18,23	<i>cielo</i>	2,54	<i>tener</i>	0,02
			<i>tener</i>	10,00	<i>mortal</i>	17,14	<i>fecundo</i>	2,49	<i>ir</i>	0,02
			<i>ir</i>	10,00	<i>centro</i>	13,56	<i>rodrigar</i>	2,44	<i>bien</i>	0,02

Tabla 2-66 Colocaciones en el Corpus de Galdós, sustantivo tierra

	<b>ZCh</b>		<b>MA</b>		<b>Box</b>		<b>Inf.</b>	<b>FrecR</b>		
			<b>D</b>		<b>Plot</b>					
<b>público</b>	<i>orden</i>	4,11	<i>órdenes</i>	23	<i>orden</i>	5,4	<i>privada</i>	4,80	<i>órdenes</i>	0,06
	<i>órdenes</i>	4,11	<i>orden</i>	23	<i>órdenes</i>	5,4	<i>instrucciones</i>	4,66	<i>orden</i>	0,06
			<i>público</i>	13			<i>instrucción</i>	4,66	<i>hacer</i>	0,039
			<i>público</i>	11			<i>implorar</i>	4,49	<i>ver</i>	0,03
			<i>público</i>	11			<i>notorio</i>	3,97	<i>dar</i>	0,03
			<i>instrucciones</i>	8			<i>vía</i>	3,94	<i>instrucción</i>	0,02
			<i>instrucción</i>	8			<i>privado</i>	3,93	<i>instrucciones</i>	0,02
			<i>decir</i>	7			<i>sitios</i>	3,27	<i>decir</i>	0,02
			<i>hombre</i>	7			<i>privar</i>	3,00	<i>hombre</i>	0,02
		<i>voz</i>	6			<i>hacienda</i>	2,93	<i>voz</i>	0,02	

Tabla 2-67 Colocaciones en el Corpus de Galdós, adjetivo público

Se manifiesta la dificultad a la hora de encontrar muestras de las combinaciones que se ajusten a estructuras colocacionales en español cuando el corpus es significativamente menos extenso: *contrato* y *estrechamente* solo producen una, *omiso* dos y *frugalmente* ninguna. Es por ello que *ZChebyshev*, *MAD* y *Box-Plot* no pueden seleccionar combinaciones de uso preferente, puesto que la ausencia de elementos donde evaluar la variación impide captar anomalías en las frecuencias relativas (Tabla 2-68).

		<b>Inf.</b>	<b>Frec.</b>
		<b>Mutua</b>	<b>Relativa</b>
<i>contrato</i>	<i>hacer</i>	-0,22	0,15
<i>omiso</i>	<i>caso</i>	5,82	1
<i>omiso</i>	<i>casar</i>	3,26	1
<i>omiso</i>	<i>hacer</i>	2,512	1
<i>estrechamente</i>	<i>abrazar</i>	6,65	0,44

Tabla 2-68 Ejemplos en los que no se encuentra variación

## 3. Detección de nexos con grupos semánticos

### 3.1. Análisis cualitativo de los datos y selección léxica

El análisis estadístico del corpus permite extraer una gran cantidad de combinaciones que representan las propiedades combinatorias de una amplia muestra del español. Las estrategias seguidas para la obtención de las colocaciones se han centrado en determinar aquellas combinaciones que reflejan el uso preferente de las mismas frente a cualquier otra combinación posible. Sin embargo, se hace patente en los conjuntos extraídos para cualquier forma canónica, que se pueden establecer grupos o clases a partir de criterios semánticos. Se observa que las bases con las que combina un determinado verbo, adverbio o sustantivo en las colocaciones *verbo + sustantivo*, *adverbio + verbo* o *sustantivo + adjetivo* pertenecen a grupos semánticos, los elementos recuperados mediante las técnicas de detección de outliers representan los de uso preferente. Se presenta como ejemplo de este fenómeno algunos de los sustantivos que se obtienen con los verbos *sentir* y *palpitar*, los adverbios *enérgicamente* y *abundamente* y los adjetivos *brusco* y *doloroso* (Tabla 3-1).

Estos resultados se pueden explicar a partir del fenómeno de la selección léxica. En Bosque (2001) se resalta la interpretación de las colocaciones como un caso particular del concepto más amplio de selección léxica. Una determinada base de forma aislada no constituye el único caso compatible para un determinado colocativo. En contraposición, éste la selecciona por sus rasgos semánticos, que son compartidos por otros sustantivos en el caso *verbo + sustantivo*, otros verbos para el tipo *adverbio + verbo* o bien sustantivos cuando se trata del caso *sustantivo + adjetivo*. Este enfoque también está presente en (Serra, 2009), quien defiende la idea de que los grupos que se aprecian en el conjunto de bases dado un colocativo corresponden a casos de la relación *predicado - argumentos*. El DCECR, utilizado como conjunto de ensayo en este trabajo, presenta las propiedades combinatorias de verbos, adjetivos y adverbios de forma intencional y extensional. La información que proporciona este diccionario, correspondería a los predicados, y no se limita a la pura enumeración de combinaciones de uso más o menos frecuente, sino que se muestran grupos de argumentos integrantes de una clase semántica, junto con los rasgos semánticos que comparten. La frecuencia en este caso se considera irrelevante. Si se consulta en el DCECR sobre los verbos *sentir* y *palpitar* se encuentra la información en la Tabla 3-2 y la Tabla 3-3 –aparecen resaltados en negrilla los casos detectados como colocaciones según los criterios estadísticos que se han aplicado en el capítulo precedente de este trabajo.

<b>Forma Canónica</b>	<b>Grupos</b>
<i>sentir</i>	<i>náusea, escalofrío, hambre.</i>  <i>roce, caricia.</i>  <i>estómago, garganta, pierna.</i>  <i>pánico, terror, horror.</i>  <i>alivio, cansancio, presión, angustia.</i>  <i>compasión, lástima.</i>  <i>envidia, desprecio.</i>  <i>simpatía, admiración, orgullo, respeto, confianza, felicidad, alegría.</i>  <i>remordimiento, vergüenza.</i>
<i>palpitar</i>	<i>pulso, latido.</i>  <i>compás, ritmo.</i>  <i>sién, vena, vientre, entraña, mejilla, cuello.</i>  <i>motor.</i>  <i>miedo, esperanza.</i>  <i>placer, pasión.</i>
<i>enérgicamente</i>	<i>protestar, decir.</i>
<i>abundantemente</i>	<i>proveer, producir, surtir.</i>  <i>llorar, sudar.</i>  <i>comer, beber.</i>
<i>brusco</i>	<i>gesto, movimiento.</i>  <i>modo, manera.</i>
<i>doloroso</i>	<i>experiencia, momento, vida, días.</i>  <i>corazón, alma.</i>

Tabla 3-1 Colocaciones agrupadas por la semántica

Identificador	Rasgo Semántico	Clase
(I)	SUSTANTIVOS QUE DESIGNAN SENTIMIENTOS DE INCLINACIÓN APRECIO PREFERENCIA SATISFACCIÓN Y OTROS ANÁLOGOS	<i>admiración, afecto, amor, aprecio, cariño, consideración, debilidad, estima, orgullo, predilección, respeto, simpatía</i>
(II)	SUSTANTIVOS QUE DENOTAN PASIÓN O DEVOCIÓN	<i>adoración, delirio, devoción, fascinación, locura, pasión, veneración</i>
(III)	SUSTANTIVOS QUE DENOTAN ATRACCIÓN HACIA ALGO O ALGUIEN O DESEO DE CONOCER O EXPERIMENTAR ALGUNA COSA	<i>afición, atracción, curiosidad, interés, seducción</i>
(IV)	SUSTANTIVOS QUE DENOTAN ODIO, RECHAZO O DESEO DE ALEJAMIENTO	<i>aversión, desprecio, distanciamiento, fobia, odio, rechazo</i>
(V)	SUSTANTIVOS QUE DENOTAN TRISTEZA O CONMISERACIÓN	<i>compasión, conmiseración, lástima, pena, tristeza</i>
(VI)	SUSTANTIVOS QUE DENOTAN PREOCUPACIÓN, INCERTIDUMBRE O RESERVA	<i>inquietud, miedo, preocupación, reparo, reserva, temor, terror, zozobra</i>
(VII)	SUSTANTIVOS QUE DENOTAN CULPA	<i>culpa, culpabilidad</i>

Tabla 3-2 Clases semánticas en el DCECR para el verbo sentir

Listados similares se obtienen en las entradas de los adverbios *enérgicamente* y *abundantemente*, así como el adjetivo brusco, que quedan recogidos en la Tabla 3-4, Tabla 3-6 y Tabla 3-5, respectivamente. Sin embargo, no hay entrada para el adjetivo *doloroso*.

Esta teoría resulta especialmente útil de cara a afrontar la tarea de establecer grupos semánticos entre las combinaciones recuperadas para una forma canónica dada, puesto que explica el fenómeno observado. Por una parte, determina la orientación que se debe dar a la explotación de los resultados, indica qué elemento se debe fijar, y en qué categoría gramatical establecer los grupos. Son tres los casos que se contemplan: fijar un verbo y agrupar los sustantivos con los que se combina, fijar un adverbio y agrupar los verbos a los que modifica, o fijar un adjetivo y agrupar sus sustantivos. Por tanto, se hace necesario diseñar alguna estrategia

que permita agrupar mediante criterios semánticos las bases; se trata de establecer el nexo entre el predicado y sus argumentos de forma extensional.

El problema consiste en delimitar el subconjunto formado por aquellas que responden al fenómeno de la selección léxica dentro del conjunto de combinaciones posibles para un predicado dado. Desde un punto de vista lingüístico, el subconjunto de las colocaciones se encuentra en un nivel más profundo de las restricciones combinatorias. Si se traslada este esquema a la automatización del proceso, se parte de todas las combinaciones extraídas del corpus, entre las que se incluyen grupos semánticos y en el nivel más profundo, los casos que destacan por su uso frecuente o outliers (Figura 3-1).

<b>Rasgo Semántico</b>	<b>Clase</b>
SUSTANTIVOS QUE DESIGNAN EMOCIONES Y SENTIMIENTOS PRINCIPALMENTE LOS DE NATURALEZA PASIONAL	<i>amor, deseo, emoción, pasión, pulsión</i>
SUSTANTIVOS QUE DENOTAN MIEDO ODIOS O PREOCUPACIÓN TAMBIÉN CON ALGUNOS QUE DESIGNAN OTRAS ACTITUDES QUE EXPRESAN LA INTRANQUILIDAD MANIFESTADA CON RELACIÓN A ALGO	<i>angustia, ansiedad, desconfianza, <b>miedo</b>, odio, preocupación, resentimiento, temor</i>
SUSTANTIVOS QUE DENOTAN FUERZA O VITALIDAD	<i>fuerza, vida, vitalidad</i>
SUSTANTIVOS QUE DESIGNAN ARTES TÉCNICAS Y TENDENCIAS ARTÍSTICAS ASÍ COMO ALGUNAS DE SUS MANIFESTACIONES	<i>acorde, creación, cubismo, flamenco, literatura, melodía, obra, partitura, poesía, tendencia</i>
SUSTANTIVOS QUE DENOTAN LUGAR CASI SIEMPRE USADOS EN SENTIDO METONÍMICO CON LA INTERPRETACIÓN DE GRUPO HUMANO	<i>banquillo, ciudad, país</i>

Tabla 3-3 Clases semánticas en el DCECR para el verbo palpar



Rasgo Semántico	Clase
VERBOS QUE DENOTAN OPOSICIÓN O RECHAZO A MENUDO MANIFESTADAS VERBALMENTE	<i>condenar, criticar, denunciar, descartar, oponerse, protestar, rechazar, reprobado, repudiar</i>
LOS VERBOS DEFENDER Y APOYAR	<i>apoyar, defender</i>
VERBOS QUE DENOTAN LA ACCIÓN DE AFIRMAR O NEGAR ALGUNA COSA	<i>afirmar, asegurar, desmentir, negar</i>
VERBOS QUE DENOTAN SOLICITUD	<i>demandar, exigir, pedir, preguntar, reclamar</i>
VERBOS QUE DENOTAN ENFRENTAMIENTO O CONFRONTACIÓN, A MENUDO NO FÍSICA	<i>combatir, discutir, luchar</i>
VERBOS QUE DESIGNAN DIVERSAS MANIFESTACIONES VERBALES CARACTERIZADAS POR LA REITERACIÓN O EL ÉNFASIS	<i>recalcar, reiterar, repetir</i>
VERBOS QUE DENOTAN AMENAZA O ADMONICIÓN	<i>advertir, amenazar, recomendar</i>

Tabla 3-4 Clases Semánticas en el DCECR. para el adverbio enérgicamente



Figura 3-1 Esquema del problema en el nivel lingüístico y en el nivel computacional.

<b>Rasgo Semántico</b>	<b>Clase</b>
SUSTANTIVOS QUE DENOTAN CAMBIO REFORMA O MODIFICACIÓN TAMBIÉN CON OTROS QUE DESIGNAN DIVERSOS ESTADOS SUCESIVOS DE ALTERNANCIA	<i>alteración, altibajo, corrección, fluctuación, mutación, oscilación, reforma, transformación, vaivén, variación</i>
SUSTANTIVOS QUE DENOTAN AUMENTO O ASCENSO	<i>alza, ascensión, ascenso, aumento, crecimiento, elevación, incremento, subida</i>
SUSTANTIVOS QUE DENOTAN DISMINUCIÓN O DESCENSO	<i>bajada, bajada, caída, derrumbe, descenso, disminución, rebaja, recorte, reducción</i>
SUSTANTIVOS QUE DENOTAN MOVIMIENTO, Y A MENUDO CAMBIO RÁPIDO DE DIRECCIÓN	<i>giro, maniobra, movimiento, quiebro, viraje, volantazo, vuelco</i>
SUSTANTIVOS QUE DENOTAN DETENCIÓN O VARIACIÓN REPENTINA DE VELOCIDAD	<i>aceleración, acelerón, desaceleración, detención, frenada, frenazo, paralización, parón</i>
SUSTANTIVOS QUE DENOTAN DIVISIÓN, CORTE O RUPTURA	<i>corte, escisión, interrupción, quiebra, rotura, ruptura, separación</i>
EL SUSTANTIVO GESTO Y CON OTROS QUE DESIGNAN ACCIONES Y MOVIMIENTOS CORPORALES, A MENUDO RÁPIDOS O REPENTINOS	<i>abrazo, ademán, brinco, cabezada, gesto, pirueta</i>
SUSTANTIVOS QUE DENOTAN GOLPE, CHOQUE, ENCUENTRO VIOLENTO O MOVIMIENTO IMPULSIVO O IMPETUOSO	<i>chapuzón, choque, embestida, empujón, encontronazo, golpe, remezón, sacudida, tirón, vapuleo, zapatillazo, zarpazo</i>
SUSTANTIVOS QUE DENOTAN APARICIÓN, INICIO O SURGIMIENTO DE ALGO, O DESIGNAN ALGUNAS DE LAS ACCIONES QUE LLEVAN A CABO ESOS PROCESOS	<i>aparición, apertura, arrancada, arranque, eclosión, irrupción, rebrote, surgimiento</i>

Tabla 3-5 clases semánticas en el DCECR. para el adjetivo brusco (I)

<b>Rasgo Semántico</b>	<b>Clase</b>
VERBOS DE CAMBIO DE ESTADO ESPECIALMENTE LOS QUE DESIGNAN LA ACCIÓN DE PROPORCIONAR HACER ADQUIRIR O LOGRAR QUE ALGO O ALGUIEN PASE A ESTAR EN DISPOSICIÓN DE ALGUNA COSA	<i>dotar, engrasar, financiar, pintar, premiar, proporcionar, suministrar, untar</i>
VERBOS QUE DENOTAN CONSUMICIÓN GENERALMENTE DE ALIMENTOS	<i>abreviar, beber, cenar, comer, consumir, leer</i>
VERBOS QUE DENOTAN EXISTENCIA O INCREMENTO DE ALGO	<i>almacenar(se), crecer, dar(se), encontrar(se), existir, extender(se), practicar(se), presentar(se), prodigar(se), producir(se), realizar(se), sobrar</i>
VERBOS QUE DESIGNAN EL PROCESO DE SURGIR O FLUIR UN LÍQUIDO U OTRAS MATERIAS QUE SE LE ASIMILAN FÍSICA O FIGURADAMENTE TAMBIÉN CON OTROS VERBOS QUE EXPRESAN LA ACCIÓN QUE DESENCADENA ESE PROCESO	<i>brotar, correr, desbordar, fluir, llover, manar, nevar, perder sangre, regar, remojar, rociar, sangrar, sudar, verter</i>
VERBOS QUE DESIGNAN LA ACCIÓN DE EXPRESAR O MANIFESTAR ALGO GENERALMENTE COMO MEDIOS VERBALES PERO TAMBIÉN CON OTROS	<i>anunciar, citar, comentar, contar, describir, dialogar, discutir, escribir, hablar, informar, mencionar</i>
VERBOS QUE DESIGNAN LA ACCIÓN DE PROBAR ALGO	<i>demostrar, probar, dar propaganda</i>
VERBOS QUE DESIGNAN LA ACCIÓN DE REPRODUCIR DIFUNDIR O HACER PÚBLICA ALGUNA COSA	<i>exhibir, filmar, grabar, ilustrar, representar, reproducir, reseñar, traducir, versionear</i>
VERBOS QUE DESIGNAN LA ACCIÓN DE USAR ALGO	<i>emplear, recurrir, usar, utilizar</i>

Tabla 3-6 Clases Semánticas en el DCECR. para el adverbio abundantemente.

### 3.2. Similitud semántica

Se propone aplicar técnicas que permitan automatizar el proceso de generación de clases léxicas; deben estar formadas por palabras que estén próximas en el espacio semántico, de manera que se puedan construir grupos formados por elementos similares en cuanto a su significado. La generación automática de grupos semánticos obliga a cuantificar las relaciones semánticas existentes entre las palabras del texto que se analiza. Se debe hacer uso de medidas que reflejen la relación entre palabras que representan un mismo concepto, y así conformar grupos que hagan referencia a las mismas entidades conceptuales.

El problema de cuantificar las relaciones semánticas ha sido ampliamente estudiado por su utilidad en distintas tareas del Procesamiento del Lenguaje Natural. Se emplean diferentes términos concernientes a métricas que evalúan tanto la existencia de una relación semántica entre dos palabras como su fuerza. En este sentido se pueden encontrar en la literatura referencias a *relación*, *similitud* o *distancia semántica*.

#### 3.2.1. Técnicas para determinar la similitud semántica

Desde el punto de vista algorítmico, los métodos comparan los términos en la representación topológica de algún recurso léxico que incluya relaciones semánticas, o bien se analizan contextos en los que aparecen –en muchas ocasiones se representan mediante vectores.

En cualquier caso, una medida de similitud debe verificar los siguientes requisitos (Lin, 1998):

- i. Cuando las dos unidades son idénticas, el valor de la similitud debe ser máximo.
- ii. Dos unidades son más similares cuantos más elementos comparten. Por tanto, la similitud está relacionada con los elementos en común.
- iii. Dos unidades son menos similares cuantas más diferencias haya entre ellos. La similitud está relacionada con las diferencias entre ellas.

##### 3.2.1.1. Métodos basados en Bases de Conocimientos léxicas.

Se utilizan bases de conocimiento léxicas, ya sean diccionarios, tesauros, ontologías o redes semánticas: algún recurso léxico con relaciones semánticas bien de forma explícita o bien que se puedan inferir como en el caso de los diccionarios convencionales.

Kozima y Furugori (1993) evalúan la similitud semántica a partir de una red neuronal en la que los nodos son entradas del diccionario. Dos entradas tienen relación si una de ellas se usa para definir a la otra. Sobre esta red se calcula la similitud a partir de la activación de los nodos.

Otros autores como Morris y Hirst (1998) utilizan un tesauro, que agrupa las palabras sin especificar qué tipo de relación motiva el agrupamiento. Dos palabras se consideran cercanas semánticamente si comparten alguna categoría, están en la misma subcategoría, etc.

Existe un grupo importante de medidas de similitud que toman como espacio topológico en el que proyectar las palabras a *WordNet* –red semántica de la lengua inglesa. Se agrupan las palabras en conjuntos de sinónimos o *SynSets*, que proporcionan definiciones cortas y generales, y almacenan las relaciones semánticas entre estos conjuntos de sinónimos. De esta base de datos también se ha creado una versión multilingüe para varios idiomas europeos (Holandés, Italiano, Español, Alemán, Francés, Checo y Estonio) llamada *EuroWordNet*.

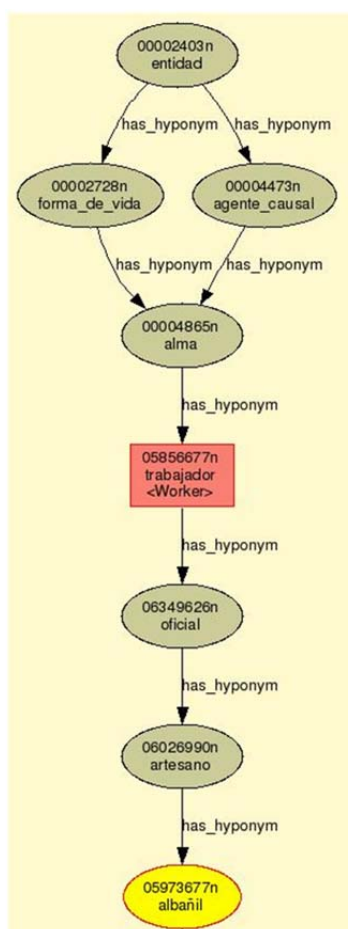


Figura 3-2 Grafos de las consultas albañil y empresa en EuroWordNet.

Si bien este grupo de técnicas considera la similitud semántica en función de la longitud del camino que las une en el grafo que representa la red, esta aproximación es demasiado simple. Considerar exclusivamente el número de arcos que unen los synsets equivale a asumir que todas las taxonomías registradas son equidistantes, sin embargo, la estructura de la red no cumple con tal suposición. Es necesario determinar pesos que reequilibren su naturaleza no homogénea, de forma que los parámetros a contemplar incluyan tanto la longitud del camino

como la especificidad de los sentidos de las palabras en el camino (Agirre, 1995). Este problema se ilustra en (Resnik, 1999) en donde se presentan los enlaces:

Rabbit Ears is-a Television Antenna, o

Phytoplankton is-a Living Thing.

Para resolver este problema existen diferentes propuestas que incorporan datos adicionales, como la profundidad del concepto en la jerarquía. A profundidades mayores se deben obtener mayores puntuaciones; también deben ser más próximos semánticamente hablando aquellos que se encuentren en zonas de mayor concentración de conceptos frente a los que estén en zonas dispersas (Agirre, 1995).

### **3.2.1.2. A partir de los contextos de uso**

El segundo grupo de técnicas para la detección de similitud semántica se basa en la información que se puede extraer de las palabras a partir de los contextos en que se emplean. Su uso exige como requisito indispensable disponer de corpus extensos para que las medidas sean fiables y precisas.

Se evalúa la relación de las palabras mediante distancias que representan la similitud o la fuerza de la relación existente –se basan en la frecuencia de aparición conjunta en los textos. En general, la idea subyacente es que la relación semántica es más fuerte en aquellos grupos de palabras cuyos contextos estén próximos.

Se proyecta la información lingüística sobre un espacio vectorial N dimensional, habitualmente vectores de coocurrencias en el corpus. Sobre dicho espacio se calculan distancias geométricas entre sus elementos: distancia euclídea, distancia de Manhattan, distancia del coseno, etc. Autores como Manning y Schütze (1999) indican una representación de vectores binarios, a los que se pueden aplicar la distancia de Dice, Jaccard, etc.

Además de poder aplicar las distancias anteriores, cuando los vectores representan distribuciones de probabilidad de los objetos lingüísticos en el corpus, surgen otras distancias que explotan las propiedades inherentes a las distribuciones de probabilidad: la Divergencia de Jensen-Shannon, el coeficiente de correlación de Pearson, la divergencia de Kullback-Leibler, o el coeficiente de  $\tau$ -Kendall, entre otras.

## **3.3. Obtención de clases léxicas**

Se pretende construir los grupos semánticos a partir de los datos registrados en la BDD de combinaciones léxicas. El problema que se pretende resolver sobre nuestros datos es el

siguiente: *dada una forma canónica, agrupar en clases semánticas el conjunto de formas canónicas con las que aparece en alguna de las estructuras admisibles para colocaciones del español*. Por ejemplo, dada una forma canónica con categoría gramatical verbo, se deben obtener, en caso de que existan, grupos dentro de la bolsa de sustantivos con los que aparece en las combinaciones *verbo + sustantivo*. El criterio que regirá su formación es la similitud semántica entre sus miembros. Otra posibilidad, que consideraremos de forma independiente sería realizar agrupamientos entre los adverbios con los que se ha registrado en el corpus.

### 3.3.1. Agrupamiento mediante técnicas estadísticas

En este epígrafe revisamos dos técnicas estadísticas que permiten generar grupos semánticos, analizando si su aplicación sobre los datos que manejamos es recomendable. Por una parte se presenta el análisis clúster, dentro de las técnicas convencionales de clasificación estadística, y por otro el análisis de semántica latente.

#### 3.3.1.1. Análisis clúster

La generación de grupos semánticos o clases léxicas puede resolverse a partir de técnicas basadas en análisis *clúster*, tradicionalmente utilizadas para la creación de clasificaciones. Se conoce como *análisis clúster* al conjunto de técnicas multivalentes diseñadas para extraer grupos o *clústeres* de un conjunto de observaciones de un conjunto de variables aleatorias. Estos métodos dividen un conjunto de objetos en grupos o clústeres con cohesión interna, es decir, los objetos en un mismo grupo serán muy similares entre sí. Por otra parte, también se impone la restricción de aislamiento externo del grupo, es decir, que los elementos en clústeres distintos sean muy diferentes.

En esta técnica se hace fundamental la selección de la medida de similitud, siendo importante que se determine en función de una distancia  $d: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^+$  debe verificar (Rodríguez, 2003):

$$\begin{aligned} d(\vec{x}_i, \vec{x}_j) &\geq 0 \\ d(\vec{x}_i, \vec{x}_i) &= 0 \\ d(\vec{x}_i, \vec{x}_j) &= d(\vec{x}_j, \vec{x}_i) \\ d(\vec{x}_i, \vec{x}_j) &\leq d(\vec{x}_i, \vec{x}_k) + d(\vec{x}_k, \vec{x}_j) \end{aligned}$$

Generalmente, los métodos de análisis clúster se clasifican en *jerárquicos* o no *jerárquicos*. En el primer caso se presupone una estructura jerárquica en la clasificación, además los resultados son muy sensibles a los outliers. Dadas las características de la distribución de las combinaciones en las que interviene una forma canónica, totalmente asimétrica y con presencia de valores extremos hace que se descarten estas técnicas en este problema. En el segundo, no se

impone ningún tipo de estructura de antemano, pero se debe prefijar el número de clases que se van a generar. Lo más común en estos casos es aplicar algún método de reasignación: se asignan individuos a una configuración inicial de grupos y a lo largo de las sucesivas iteraciones se permite reasignar individuos a grupos, de manera que se optimice el criterio de selección; el proceso acaba cuando no existan objetos que al ser reasignados provoquen mejoras. Es necesario por tanto, determinar algún mecanismo de evaluación de la bondad de la configuración de los grupos que permita establecer la reasignación de individuos orientada a la optimización. Con este fin se define el *centroide* del clúster como representante de la clase en cada iteración, consistiendo en el vector formado por la media aritmética de los valores en el grupo de cada dimensión. Cada componente del *centroide* del grupo  $i$  es el promedio de los valores de la variable correspondiente a esa componente para todos los elementos que conforman dicho grupo. Es decir:

$$\vec{c}_i = \frac{1}{n_i} \sum_{j=0}^{n_i} \vec{x}_j$$

Cuánto de bueno es el agrupamiento se mide a partir de este vector, los elementos del grupo tienen que tener distancias pequeñas al *centroide* y los grupos deben estar lo más diferenciados posible, lo cual se puede lograr maximizando las distancias entre *centroides*. En cualquier caso, se hace necesario medir las distancias entre los objetos que se pretende agrupar, para lo que se dispone de diferentes funciones tanto para variables cuantitativas como categóricas.

En este caso se establece como contexto de una forma canónica dada, todas las formas canónicas con las que aparece en alguna combinación en el corpus. La Tabla 3-7 ilustra los contextos de algunos sustantivos que se revelan como colocaciones del verbo *sentir*. Cada entrada corresponde a la frecuencia de coocurrencia de las dos formas canónicas, también se pueden utilizar las de frecuencias relativas. En el ejemplo seleccionado, la construcción de la matriz de contextos conlleva extraer todas las bases del verbo *sentir*, éstas constituirán las filas, es decir, los objetos a clasificar. Las columnas, o variables utilizadas para clasificarlas, corresponden a cada uno de los verbos o colocativos con los que aparecen esas bases. Algoritmos de análisis clúster sobre esta matriz permiten obtener agrupamientos bajo el principio de que sustantivos similares se usan en contextos similares: los grupos se construyen bajo la premisa de que *sustantivos similares combinan con verbos similares*.



Sustantivo	<i>abajar</i>	<i>abalar</i>	<i>abandoner</i>	<i>abarcar</i>	<i>abatir</i>	<i>abismar</i>	<i>ablandar</i>	<i>abordar</i>	...
<i>escalofrío</i>	12								
<i>miedo</i>	79	11	322	4	38	29	4	6	
<i>náusea</i>			6						
<i>terror</i>	34	3	111		9	35			

Tabla 3-7 Extracto de vectores de contexto de algunos sustantivos en las colocaciones sentir + Sust.

La matriz de vectores de contexto del verbo *sentir*, restringida a los sustantivos que se detectan como colocaciones desde el punto de vista frecuentista, es de  $12545 \times 7494$ ; cuando se trata del verbo *palpitar*, la matriz es de rango  $661 \times 6955$ . Si se restringen los casos a aquellos que aparecen al menos 10 veces en el corpus se obtiene una matriz de  $8285 \times 5606$  y  $224 \times 4955$  respectivamente.

La construcción de dicha matriz a partir de la BDD es inviable en tiempo real, para cada caso es necesario extraer la bolsa de palabras base y cada una de las formas canónicas que combinan con ellas, luego se asignará la frecuencia de coocurrencia en caso de que la combinación exista en el corpus o 0 en caso contrario. Además se aconseja aplicar alguna técnica de reducción de dimensiones para poder aplicar los algoritmos de análisis clúster, en caso contrario sobre matrices tan extensas sería muy costoso computacionalmente hablando. Nuevamente surgen las restricciones marcadas por el gran volumen de datos que se maneja, lo que obliga a elegir una técnica viable desde el punto de vista del tratamiento computacional. Dado que el objetivo es obtener clases léxicas se opta por valorar el Análisis de Semántica Latente, que es una técnica de reducción de dimensiones capaz de extraer conceptos y agrupar términos en torno a ellos.

### 3.3.1.2. Análisis de Semántica Latente

El Análisis de Semántica Latente (LSA, del inglés Latent Semantic Analysis) utiliza la representación de contextos de términos o palabras para determinar la similaridad de significados de palabras y conjuntos de palabras (Landauer, T. K., Foltz, P.W., Laham, 1998). Esta técnica, diseñada para aplicarse a grandes corpus textuales, revela conceptos implícitos en los datos, y permite representar relaciones semánticas entre documentos, palabras o palabras con documentos. En este caso documentos hace referencia al contexto que se considere para cada término o palabra.

En LSA se parte de la matriz de términos-documentos o matriz de contextos, en la que cada fila representa las frecuencias de aparición de un término en el contexto (documento)

correspondiente. Cada columna por tanto representa un documento o contexto, y los valores en la columna la frecuencia de aparición de cada uno de los términos contenidos en él. Los elementos que se tomen como contextos pueden variar: pueden ser concordancias extraídas para una palabra, frases, párrafos o documentos enteros (Lancia, 2007). En cualquier caso, serán considerados como documentos cuando se aplica el algoritmo y los términos hacen referencia a las palabras que se analizan. En la presente apuesta, los documentos corresponden a las líneas de concordancia simuladas por la aplicación desarrollada para recuperar la información que se ha registrado en la BD y los términos son las formas canónicas registradas en ella. De esta manera, dada una FC se puede reconstruir la matriz de contextos que se obtendría. Las combinaciones en la Tabla 3-7 sería el punto de partida si se quiere obtener los conceptos o grupos de palabras entre las bases del verbo *sentir*. La extracción de conceptos se realiza aplicando la descomposición en valores singulares (SVD, del inglés Singular Value Decomposition) de la matriz de contextos.

### 3.3.1.3. Descomposición SVD

La descomposición SVD se aplica a matrices rectangulares, dada una matriz  $A_{m \times n}$  con  $m > n$  se determinan las matrices  $S_1, S_2, \Sigma$  tales que:

$$A = S_2 \Sigma S_1^T$$

La matriz  $\Sigma$  es la matriz diagonal con los autovalores de  $A^T A$ , existe un índice  $r$  para el que las raíces cuadradas de los primeros  $r$  autovalores ordenados decrecientemente:  $\sigma_1, \sigma_2, \dots, \sigma_r$  son positivos, y el resto es cero, generalmente coincide con  $n$ . Este conjunto de valores se denominan **valores singulares**.  $S_2$  es la matriz de autovectores de  $A^T A$  y  $S_1$  la de autovectores de  $AA^T$ , cada entrada de  $A^T A$  es el número de palabras en común entre los contextos  $i, j$ . Cada entrada de  $AA^T$  es el número de contextos en los que aparecen conjuntamente las palabras  $i, j$ . La descomposición de  $A$  se denomina **descomposición en valores singulares**. Para reducir la dimensión, a  $k$  factores o conceptos ( $k < n$ ) se hacen 0 los correspondientes valores singulares, es decir, los autovalores  $\sigma_{k+1}, \sigma_{k+2}, \dots, \sigma_r$  se obtiene una aproximación de  $A$  utilizando la nueva versión de  $\Sigma$  para su reconstrucción (de Jorge Botana, G., 2010):

$$A \approx A_k = S_{2k} \Sigma_k S_{1k}^T$$

La elección de  $k$  se realiza empíricamente, estableciendo qué valores singulares:  $\sigma_r, \sigma_{r-1}, \dots, \sigma_{k+1}$  se consideran despreciables respecto a los restantes  $k$ . Al estar ordenados los valores singulares de mayor a menor, se considerarán las  $k$  primeras filas en la matriz. Estas  $k$

dimensiones seleccionadas representan los diferentes conceptos o grupos que se extraen de los datos.

En la nueva matriz reconstruida, cada entrada de coocurrencias observadas es sustituida por valores estimados que marcan las diferencias entre las clases de palabras, de forma que las puntuaciones que obtienen serán mayores o menores, incluso se infiere la relación de palabras con contextos que no son plausibles en la matriz original, pasando de tener frecuencia de coocurrencia 0 a valores mayores que 0.

Se ha aplicado LSA para clasificar las bases que combinan con el verbo *perpetrar*, que tiene un único significado, *palpitar* con 4 acepciones en el DRAE y el verbo *sentir*, que es polisémica con numerosas acepciones y todas ellas con distintos rangos de frecuencia en el corpus. Se ha utilizado el paquete *lsa* del software estadístico R (2012).

#### 3.3.1.4. Semántica Latente en las bases de *perpetrar*:

En el corpus se encuentran 23 sustantivos que combinan con el verbo *perpetrar*, la matriz de contextos es de dimensión  $23 \times 4046$ .

Los autovalores que se obtienen se listan a continuación, su representación gráfica ayuda a elegir la cantidad de dimensiones, que se traduce en los conceptos o clases léxicas que se extraen de las bases:

```
[1] 1367.100581  552.573877  274.523182  242.149528  209.576341  185.234976
[7]  179.162088  165.281073  164.306348  149.203851  144.902200  140.766047
[13]  138.140623  131.904095  114.926038  111.266192   88.563222   60.356442
[19]   37.929552   20.078836    4.319880    1.031975
```

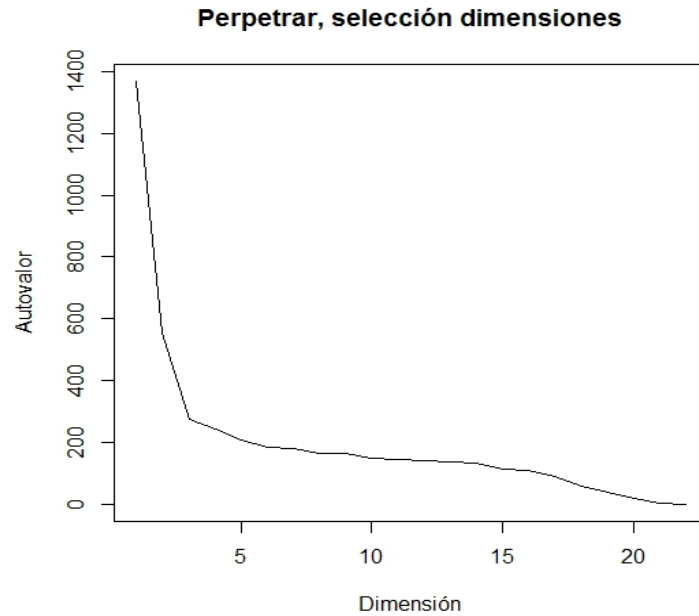


Gráfico 3-18 Autovalores en la descomposición SVD del verbo perpetrar

Se observa el peso diferenciado de la primera dimensión respecto a las restantes. En el Gráfico 3-19 se muestra la representación de los términos en el espacio semántico latente, utilizando la matriz:  $S_{2k}\Sigma_k$ , se seleccionó  $k=2$ . Se observa el grupo. *acto, asesinato, ataque* y el resto de sustantivos como *masacre, delito, robo*, etc. que se ven entremezclados con sustantivos de uso muy frecuentes como *hombre, mujer*.

#### 3.3.1.5. Semántica Latente en las bases de palpar:

Se repite el experimento sobre las bases del colocativo *palpar*, la construcción de la matriz de contextos consume un intervalo de tiempo que se mide en horas. Aún así se ha procesado la matriz de contextos generada por este verbo, los resultados se muestran en el Gráfico 3-20.

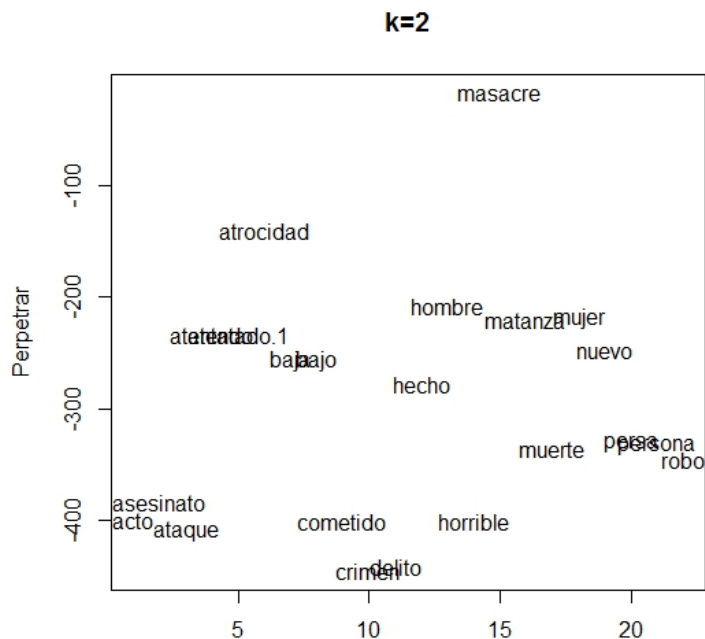


Gráfico 3-19 Representación de los términos en el espacio semántico latente del verbo perpetrar.

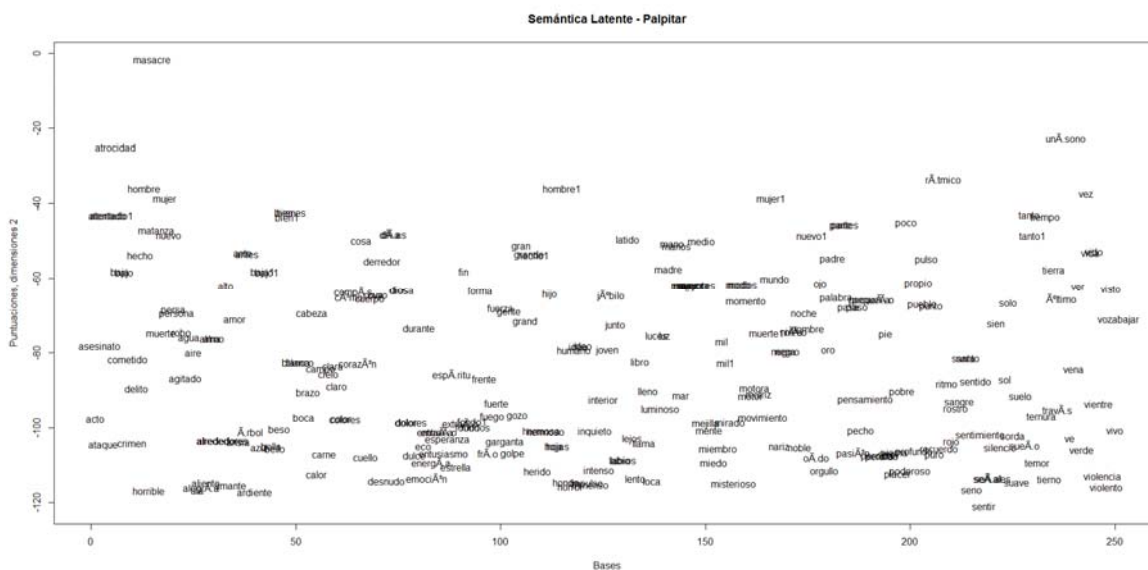


Gráfico 3-20 Representación de los términos en el espacio semántico latente del verbo palpitar.

Las formas canónicas cercanas en este gráfico reflejan similitud semántica (Jorge-Botana, G. *et al.*, 2010), (Olmos, R. *et al.*, 2014); para apreciar los resultados se ha realizado un zoom del Gráfico 3-20 que se presenta en 3 piezas separadas en zonas de izquierda a derecha del gráfico original: Gráfico 3-21, Gráfico 3-22 y Gráfico 3-23. Sobre ellos se han marcado conjuntos de bases cuya cercanía en el espacio semántico latente revela cierta similitud, o más bien relación semántica. La definición en el DRAE del verbo *palpitar* pone de manifiesto su vínculo con sustantivos relacionados con *corazón* y *afectos* o *pasiones*, también es lógico encontrar

fuertes vínculos con sustantivos *referidos a alguna parte del cuerpo*. Los grupos que se han marcado reflejan el mismo problema que se detectó en el análisis de *perpetrar*, bases de frecuencia muy alta como *hombre, mujer, cosa, tiempo, vez, ...* se entremezclan con elementos en los grupos en los que hay una motivación semántica en su configuración. Se han marcado 5 grupos:

G1: *ataque, crimen, asesinato, muerte, delito, matanza, hombre, mujer, hecho, ...*

G2: *boca, brazo, cuello, corazón, ...*

G3: *espíritu, esperanza, entusiasmo, emoción, fuerte, fuego, gozo, impulso, dolor, ...*

G4: *unísono, rítmico, tiempo, vez, ...*

G5: *pasión, orgullo, placer, ternura, violencia, profundo, puro, ...*

G3 se ajusta al grupo del DCECR:

SUSTANTIVOS QUE DENOTAN FUERZA O VITALIDAD.

G4 se ajusta al grupo del DCECR:

SUSTANTIVOS QUE DESIGNAN EMOCIONES Y SENTIMIENTOS PRINCIPALMENTE LOS DE NATURALEZA PASIONAL.

### **DRAE palpitar.**

(Del lat. *palpitāre*).

1. *intr.* Dicho del corazón: Contraerse y dilatarse alternativamente.
2. *intr.* Dicho del corazón: Aumentar su palpitación natural a causa de una emoción.
3. *intr.* Dicho de una parte del cuerpo: Moverse o agitarse interiormente con movimiento trémulo e involuntario.
4. *intr.* Dicho de algún afecto o pasión: Manifestarse vehementemente. *En sus gestos y palabras palpita el rencor.*

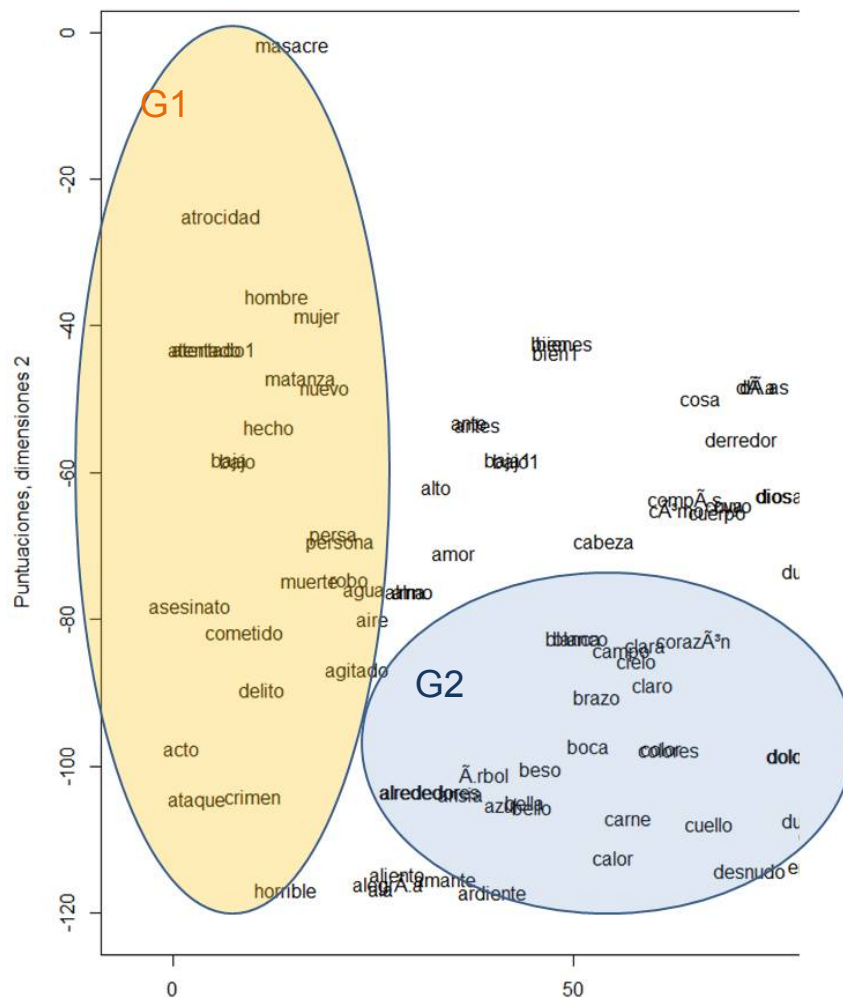


Gráfico 3-21 Zoom 1, resultados de palpar.

Semántica Latente - Palpitar

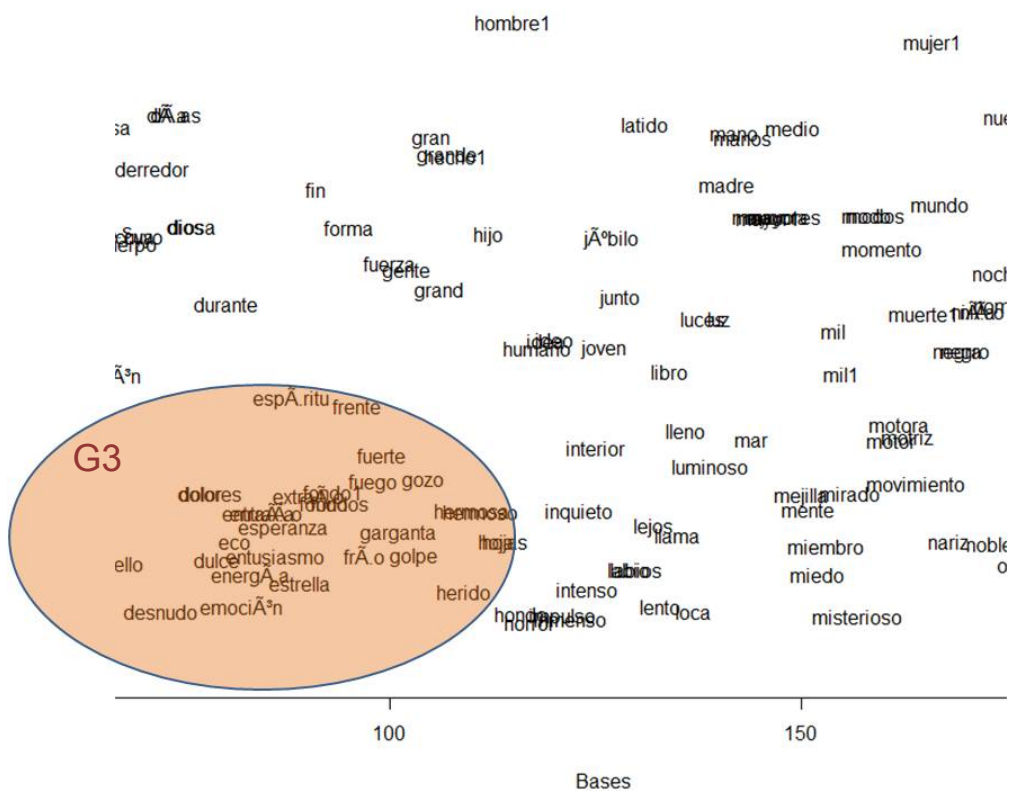


Gráfico 3-22 Zoom 2, resultados de palpitar



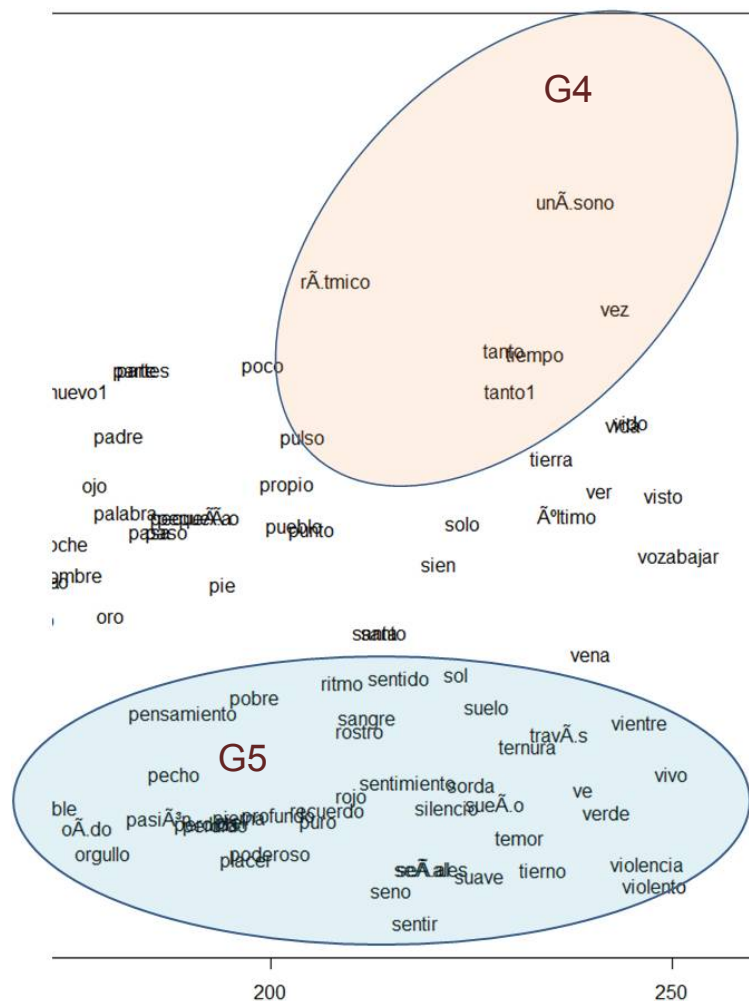


Gráfico 3-23 Zoom 3, resultados de palpar

Si bien en los grupos obtenidos, el método del análisis de la semántica latente resulta adecuado, hay varias cuestiones que no aconsejan aplicarlo sobre los registros extraídos del corpus. Por una parte, generar la matriz de datos que se requiere consume un tiempo inasumible en una aplicación de consulta de información de características combinatorias del español. Por otra parte, los resultados conllevan la supervisión manual para establecer la configuración de grupos.

### 3.3.2. Obtención de grupos semánticos a partir de la relación predicado-argumentos

La obtención de clases léxicas mediante las técnicas expuestas anteriormente requiere un tiempo de cómputo y recursos de memoria excesivos, además, se hace necesaria la toma de decisiones como el número de dimensiones o conceptos presentes, para lo que no hay un criterio claro. Abundando en los aspectos negativos, no quedan claramente marcados los grupos, por ejemplo, en el caso de *perpetrar* (Gráfico 3-19) se observan muy próximas las formas canónicas

*hombre-mujer, nuevo y matanza, o persona y robo*. También se suma a estas dificultades la necesidad de reconstruir en cada caso la matriz de contextos a partir de los datos registrados en la BD. Por esta razón se acomete la solución del problema a partir del modelo lingüístico de la **selección predicado-argumentos**, mediante las **preferencias de selección**.

En algunos trabajos de desambiguación semántica de palabras se utilizan medidas de las relaciones entre un predicado y sus argumentos. En Resnik (1999) se define como asociación de selección una medida entre verbos y sustantivos que se formaliza según la teoría de la información. Utiliza probabilidades que se estiman a partir de las coocurrencias de un predicado dado con los elementos de una determinada clase de sustantivos. Las clases están formadas por los elementos en synsets de *WordNet* que contienen hipónimos de un sustantivo o por synsets correspondientes a hiperónimos del mismo. También en Agirre (1994) se evalúa una medida similar, aunque establece la asociación entre clases de predicados con clases de argumentos.

El método de Resnik requiere disponer de un recurso léxico que provea de clases semánticas ya establecidas, para tal fin se propone utilizar algún Diccionario Ideológico del español como base de conocimiento léxico contrastada por los lingüistas.

### **3.3.3. Diccionarios Ideológicos, bases de conocimientos léxicas.**

Un diccionario ideológico agrupa las palabras en torno a una idea, según un orden conceptual. Representan, por tanto, una fuente de información sobre relaciones semánticas similares a los tesauros. En español se dispone del Diccionario Ideológico Vox (DIV) (Vox, 1995) y del Diccionario Ideológico de la Lengua Española de Julio Casares (DILE) (Casares, 1959). Ya que se propone el uso de este tipo de recursos en la extracción de preferencias de selección en sustitución de *WordNet* o su equivalente *EuroWordNet* en español, en los siguientes epígrafes se describe la estructura de ambos diccionarios.

#### **3.3.3.1. Estructura del Diccionario Ideológico Vox.**

El Diccionario Ideológico Vox se estructura en varios niveles diferentes (Figura 3-3):

*Cuadro General*, recoge los grandes grupos de ideas: Relaciones Generales, Materia, Naturaleza, Seres Vivos, Individuo y Modos de Vida.

*Parte Sinóptica*: Incluye los cuadros, corresponde a otro nivel de agrupación y representan identificadores de títulos de grupos conceptuales (Tabla 3-8).

*Parte Analógica*: Aparecen todas las palabras ordenadas en grupos según un criterio conceptual. Dentro de los grupos pueden aparecer subgrupos, motivados por una mayor afinidad entre sus componentes y también marcadores que representan relaciones entre

grupos y subgrupos, sin especificar el tipo. Los grupos a su vez se dividen en bloques por categoría gramatical de los elementos (Tabla 3-9). Se presentan las palabras como se asocian en la mente del hablante (Vox, 1 995).

<b>Cuadro</b>	<b>Grupo</b>
<i>dolor</i>	<i>cantar, doler, entortijarse, entrepunzar, latir, punzar, rabiar, suposición, traspasar, ver las estrellas, sentir</i>
<i>juicio</i>	<i>cantar, doler, entortijarse, entrepunzar, latir, punzar, rabiar, suposición, traspasar, ver las estrellas, sentir</i>
<i>juicio</i>	<i>conceptualizar, conceptuar, conocerse, considerar, contar, creer, dar por, decir, encasillar, entender, estimar, juzgar, opinar, pasar, reputar, sentir, tenerse, tomar por, ver, reconocerse</i>
<i>sentimiento</i>	<i>afección, corazón, entretelas, podredumbre, psicología, sentir, sentimiento</i>
<i>sentimiento</i>	<i>acoger, cobrar, coger, concebir, dar, morir, anidar</i>
<i>tristeza</i>	<i>deplorar, llorar, sentir</i>
<i>balbuceo</i>	<i>declamar, monologar, recitar, responder, sentir</i>

Tabla 3-8 Cuadros y grupos en los que aparece sentir

#### *Dolor*

*acritud, analgesia, clavo, dolor, puntada, punzada, ramalazo, rayo, sinalgia, tormento, yaya*

*agujetas, ciática, entuertos, ijada, lumbago, mastalgia, mialgia, ostealgia, retortijón*

*cálculos, cólico, cólico nefrítico, cólico renal, litiasis, mal de piedra*

*arenillas, bezoar, cálculo, piedra*

*ardor, comecome, comezón, escozor, fogaje, fuego, hormiguilla, hormiguillo, picazón, picor, piquiña, prurito, quemazón, rascazón*

*agudo, dolorido, doloroso, lacerante, lancinante, pedregoso, pruriginoso, terebrante, urente*

*indolente, indoloro*

*cantar, entortijarse, entrepunzar, latir, punzar, rabiar, sentir, traspasar, ver las estrellas*

*abrasar, comer, enardecerse, encender, escaldarse, escarabajear, irritar, picar, quemar*

Tabla 3-9 Grupos en el Cuadro dolor

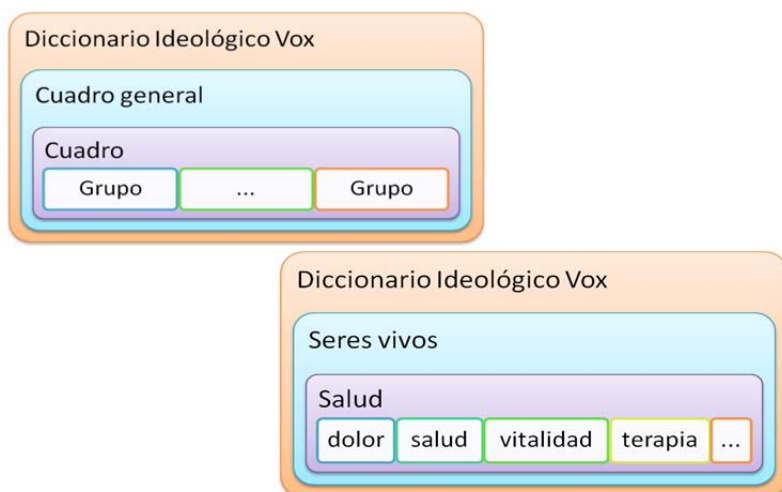


Figura 3-3 Estructura del Diccionario Ideológico VOX

### 3.3.3.2. Diccionario Ideológico Casares

En la parte sinóptica de este diccionario, se establecen los distintos niveles de agrupamiento que se emplean. En un primer nivel se tienen las grandes materias que estructuran la clasificación (Figura 3-4), donde se encuadran las 38 categorías básicas en las que se distribuyen las palabras; estas categorías se denominan cuadros y se subdividen a su vez en 2000 grupos de palabras conceptualmente relacionadas (Tabla 3-10).

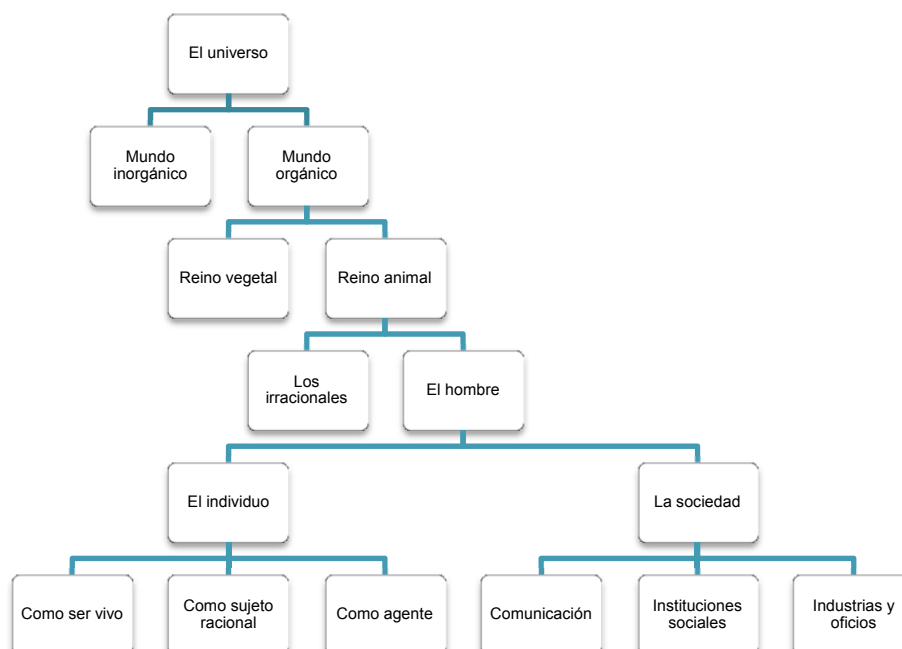


Figura 3-4 Grandes materias, Diccionario Ideológico Casares

<b>Cuadro</b>	<b>Cuadro</b>
<i>Religión</i>	<i>Colocación</i>
<i>Física y Química</i>	<i>Tiempo</i>
<i>Geografía, Astronomía, Meteorología</i>	<i>Cantidad</i>
<i>Geología, Mineralogía</i>	<i>Intelección</i>
<i>Botánica</i>	<i>Apreciación, juicio</i>
<i>Zoología</i>	<i>Voluntad</i>
<i>Anatomía</i>	<i>Conducta</i>
<i>Fisiología</i>	<i>Acción</i>
<i>Medicina</i>	<i>Lenguaje</i>
<i>Alimentación</i>	<i>Arte</i>
<i>Vestido</i>	<i>Estado, Nación</i>
<i>Vivienda</i>	<i>Costumbres</i>
<i>Sensibilidad, Sentidos</i>	<i>Derecho y Justicia</i>
<i>Sentimientos</i>	<i>Propiedad</i>
<i>Existencia, Cambio</i>	<i>Milicia</i>
<i>Relación, Orden, Causalidad</i>	<i>Comercio, Banca y Bolsa</i>
<i>Espacio</i>	<i>Agricultura, Zootecnia</i>
<i>Forma</i>	<i>Zootecnia</i>
<i>Movimiento</i>	<i>Transportes</i>

Tabla 3-10 Cuadros de la clasificación del DIC .

En la parte analógica se enumeran los grupos, divididos en subgrupos, y encabezados por un término, normalmente sustantivo, que corresponde con los que se detallan en los diferentes cuadros. Puede incluir remisiones a otros grupos relacionados conceptualmente a través de elementos marcados en negrilla para indicar que son otras cabeceras. También se diferencian subgrupos según la categoría gramatical de sus componentes: sustantivos, verbos, adjetivos, adverbios, modos adverbiales, preposiciones e interjecciones.

En la parte alfabética se recogen los significados de todas las palabras que aparecen en la parte analógica y se muestran referencias a los grupos a los que pertenece.

<b>Cuadro</b>	<b>Encabezado</b>	<b>Grupo</b>
<i>Conducta, lenguaje</i>	<i>ademán</i>	<i>accionar, amanerarse, chistar, cocar, enzainarse, expresar, gestear, gesticular, gesticular, guiñar, guiznar, hablar con los ojos, hablar por señas, hacer cocos, hacer del ojo, hacer figuras, hacer telégrafos, momear, pujar, remilgarse, sentir, timarse, torcer el gesto</i>
<i>Sentimientos</i>	<i>aflicción</i>	<i>ahogarse en poco agua, añorar, apurarse, atravesarse un nudo en la garganta, atribularse, caerse, cariñar, clamar a Dios, cubrirse el corazón, dar el pésame, deplorar, deprimirse, desolarse, dolerse, echar de menos, endecharse, engurrñarse, enmantarse, ensombrecerse, entristecer, lamentar, llorar lágrimas de sangre, llorar lágrimas de sangre, no caber el corazón en el pecho, nublársele el cielo, padecer, partirse el alma, pasar la pena negra, pasar las penas del purgatorio, penarse, poderse ahogar con un cabello, sentir, sentir de muerte</i>
<i>Sentimientos, conducta</i>	<i>arrepentimiento</i>	<i>acusar la conciencia, arrepentirse, compungir, compungirse, concomerse, corroer, deplorar, dolerse, escarabajar la conciencia, lamentar, llorar, llorar con lágrimas de sangre, morderse las manos, pesarle a uno, recomerse, sentir</i>
<i>Sensibilidad, sentidos</i>	<i>audición</i>	<i>abrir el oído, abrir los oídos, abrir tanto el oído, aguzar el sentido, aguzar las orejas, aplicar el oído, atender, auscultar, beber las palabras, dar oídos, dejarse de oír, entreoír, escuchar, estar colgado de los labios, estar pendiente de la boca, herir el oído, no perder ripio, oír, perderse, sentir, taparse los oídos, trasoír</i>
<i>Sentimientos, costumbres</i>	<i>compasión</i>	<i>ablandarse, adolecer, apiadarse, arrancársele a uno las entrañas, compadecer, compadecerse, compartir, compungirse, condolerse, conmoverse, contristarse, deplorar, dolerse, emblandecerse, lamentar, lastimarse, sentir</i>

Tabla 3-11 Cuadros y Encabezados de Grupo para el verbo sentir en D.I.C (I)

<b>Cuadro</b>	<b>Encabezado</b>	<b>Grupo</b>
<i>Intelección, conducta</i>	<i>conciencia</i>	<i>acusar la conciencia, advertir, ajustarse con la conciencia, apereibir, argüir la conciencia, caer en la cuenta, cargar la conciencia, comprender, conocerse, darse cuenta, encargar la conciencia, enconar, entenderse, escarabajear, escarbar, escrupulizar, percibir, reflexionar, remorder, sentir, tener conciencia de</i>
<i>Intelección</i>	<i>experiencia</i>	<i>advertir, aprender, ensayar, experimentar, foguear, notar, sentir, sufrir, tocar, ver</i>
<i>Apreciación- Juicio</i>	<i>juicio</i>	<i>apreciar, calificar, censurar, conocer de, conocerse, considerar, creer, criticar, dictaminar, discernir, echar el fallo, emitir, enjuiciar, entender en, estimar, existimar, fallar, hacer distinción, informar, juzgar, librar, opinar, parecerle a uno, reputar, sentir, someter, tener por, tomarle las medidas (a uno), valorar</i>
<i>Acción</i>	<i>pasividad</i>	<i>aceptar, aguantar, experimentar, expiar, incurrir en, lastar, padecer, pasar, pasar por las picas, penar, percibir, permitir, recibir, sentir, ser objeto de, soportar, sufrir, tener buenas espaldas, tolerar, tomar, tomarle a uno</i>
<i>Intelección, conducta</i>	<i>previsión</i>	<i>agorar, antever, anunciar, aprevenir, augurar, barruntar, darle el aire de, darle el corazón, decirle el corazón, haberse tragado, mirar las cosas con anteojo de aumento, mirar las cosas con anteojo de larga vista, ominar, preconocer, predecir, prenotar, presentir, presumir, prevenir, prever, pronosticar, remusgar, sentir, sospechar, tenerse tragado, ver, ver las cosas con anteojo de aumento, ver las cosas con anteojo de larga vista</i>
<i>Fisiología, Sensibilidad- sentidos, sentimientos</i>	<i>sensibilidad</i>	<i>advertir, apreciar, entrar en, experimentar, impresionarse, notar, observar, padecer, percibir, sentir, sufrir</i>
<i>Sentimientos</i>	<i>sentimiento</i>	<i>abrigar, cobrar, coger, concebir, experimentar, profesar, sentir, tener, tomar</i>

Tabla 3-12 Cuadros y Encabezados de Grupo para el verbo sentir en D.I.C (II)

### 3.3.4. Extracción de preferencias de selección

En este apartado se expone la medida de asociación de selección (Resnik, 1993), su adaptación a las Bases de Conocimiento sobre las que se trabaja y los resultados obtenidos.

#### 3.3.4.1. Asociación de selección

La medida de asociación de selección de Resnik cuantifica las restricciones de selección de un predicado respecto a una clase de argumentos. Este valor surge a partir de la entropía de una variable aleatoria que se utiliza como herramienta para determinar qué cantidad de incertidumbre se tiene sobre los argumentos que pudieran aparecer junto a un predicado dado, es decir, qué “ha sucedido”.

Este artificio se construye sobre la Teoría de la Información, que se fundamenta en la idea de que un suceso poco probable aporta más información que uno altamente probable. Por ejemplo, el verbo *palpitar* aporta más información sobre qué palabras se pueden esperar en su contexto que el verbo *dar*. La **cantidad de información** de un suceso dada por:

$$I(X = x) = \log_2 \frac{1}{p(x)}$$

La **entropía** de una variable aleatoria es la información media que aporta la misma:

$$H(X) = \sum_x p(x) \cdot \log_2 \frac{1}{p(x)}$$

La **entropía relativa**, dadas dos distribuciones de probabilidad  $p, q$ :

$$D(p \parallel q) = \sum_x p(x) \cdot \log_2 \frac{p(x)}{q(x)}$$

También es conocida como la distancia de Kullback-Leibler y verifica:

$$D(p \parallel q) \geq 0$$

$$D(p \parallel q) = 0 \Leftrightarrow p, q \text{ idénticas}$$

$$D(p \parallel q) = \sum_x p(x) \cdot \left[ \log_2 \frac{1}{q(x)} - \log_2 \frac{1}{p(x)} \right] = \sum_x p(x) \cdot [I_q(X = x) - I_p(X = x)]$$

es el promedio de diferencias de cantidad de información.



Se interpreta como el coste, en bits de información, que tiene el hecho de asumir que la distribución de probabilidad es  $q$  cuando en realidad es  $p$ .

Resnik recurre a la entropía relativa para evaluar si existen diferencias entre la distribución de probabilidades de una clase léxica de argumentos  $C$ , supuesto que ha sucedido el predicado  $p_i$ , frente a la probabilidad de la clase  $C$  sin tal supuesto, es decir:

$$D(p(C|p_i) \parallel p(C)) = \sum_C p(C|p_i) \cdot \log_2 \frac{p(C|p_i)}{p(C)}$$

Siendo:

$p(C)$  la distribución de la clase  $C$ , o **distribución a priori**

$p(C|p_i)$  la distribución de la clase  $C$ , conocido  $p_i$ , o **distribución a posteriori**.

Si un predicado restringe sus argumentos, selecciona una clase léxica, por tanto, las distribuciones a priori y posteriori deben diferir. Además, el rango de selección puede ser amplio o estrecho según la fuerza de la restricción que imponga a sus argumentos. En caso de rangos estrechos, la diferencia entre las dos distribuciones será grande, la distribución a priori debe ser muy diferente a la distribución a posteriori. Por este motivo se denomina **fuerza de la selección** a  $D(p(C|p_i) \parallel p(C))$ , y es el coste de asumir que  $C$  no tiene ninguna motivación para aparecer en el contexto de  $p_i$ .

Las clases  $C$  a las que hace referencia son synsets de hipónimos o hiperónimos del lexicon. Sin embargo, un sustantivo consta en un conjunto de synsets en la taxonomía, correspondientes a los diferentes sentidos.

La **asociación de selección**,  $A(p_i, C)$ , de un predicado,  $p_i$ , se define para cada clase,  $C$ , como la contribución de la clase en la fuerza de la selección, ponderada por el factor de escala respecto al total –normaliza el indicador y permite comparar los valores independientemente de la fuerza de la selección.

$$A(p_i, C) = \frac{p(C|p_i) \log_2 \frac{p(C|p_i)}{p(C)}}{\sum_C p(C|p_i) \cdot \log_2 \frac{p(C|p_i)}{p(C)}}$$

El numerador corresponde a cada sumando de  $D(p(C|p_i) \parallel p(C))$  y el denominador es dicho valor.

En la práctica, en lugar de trabajar con distribuciones de probabilidad, se emplean sus estimaciones obtenidas a partir de las frecuencias registradas en el corpus:

$$\hat{p}(C|p_i) = \frac{\sum_{w \in C} \frac{frec(p_i, w)}{frec(w)}}{frec(p_i)}$$

$$\hat{p}(C) = \frac{\sum_{w \in C} \frac{1}{|\mathbb{C}_w|}}{N}$$

siendo:  $p_i$  el predicado,  $C$  la clase,  $w$  una palabra,  $frec(p_i, w)$  la frecuencia conjunta en el corpus del predicado con la palabra  $w$ ,  $frec(p_i)$  la frecuencia del predicado en el corpus, y  $\mathbb{C}_w$  el conjunto de todas las clases en las que aparece  $w$  en la clasificación,  $N$  es el número de palabras en la clasificación.

### 3.3.4.2. La asociación de selección sobre el Diccionario Ideológico Vox

El Diccionario Ideológico Vox, DIV, proporciona las clases de palabras a partir de las que inducir grupos semánticos en los argumentos que seleccionan los predicados en el corpus. Dada la estructura de ambos diccionarios del mismo modo se podría utilizar el DILE, sin embargo, al disponer de la información de la categoría de las palabras del grupo se escoge el primero para realizar los cálculos. Las suposiciones que se hacen son las siguientes:

- Las clases que se utilizarán son los grupos en el nivel de la parte analógica.
- Se diferencian las tres estructuras:
  - *Verbo + Sustantivo*: El predicado tiene que ser un verbo.
  - *Sustantivo + Adjetivo*: El predicado tiene que ser un adjetivo.
  - *Verbo + Adverbio*: El predicado tiene que ser un adverbio.

En cada caso, se consideran solo los grupos del ideológico cuya categoría gramatical coincida con la de los argumentos.

Se ilustran los cálculos utilizando como predicado el verbo *sentir*, y el grupo de sustantivos:

$$G = \left\{ \begin{array}{l} \text{analgesia, sinalgia, yaya, acritud, rayo, clavo, puntada,} \\ \text{ramalazo, tormento, punzada} \end{array} \right\}$$

incluido en la cabecera dolor y que tiene elementos que aparecen en el corpus en el contexto del verbo *sentir*.

$$A(\text{sentir}, G) = \frac{\hat{p}(\text{sentir}|G) \cdot \log_2 \frac{\hat{p}(\text{sentir}|G)}{\hat{p}(G)}}{\sum_{G_{\text{Sust}_{\text{sentir}}}} \hat{p}(G_{\text{Sust}_{\text{sentir}}}) \cdot \log_2 \frac{\hat{p}(\text{sentir}|G_{\text{Sust}_{\text{sentir}}})}{\hat{p}(G_{\text{Sust}_{\text{sentir}}})}}$$

$$\hat{p}(\text{sentir}|G) = \frac{\sum_{\text{palabra} \in G} \frac{\text{frec}(\text{sentir}, \text{palabra})}{\text{frec}(\text{palabra})}}{\text{frec}(\text{sentir})} = \sum_{\text{palabra} \in G} \frac{\text{frecRel}_{\text{palabra}}(\text{sentir})}{\text{frec}(\text{sentir})}$$

$$\hat{p}(G) = \frac{\sum_{\text{palabra} \in G} \frac{1}{|\mathbb{C}_{\text{palabra}}|}}{N}$$

$|\mathbb{C}_{\text{palabra}}|$  es el cardinal del conjunto de todos los grupos en los que aparece *palabra*.

Los grupos en la cabecera *dolor* que aparecen en el contexto del verbo *sentir*, y el orden inducido por la asociación de selección respecto al total de grupos de sustantivos de los que hay muestras en el corpus en combinación con el verbo *sentir* se muestran en la Tabla 3-13.

Orden	Grupo	$A(\text{sentir}, G)$
2	<i>acritud, analgesia, clavo, dolor, puntada, punzada, ramalazo, rayo, sinalgia, tormento, yaya</i>	0.00084056
5	<i>ardor, comecome, comezón, escozor, fogaje, fuego, hormiguilla, hormiguillo, picazón, picor, piquiña, prurito, quemazón, rascazón</i>	0.00056183
49	<i>agujetas, cefalalgia, cefalea, ciática, entuertos, ijada, jaqueca, lumbago, mastalgia, mialgia, migraña, neuralgia, odontalgia, ostealgia, otalgia, retortijón, silbido de oídos, torticolis, tortícolis</i>	0.00025726
2640	<i>arenillas, bezoar, cálculo, piedra</i>	0.000019488

Tabla 3-13 Resultados para sentir + "los grupos de la cabecera dolor"

La Tabla 3-14 recoge los 25 grupos con mayor asociación de selección cuando ésta se evalúa en grupos de sustantivos en el DIV. Se incluyen todos los elementos del grupo que han sido registrados en el corpus –se resaltan en negrilla los que previamente se catalogaron como colocaciones, mostrándose también las cabeceras. Nótese que se pueden obtener diferentes grupos de una misma cabecera con desiguales valores de asociación de selección. Se muestra la capacidad de la asociación de selección para detectar grupos de argumentos que no incluyen elementos identificados como colocaciones –a pesar de que no existe una marcada preferencia individual, el conjunto de combinaciones aporta la suficiente fuerza de selección para hacer destacar al grupo.

<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(sentir, G)</b>
<i>agitación, alarma, ansia, ansiedad, batalla, combate, comezón, conflicto, desasosiego, guerra, hormigueo, inquietud, intranquilidad, malestar, pesadilla, preocupación, taco, zozobra</i>	<i>intranquilidad</i>	<i>0,00103141</i>
<i>ardor, comezón, escozor, fuego, picazón, picor, prurito, quemazón</i>	<i>dolor</i>	<i>0,00084056</i>
<i>mojada, picadura, pinchazo, puntada, punzada</i>	<i>herida</i>	<i>0,00061050</i>
<i>aflicción, agonía, ahogo, amargura, angustia, apuro, congoja, desconsuelo, dolor, duelo, hiel, pasión, pena, pesar, quebranto, sentimiento, sufrimiento</i>	<i>sufrimiento</i>	<i>0,00056807</i>
<i>clavo, dolor, puntada, punzada, ramalazo, rayo, tormento, yaya</i>	<i>dolor</i>	<i>0,00056182</i>
<i>arcada, asco, basca, escrúpulo, fatiga, náusea, repugnancia</i>	<i>indigestión</i>	<i>0,00053228</i>
<i>desazón, disgusto, pesadumbre, punzada, puñalada, sinsabor, tribulación</i>	<i>sufrimiento</i>	<i>0,00052641</i>
<i>enfado, enojo, escozor, indignación, mosca, ofensa, picazón</i>	<i>enojo</i>	<i>0,0005178</i>
<i>afán, agonía, anhelo, ansia, apetito, deseo, empeño, emulación, envidia, gana, hambre, pelota, pío, prurito, sed, voto</i>	<i>deseo</i>	<i>0,00050270</i>
<i>azote, calamidad, catástrofe, desastre, desdicha, desgracia, desventura, drama, fatalidad, infortunio, mal, miseria, ramalazo, rayo, suceso, suspenso, tormenta, tragedia, trago, través</i>	<i>revés</i>	<i>0,00047396</i>
<i>asco, horror, náusea, repugnancia, repulsión</i>	<i>desagrado</i>	<i>0,00047321</i>
<i>disgusto, enojo, fastidio, gravedad, incomodidad, incomodo, mareo, molestia</i>	<i>desagrado</i>	<i>0,00045289</i>
<i>agujón, aguijonazo, espuela, pinchazo</i>	<i>inducción</i>	<i>0,00044326</i>
<i>aberración, delirio, desatino, desvarío, devaneo, enajenación, locura, luna, ramalazo, vértigo</i>	<i>demencia</i>	<i>0,00044148</i>

Tabla 3-14 Los 25 grupos con mayor asociación de selección para el verbo sentir (I)

<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(sentir, G)</b>
<i>alborozo, <b>alegría</b>, dicha, euforia, felicidad, gozo, <b>humor</b>, júbilo, <b>placer</b>, regocijo</i>	<i>alegría</i>	<i>0,00042495</i>
<i>aprensión, cuidado, desconfianza, escrúpulo, espina, inseguridad, malicia, <b>miedo</b>, <b>pensamiento</b>, prejuicio, recelo, sospecha, <b>temor</b></i>	<i>recelo</i>	<i>0,00041306</i>
<i>aguijonazo, dardo, picante, quemazón</i>	<i>ironía</i>	<i>0,00040895</i>
<i>añoranza, depresión, melancolía, nostalgia, soledad, <b>vacío</b></i>	<i>tristeza</i>	<i>0,00038088</i>
<i>animadversión, animosidad, antipatía, aversión, discordia, enemistad, rivalidad</i>	<i>enemistad</i>	<i>0,00036861</i>
<i>abatimiento, aturdimiento, hormigueo</i>	<i>enfermedad</i>	<i>0,00036835</i>
<i>agrado, deleite, delicia, dulzura, fruición, gloria, goce, gozo, gusto, <b>placer</b>, regalo, satisfacción</i>	<i>placer</i>	<i>0,00036707</i>
<i>aleteo, contracción, <b>golpe</b>, latido, palpitación, pulsación, pulso, salto</i>	<i>agitación</i>	<i>0,00036645</i>
<i>cicatriz, deje, <b>efecto</b>, excitación, huella, impaciencia, impacto, <b>impresión</b>, regusto, sabor, <b>sensación</b></i>	<i>excitación</i>	<i>0,00036569</i>
<i>afección, <b>afecto</b>, <b>amor</b>, apego, cariño, debilidad, querencia, querer, ternura</i>	<i>amor</i>	<i>0,00036494</i>
<i>efusión, <b>emoción</b>, entusiasmo, exaltación, expansión, inspiración, intensidad, lirismo, <b>pasión</b>, tensión, vehemencia</i>	<i>exaltación</i>	<i>0,00036076</i>

Tabla 3-15 Los 25 grupos con mayor asociación de selección para el verbo sentir (II)

La representación gráfica de los resultados que se registran para cada uno de los ejemplos tratados en este capítulo sigue un mismo patrón. Para la mayor parte de los grupos que contienen alguna de las bases con las que combina el colocativo, se obtiene una puntuación casi nula; sin embargo, en todos los casos se observa un conjunto de grupos del DIV que presentan un valor alto de asociación de selección.

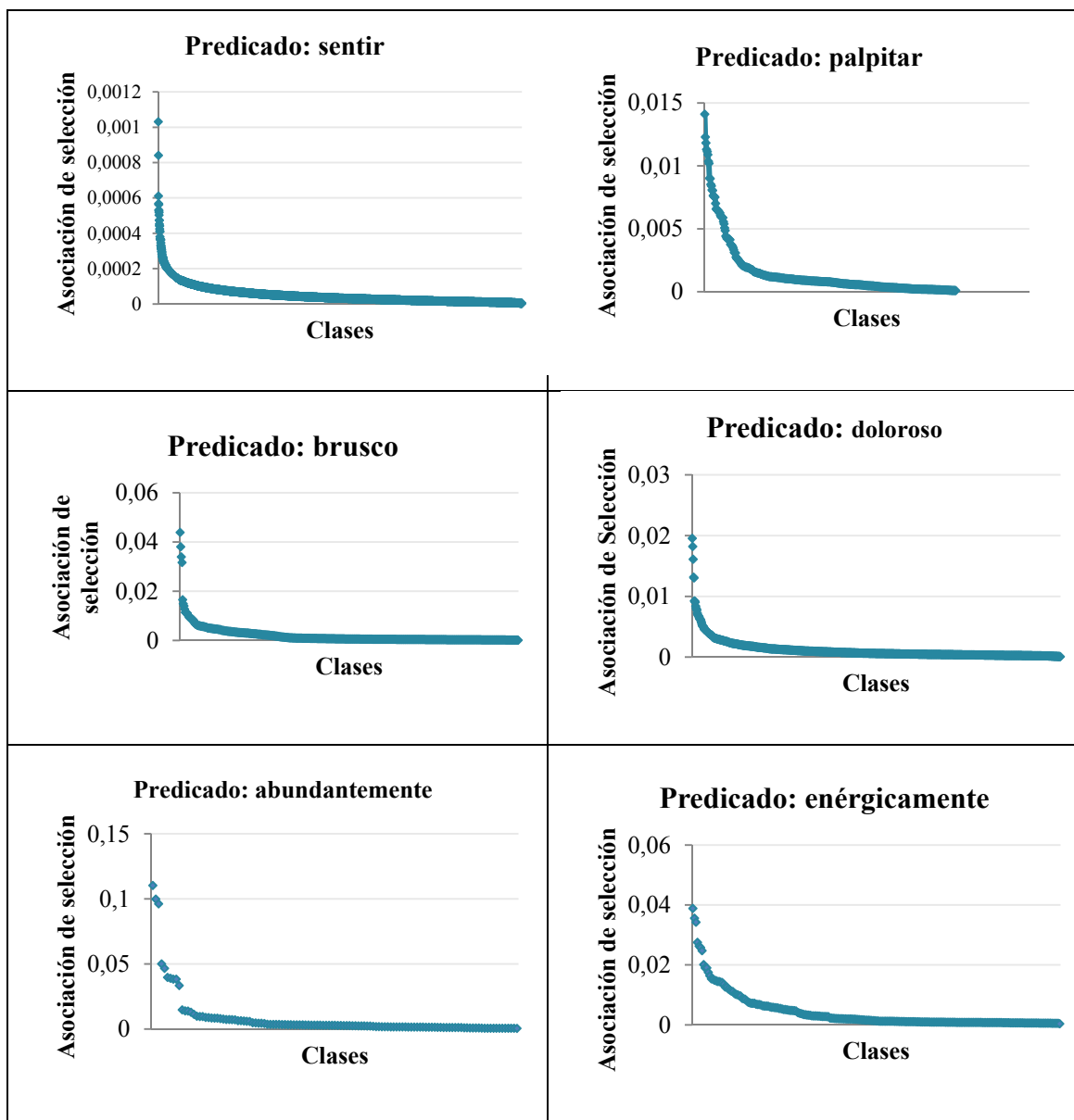


Gráfico 3-24 Asociación de Selección para los ejemplos tratados.

Los elementos pertenecientes a los grupos mejor puntuados para las restantes formas canónicas que han servido de ejemplo en este capítulo se enumeran en el Anexo 6.

### 3.4. Evaluación de la asociación de selección en el DCECR.

Los resultados obtenidos en los ejemplos anteriores muestran el buen comportamiento de la asociación de selección. Como indicio de la idoneidad del método se aportan las clases léxicas que describen las capacidades de selección de argumentos del verbo *sentir* descritas en DCECR. (Tabla 3-2) junto con las clases léxicas del DIV con mejores puntuaciones, se han resaltado las coincidencias en los grupos. En el primer caso se puede considerar que el grupo

proporcionado por el DCECR es prácticamente un subconjunto del grupo del DIV (Tabla 3-16), salvo un elemento todos están incluidos en él. En el segundo caso, Tabla 3-17, se recurre a las definiciones de las palabras en cada grupo para verificar que sus elementos responden al rasgo semántico *tristeza* que define al grupo del DCECR. Todos los elementos del grupo tienen una acepción con significado de *tristeza* o *aflicción*.

---

**DCECR.: SUSTANTIVOS QUE DENOTAN PREOCUPACIÓN, INCERTIDUMBRE O RESERVA = {inquietud, miedo, preocupación, reparo, temor, terror, zozobra}**

---

**DIV** agobio, alarma, ansia, comecome, comezón, conflicto, desabrimiento, hormiguelo, **inquietud**, malestar, mareta, **miedo**, pesadilla, **preocupación**, rebato, reparo, suspense, **temor**, **terror**, torozón, **zozobra**

---

Tabla 3-16 Clases léxicas para sentir en el DCECR. y DIV.

---

**DCECR SUSTANTIVOS QUE DENOTAN TRISTEZA O CONMISERACIÓN = {compasión, conmisericordia, lástima, pena, tristeza}**

---

**DIV** aflicción, agonía, angustia, congoja, desconsuelo, dolor, duelo, **pena**, pesar, sentimiento

---

Tabla 3-17 Clases léxicas para sentir en el DCECR y D.I.V.

**aflicción** (DEA)

1. Pesadumbre o tristeza

**agonía** (DRAE)

1. f. Angustia y congoja del moribundo; estado que precede a la muerte.
2. f. Pena o aflicción extremada.
3. f. Angustia o congoja provocadas por conflictos espirituales.

**angustia.** (DRAE)

(Del lat. *angustia* 'angostura', 'dificultad').

1. f. Aflicción, congoja, ansiedad.
5. f. Dolor o sufrimiento.

**angustia** (DEA)

1. Aflicción o congoja

**congoja** (DRAE)

1. f. Desmayo, fatiga, angustia y aflicción del ánimo.

**desconsuelo.** (DRAE)

1. m. Angustia y aflicción profunda por falta de consuelo.

**desconsuelo** (DAE)

- 1 Pena muy profunda y que se siente como insuperable

**dolor.** (DRAE)

(Del lat. dolor, *-ōris*).

2. m. Sentimiento de pena y congoja.

**duelo** (DAE)

1. Pena o dolor por la muerte de alguien
2. En gral: Pena o dolor

**pena** (DAE)

1. Tristeza (estado de ánimo)
2. Compasión o lástima

**pesar**<sup>2</sup>. (DRAE)

1. m. Sentimiento o dolor interior que molesta y fatiga el ánimo.

**sentimiento** (DUE)

1. («de») m. Estado afectivo de la clase que se expresa: 'Un sentimiento de angustia, de soledad, de abandono, de insatisfacción'.
2. Estado de ánimo de los que consisten en sentir atracción o aversión por una persona, o en sentir \*alegría o tristeza.

Utilizando como medida de similitud entre los grupos del DCECR y el DIV el coeficiente de Jaccard:

$$similitud(G_i, G_j) = \frac{|G_i \cap G_j|}{|G_i \cup G_j|}$$

Se obtienen los grupos con  $similitud(G_i, G_j) > 0$  ordenados por este valor, se consideran grupos similares los que presentan cumplen  $similitud(G_i, G_j) > 0,1$ . De los 3685 grupos del DIV que aparecen con el verbo *sentir* en el corpus (Gráfico 3-25), para cada caso se encuentran entre 1 y 3 grupos similares en el DIV según este criterio. En todos los casos alguno de ellos está entre los 25 con mayor asociación de selección cuando se toma como espacio semántico el DIV.



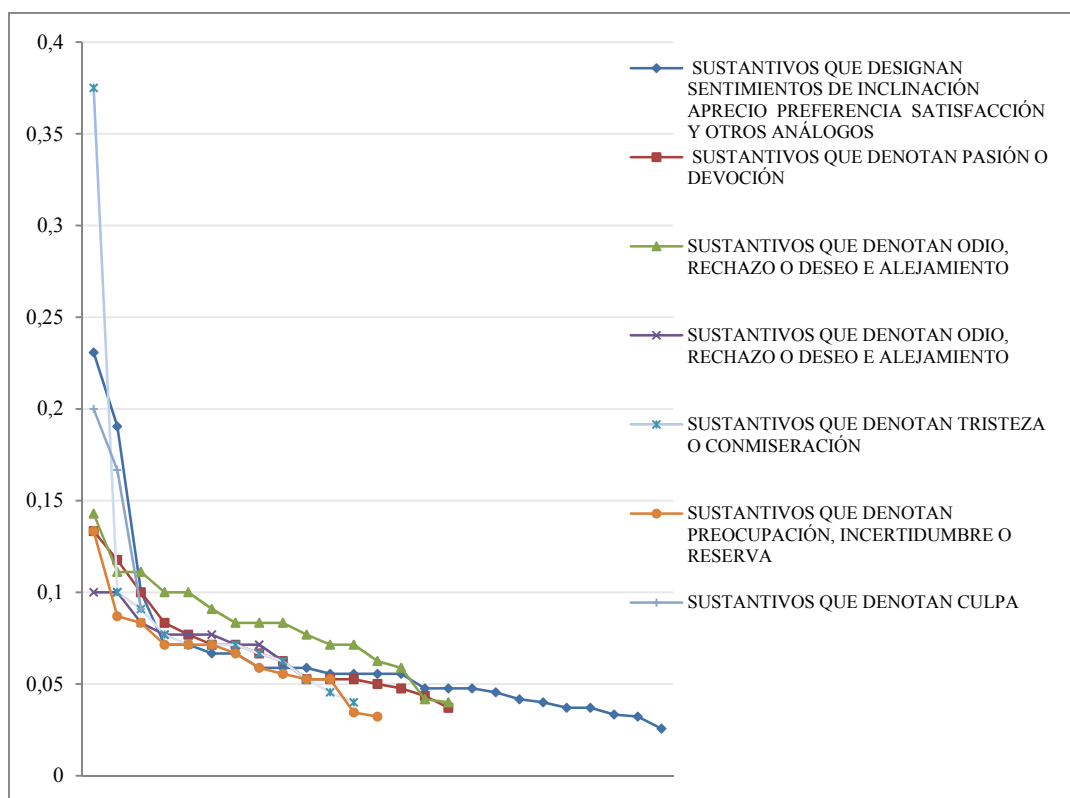


Gráfico 3-25 Grupos en el DIV similares a los grupos para el verbo sentir.

Con objeto de ratificar la validez de este indicador de forma objetiva, se realiza el experimento consistente en evaluarlo sobre el espacio semántico definido por las clases léxicas del DCECR. De este modo, se espera obtener valores altos de la asociación de selección entre un determinado predicado y las clases léxicas que selecciona como argumentos según el DCECR. Se ilustra la discusión con los ejemplos utilizados a lo largo del capítulo que tienen entrada en el citado diccionario: *sentir*, *palpitar*, *brusco*, *abundantemente* y *enérgicamente*.

En el Gráfico 3-26 se representan las puntuaciones que obtienen los grupos del DCECR cuando se utiliza éste como espacio semántico de interés. Se resaltan en las clases semánticas que combinan con el verbo *sentir*—se observa que la mayoría se posiciona entre las más altas.

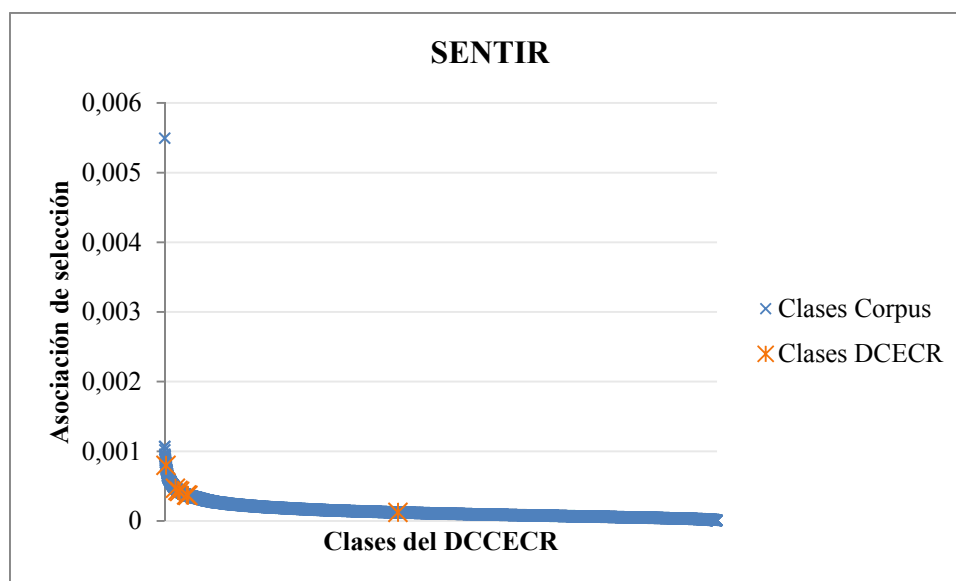


Gráfico 3-26 Asociación de Selección sobre el DCECR. para el verbo sentir

Además, otras clases con puntuaciones elevadas corresponden a grupos que hacen referencia a sustantivos que comparten rasgos semánticos con los enumerados para el verbo *sentir*. Se presentan los mejores casos en la Tabla 3-18 y Tabla 3-19 ordenados de mayor a menor asociación de selección. Se ha calculado el coeficiente de Jaccard entre éstas y las del verbo sentir en el DCECR. Se ha encontrado que en 6 de los casos se cumple el criterio establecido para considerarlas similares, indicándose los casos que lo verifican entre paréntesis. Los grupos en este espacio semántico están motivados por las capacidades combinatorias de los predicados que los determinan, y no por una clasificación semántica de las palabras del español, hecho que determina su solapamiento. De esta forma, se explica el resultado obtenido, clases léxicas seleccionadas por el verbo *sentir* entre las mejores y con una alta relación con las mejor posicionadas en el ranking.

---

**Clase DCECR. en el corpus**


---

<p>SUSTANTIVOS QUE DESIGNAN DIVERSOS SENTIMIENTOS SENSACIONES Y ESTADOS DE ÁNIMO DE SIGNO NEGATIVO EN ESPECIAL LOS REFERIDOS A LA EXCITACIÓN NERVIOSA EL ENOJO LA AVERSIÓN Y EL DESASOSIEGO TAMBIÉN CON OTROS QUE EXPRESAN ALGUNAS DE SUS MANIFESTACIONES</p>	<p><i>nervios, ánimo, ira, cólera, inquietud, excitación, angustia, descontento, disgusto, cabreo, odio, rencor, sollozo, sufrimiento, estrés, enfado, rabia, ansiedad, crispación, enojo, indignación</i></p>
<p>SUSTANTIVOS QUE DESIGNAN DIVERSAS AFLICCIONES DEL ÁNIMO MÁS FRECUENTEMENTE SI SE RELACIONAN CON EL DISGUSTO LA PESADUMBRE O LA AVERSIÓN HACIA ALGO (IV)</p>	<p><i>pesar, dolor, sufrimiento, rencor, rechazo, odio, pesimismo, desinterés, disgusto, malestar, desagrado, remordimiento, resentimiento</i></p>
<p>SUSTANTIVOS QUE DENOTAN AVERSIÓN RENCOR RECHAZO Y OTROS SENTIMIENTOS HOSTILES EN DIVERSOS GRADOS (IV)</p>	<p><i>odio, enemistad, hostilidad, antipatía, animadversión, rechazo, inquina, rivalidad, resentimiento, rencor, aversión, desprecio, aborrecimiento, repudio</i></p>
<p>SUSTANTIVOS QUE DENOTAN ANSIEDAD TEMOR Y OTRAS FORMAS DE DESEQUILIBRIO EMOCIONAL (VI)</p>	<p><i>ansiedad, nerviosismo, escalofrío, estupor, tensión, temor, miedo, irritabilidad</i></p>
<p>SUSTANTIVOS QUE DESIGNAN SENTIMIENTOS O ESTADOS DE ÁNIMO DE NATURALEZA NEGATIVA MÁS FRECUENTEMENTE SI MANIFIESTAN AGITACIÓN IRRITACIÓN Y OTRAS FORMAS DE INESTABILIDAD</p>	<p><i>nerviosismo, emoción, angustia, furia, celos, irritación, rabia, nervios, ira, ansiedad, desesperación, excitación, histeria, pánico, miedo</i></p>
<p>SUSTANTIVOS QUE DESIGNAN EL DOLOR PREOCUPACIÓN Y OTROS ESTADOS AFLICTIVOS (V)</p>	<p>LA <i>preocupación, dolor, consternación, depresión, desilusión, disgusto, malestar, pena, tristeza, angustia, cansancio, crispación, desaliento</i></p>

---

Tabla 3-18 Clases léxicas del DCECR. con mayor asociación de selección para sentir (I).

---

**Clase DCECR. en el corpus**


---

<p>SUSTANTIVOS QUE DESIGNAN OTRAS ALTERACIONES PSÍQUICAS ASÍ COMO ALGUNAS AFECCIONES LIGADAS A ELLAS MÁS FRECUENTEMENTE SI SE RELACIONAN CON LA ANSIEDAD LA INSATISFACCIÓN O EL DESABRIMIENTO</p>	<p><i>desazón, tensión, presión, malestar, preocupación, remordimiento, soledad, pena, amargura, neura, estrés, angustia</i></p>
<p>SUSTANTIVOS QUE DENOTAN AVERSIÓN REACCIÓN HOSTIL Y OTRAS FORMAS DE RECHAZO TAMBIÉN CON OTROS QUE DESIGNAN LOS SENTIMIENTOS QUE SE SUELEN ASOCIAR CON TALES REACCIONES (IV)</p>	<p><i>enemistad, antipatía, desprecio, odio, hostilidad, desconfianza, animadversión, desafecto, recelo, rechazo, oposición, envidia, ira, anatema, ojeriza</i></p>
<p>SUSTANTIVOS QUE DENOTAN MALESTAR, DISGUSTO Y OTRAS FORMAS DE PESADUMBRE, MÁS FRECUENTEMENTE LAS QUE DESIGNAN EL EFECTO ANÍMICO QUE PROVOCA LO QUE SE MALOGRA. TAMBIÉN CON OTROS QUE DESIGNAN EL QUE PRODUCE LA PÉRDIDA DE LAS FUERZAS O LA VOLUNTAD</p>	<p><i>frustración, decepción, desánimo, desgana, amargura, cansancio, debilidad, agotamiento, inapetencia, aburrimiento, malestar</i></p>
<p>SUSTANTIVOS QUE DENOTAN INCLINACIÓN DEL ÁNIMO A MENUDO INTENSA Y FAVORABLE HACIA ALGUNA COSA (I)</p>	<p><i>predilección, fascinación, pasión, afición, simpatía, amor, cariño, adicción, debilidad, parcialidad, admiración, respeto, identificación</i></p>

---

Tabla 3-19 Clases léxicas del DCECR. con mayor asociación de selección para sentir (II).

El análisis de los resultados del verbo *palpitar* indica un comportamiento que se ajusta al descrito para el verbo *sentir*, a pesar de que en el gráfico aparecen numerosas clases léxicas mejor puntuadas que las propuestas para este verbo en el diccionario (Gráfico 3-27). Por ejemplo, la clase léxica en el catálogo de *palpitar*:

SUSTANTIVOS QUE DESIGNAN EMOCIONES Y SENTIMIENTOS, PRINCIPALMENTE LOS DE NATURALEZA PASIONAL

que la forman los sustantivos: *amor, deseo, emoción, pasión* y *pulsión*, tiene correspondencia con la clase:

SUSTANTIVOS QUE DENOTAN ENERGÍA, AFÁN O EMPEÑO PUESTOS EN ALGUNA COSA, ASÍ COMO OTRAS CUALIDADES Y ACTITUDES

NECESARIAS PARA LLEVAR A CABO LO QUE SE CONSIDERA DIFÍCIL O  
ARRIESGADO

Entre sus elementos se encuentran los sustantivos: *corazón, arrojo, gana*. Por una parte, *deseo* y *pasión* se encuentran entre los sustantivos sinónimos de *afán*, *deseo* también es sinónimo de *empeño*. Por otro lado, *corazón* es sinónimo de: *amor* y *sentimiento*, *gana* lo es de: *deseo, anhelo* y *arrojo* es sinónimo de: *admiración, pasión, arrebató*. Algo similar sucede entre SUSTANTIVOS QUE DENOTAN FUERZA O VITALIDAD formada por *fuerza, vitalidad, vida* y

SUSTANTIVOS QUE DESIGNAN DIVERSOS RASGOS SOBRESALIENTES DE  
LAS PERSONAS O LAS COSA, MÁS FRECUENTEMENTE SI SE RELACIONAN  
CON EL ÍMPETU, LA ENERGÍA Y OTRAS FORMAS EN QUE PUEDE  
MANIFESTARSE LA CAPACIDAD DE MOVERSE, PORFIAR, IMPULSAR O  
ATRAER.

formado por los sustantivos: *atracción, corazón, energía, estilo, fuerza, impulso, obstinación, pasión, personalidad, vena populista*, entre los que se halla el sustantivo *fuerza* y el sustantivo *energía* que figura como sinónimo de *vitalidad* y de *vida*.

Del mismo modo, existe cierto solapamiento entre:

SUSTANTIVOS QUE DESIGNAN ARTES, TÉCNICAS Y TENDENCIAS  
ARTÍSTICAS, ASÍ COMO ALGUNAS DE SUS MANIFESTACIONES

conformado por: *poesía, literatura, melodía, obra, creación, cubismo, flamenco, partitura, acorde, tendencia*

y

SUSTANTIVOS QUE DESIGNAN MANIFESTACIONES MUSICALES, TAMBIÉN  
CON OTROS QUE EXPRESAN LA REPETICIÓN REGULAR DEL SONIDO

*ritmo, melodía, canción, música, compás, estribillo, pop, rap, latido*. Se entienden las manifestaciones musicales como un caso de manifestaciones artísticas, además los sustantivos del primer caso: *acorde, melodía, partitura* encajan perfectamente en este segundo grupo.

## PALPITAR

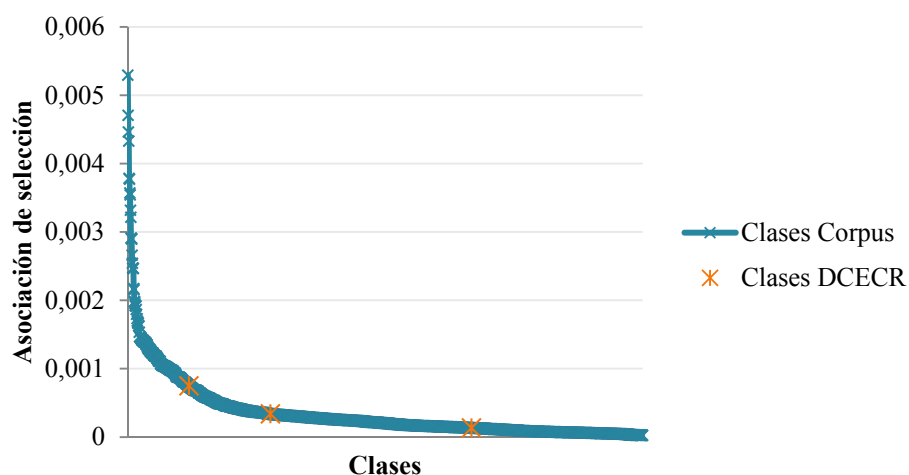


Gráfico 3-27 Asociación de selección sobre el DCECR. para el verbo palpitar.

---

**Clase DCECR. en el corpus**


---

SUSTANTIVOS QUE DESIGNAN MANIFESTACIONES MUSICALES, TAMBIÉN CON OTROS QUE EXPRESAN LA REPETICIÓN REGULAR DEL SONIDO

*compás, melodía, latido,...*

SUSTANTIVOS QUE DESIGNAN LAS PARTES INTERIORES O ESENCIALES DE ALGO O DE ALGUIEN

*corazón, entraña, interior,...*

SUSTANTIVOS QUE DENOTAN ENERGÍA, AFÁN O EMPEÑO PUESTOS EN ALGUNA COSA, ASÍ COMO OTRAS CUALIDADES Y ACTITUDES NECESARIAS PARA LLEVAR A CABO LO QUE SE CONSIDERA DIFÍCIL O ARRIESGADO

*corazón, emoción,...*

EL SUSTANTIVO RITMO Y CON OTROS QUE DESIGNAN ACCIONES O MOVIMIENTOS, ESPECIALMENTE SI SON REGULARES

*latido, ritmo,...*

---

Tabla 3-20 Clases léxicas del DCECR. con mayor asociación de selección para palpitar.

---

**Clase DCECR. en el corpus**


---

SUSTANTIVOS QUE DESIGNAN DIVERSOS RASGOS SOBRESALIENTES DE LAS PERSONAS O LAS COSA, MÁS FRECUENTEMENTE SI SE RELACIONAN CON EL ÍMPETU, LA ENERGÍA Y OTRAS FORMAS EN QUE PUEDE MANIFESTARSE LA CAPACIDAD DE MOVERSE, PORFIAR, IMPULSAR O ATRAER	<i>energía, fuerza, pasión...</i>
SUSTANTIVOS QUE DESIGNAN CAPACIDADES Y ATRIBUTOS DE LAS PERSONAS RELATIVOS A SU FORMA DE SER, DE PENSAR O DE ACTUAR MENTAL O EMOCIONALMENTE, TAMBIÉN CON OTROS QUE DESIGNAN ALGUNOS ÓRGANOS EN LOS QUE SUPUESTAMENTE RADICAN	<i>mente, corazón,...</i>
SUSTANTIVOS QUE DESIGNAN PARTES, ELEMENTOS, ASPECTOS O COMPONENTES DEL SER HUMANO QUE REPRESENTAN SENTIMIENTOS, SENSACIONES O CAPACIDADES	<i>alma, espíritu,...</i>

---

Tabla 3-21 Clases léxicas del DCECR. con mayor asociación de selección para palpar.

---

**Clase DCECR. en el corpus**


---

SUSTANTIVOS QUE DESIGNAN DIVERSOS RASGOS SOBRESALIENTES DE LAS PERSONAS O LAS COSA, MÁS FRECUENTEMENTE SI SE RELACIONAN CON EL ÍMPETU, LA ENERGÍA Y OTRAS FORMAS EN QUE PUEDE MANIFESTARSE LA CAPACIDAD DE MOVERSE, PORFIAR, IMPULSAR O ATRAER	<i>energía, fuerza, pasión...</i>
SUSTANTIVOS QUE DESIGNAN CAPACIDADES Y ATRIBUTOS DE LAS PERSONAS RELATIVOS A SU FORMA DE SER, DE PENSAR O DE ACTUAR MENTAL O EMOCIONALMENTE, TAMBIÉN CON OTROS QUE DESIGNAN ALGUNOS ÓRGANOS EN LOS QUE SUPUESTAMENTE RADICAN	<i>mente, corazón,...</i>
SUSTANTIVOS QUE DESIGNAN PARTES, ELEMENTOS, ASPECTOS O COMPONENTES DEL SER HUMANO QUE REPRESENTAN SENTIMIENTOS, SENSACIONES O CAPACIDADES	<i>alma, espíritu,...</i>

---

Tabla 3-22 Clases léxicas del DCECR con mayor asociación de selección para palpar.

El resto de ejemplos ratifica los resultados obtenidos en los casos *sustantivo + verbo*. El adjetivo *brusco* obtiene las máximas puntuaciones en las clases léxicas que selecciona según el DCECR. Si el análisis se centra en los adverbios *abundantemente* y *enérgicamente*, se observa que la mayor parte de las clases bien puntuadas tienen un alto grado de solapamiento y que tienen correspondencia con alguna de las aportadas para tales adverbios por el DCECR como clases léxicas que seleccionan como argumentos. Para el adverbio *enérgicamente* destaca la

gran variedad de clases similares que se aportan, como se puede verificar en la Tabla 3-25 y la Tabla 3-26

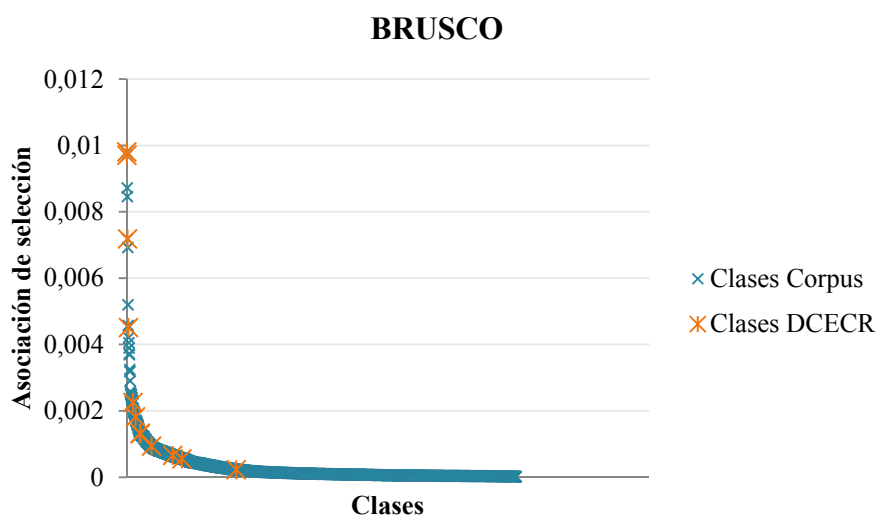


Gráfico 3-28 Asociación de selección sobre el DCECR. para el adjetivo brusco.

---

**.Clase DCECR. en el corpus**

---

SUSTANTIVOS QUE DENOTAN MOVIMIENTO, Y A MENUDO CAMBIO RÁPIDO DE DIRECCIÓN	<i>maniobra, movimiento, viraje,...</i>
SUSTANTIVOS QUE DENOTAN DETENCIÓN O VARIACIÓN REPENTINA DE VELOCIDAD	<i>aceleración, frenazo,...</i>
SUSTANTIVOS QUE DENOTAN DETENCIÓN O DISMINUCIÓN DE LA VELOCIDAD, DERIVADOS FORMALMENTE DE LOS VERBOS DEL APARTADO A	<i>frenazo, parada,...</i>
OTROS SUSTANTIVOS QUE DENOTAN REFORMA, EVOLUCIÓN, CAMBIO O TRANSFORMACIÓN DE ALGO	<i>cambio, viraje,...</i>
SUSTANTIVOS QUE DENOTAN GOLPE, CHOQUE, ENCUENTRO VIOLENTO O MOVIMIENTO IMPULSIVO O IMPETUOSO	<i>golpe, tirón, sacudida,...</i>
SUSTANTIVOS QUE DENOTAN CAMBIO DE ESTADO, MÁS FRECUENTEMENTE SI SUPONE PROGRESO, MEJORA, INCREMENTO, DESCENSO O REDUCCIÓN	<i>aumento, cambio, viraje,...</i>
SUSTANTIVOS QUE DENOTAN CAMBIO DE ESTADO, DE DIRECCIÓN O DE POSICIÓN	<i>aumento, transición, viraje,...</i>

---

Tabla 3-23 Clases léxicas del DCECR. con mayor asociación de selección para brusco.



## ABUNDANTEMENTE

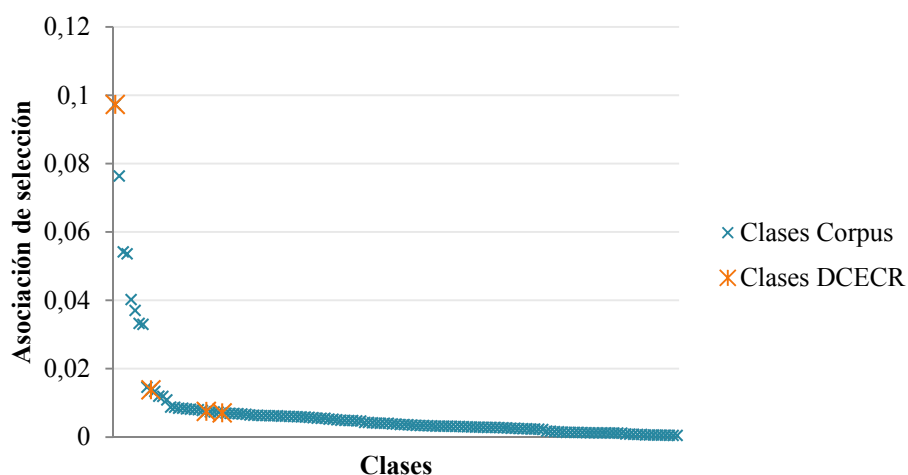


Gráfico 3-29 Asociación de selección sobre el DCECR. para el adverbio abundantemente.

### Clase DCECR. en el corpus

*correr, regar, sangrar,...*

VERBOS QUE DESIGNAN EL PROCESO DE SURGIR O FLUIR UN LÍQUIDO U OTRAS MATERIAS QUE SE LE ASIMILAN FÍSICA O FIGURADAMENTE, TAMBIÉN CON OTROS VERBOS QUE EXPRESAN LA ACCIÓN QUE DESENCADENA ESE PROCESO

VERBOS QUE DESIGNAN DIVERSAS FORMAS DE PONER ALGO A DISPOSICIÓN DE ALGUIEN *proveer,...*

EL VERBO SUDAR Y CON OTROS QUE DENOTAN EXPULSIÓN O DERRAMAMIENTO DE LÍQUIDOS, FLUIDOS U OTRAS MATERIAS QUE PUEDEN ASIMILARSELES, SE USAN A VECES EN SENTIDO FIGURADO *llorar, sudar,...*

OTROS VERBOS QUE DESIGNAN LA ACCIÓN DE FLUIR, DERRAMARSE O ESPARCIRSE ALGO *correr, llorar, sudar,...*

VERBOS QUE DENOTAN ESFUERZO O DIFICULTAD PARA EL DESEMPEÑO DE UNA ACTIVIDAD *sudar,...*

*dar, regar,...*

VERBOS QUE DENOTAN REPARTO, ENTREGA O DISTRIBUCIÓN DE ALGO, A MENUDO IMPRUDENTES O EXCESIVAS

VERBOS QUE DENOTAN EXPULSIÓN O DERRAMAMIENTO DE UN LÍQUIDO *llorar, sudar,...*

VERBOS QUE DENOTAN INGESTIÓN *beber, comer,...*

Tabla 3-24 Clases léxicas con mayor asociación de selección para abundantemente.

## ENÉRGICAMENTE

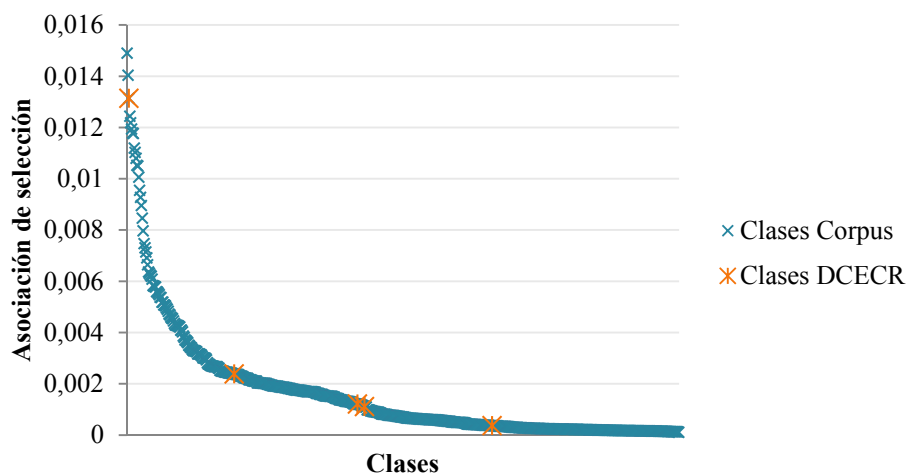


Gráfico 3-30 Asociación de selección sobre el DCECR. para el adverbio enérgicamente.

---

### Clases DCECR. en el corpus

---

<p>VERBOS QUE DESIGNAN MANIFESTACIONES DE OPOSICIÓN, DISENSO O DISCONFORMIDAD EN DIVERSOS GRADOS Y FORMAS, EN OCASIONES EXPRESADAS CON VEHEMENCIA O AGRESIVIDAD</p>	<p><i>protestar,...</i></p>
<p>VERBOS QUE DENOTAN INCUMPLIMIENTO Y CON OTROS QUE DESIGNAN DIVERSAS FORMAS DE RECHAZO DE OBSTACULIZACIÓN DE UN PROCESO</p>	<p><i>protestar,...</i></p>
<p>VERBOS QUE DENOTAN OPOSICIÓN O RECHAZO, A MENUDO MANIFESTADAS VERBALMENTE</p>	<p><i>condenar, protestar, rechazar,...</i></p>
<p>VERBOS QUE DENOTAN LANZAMIENTO (A MENUDO HOSTIL, PERO NO NECESARIAMENTE). TAMBIÉN CON OTROS QUE DESIGNAN, POR EXTENSIÓN, CIERTAS MANIFESTACIONES VERBALES DE TONO AGRESIVO</p>	<p><i>protestar,...</i></p>
<p>VERBOS QUE DESIGNAN ACCIONES VERBALES QUE EXPRESAN CRÍTICA, OPOSICIÓN O RECHAZO EN DIVERSOS GRADOS</p>	<p><i>protestar, condenar,...</i></p>
<p>VERBOS QUE DESIGNAN DIVERSAS FORMAS DE RESISTIRSE</p>	<p><i>condenar, contestar, protestar, reaccionar,...</i></p>

---

Tabla 3-25 Clases léxicas con mayor asociación de selección para enérgicamente (I).

---

**Clases DCECR. en el corpus**


---

VERBOS QUE DENOTAN NEGACIÓN, OPOSICIÓN O RECHAZO, TAMBIÉN CON OTROS QUE DESIGNAN LA ACCIÓN DE ENFRENTARSE A ALGO O A ALGUIEN O LA DE CAUSARLE ALGÚN PERJUICIO	<i>negar, combatir, oponer,...</i>
VERBOS QUE DENOTAN PARTICIPACIÓN ACTIVA, A MENUDO EN ACCIONES CONTRARIAS A ALGO O A ALGUIEN	<i>protestar,...</i>
VERBOS QUE DENOTAN ENFRENTAMIENTO U OPOSICIÓN, TAMBIÉN CON ALGUNOS QUE DESIGNAN OTRAS ACCIONES DE CARÁCTER HOSTIL	<i>atacar, protestar,...</i>
VERBOS QUE DESIGNAN DIVERSAS MANIFESTACIONES VERBALES EXPOSITIVAS O DECLARATIVAS, MÁS FRECUENTEMENTE SI SE REALIZAN ANTE ALGUIEN CON AUTORIDAD O A PETICIÓN DE OTROS. TAMBIÉN CON ALGUNOS VERBOS QUE EXPRESAN LA ACCIÓN DE PRESENTARSE O ACUDIR A REALIZAR ESAS	<i>declarar, intervenir, hablar, pronunciar,...</i>
VERBOS QUE DENOTAN EXPRESIÓN O MANIFESTACIÓN VERBAL	<i>afirmar, contestar, decir, declarar,...</i>

---

Tabla 3-26 Clases léxicas con mayor asociación de selección para enérgicamente (II).

Llegados a este punto, se plantea la necesidad de filtrar de forma objetiva y automática los grupos de argumentos del diccionario ideológico con los que los predicados establecen un nexo preferente. Dado que se pretende captar la relación predicado-argumentos, entre una determinada forma canónica en el rol de predicado y los posibles grupos del diccionario ideológico, se deduce que estamos ante casos originados mediante un mecanismo diferente al que surgen la mayoría de los grupos que combinan con un predicado fijado. Es decir, esta relación genera restricciones de selección, que se reflejarán en los valores de la asociación de selección. Este razonamiento lleva a utilizar los mismos criterios aplicados a las combinaciones simples de formas canónicas, en las que el interés se centraba en aquellas que se pudieran considerar de uso preferente utilizando técnicas de detección de outliers. Todo ello conduce a aplicar tanto el valor *MAD* como *Box-Plot* sobre los valores de la asociación de selección obtenidos por los grupos registrados con la forma canónica sobre la que se haga la consulta. Se exponen los valores máximos de *MAD* para los ejemplos utilizados en este capítulo: *sentir, palpitar, brusco, doloroso, enérgicamente, abundantemente*.

Predicado	Grupo	Asociación de selección	Argumentos
SENTIR	INTRANQUILIDAD	63,48358912	<i>mar de fondo, intranquilidad, polvareda, revuelo, repunta, taco, inquietud, desasosiego, ansia, zozobra, ansiedad, agitación, estrés, preocupación, pesadilla, malestar, rebato, alarma, agobio, suspense, conflicto, torozón, hormigueo, comecome, comezón, desabrimiento, tártago, susidio, combate, guerra, batalla, mareta</i>
	DOLOR	57,53682436	<i>fuego, hormiguillo, hormiguilla, piquiña, fogaje, ardor, picazón, picor, quemazón, prurito, comezón, comecome, rascazón, escozor</i>
	HABLA	47,67730043	<i>sibilítico, sibilino, ambiguo, impreciso, figurado, traslaticio, lato, recto, literal</i>
	DOLOR	39,47676191	<i>tormento, punzada, dolor, analgesia, sinalgia, yaya, acritud, rayo, clavo, puntada, ramalazo</i>
	DESAGRADO	38,70660407	<i>fregado, mangangá, cosijoso, lipidioso, corruptente, secante, molesto, importuno, latoso, pesado, engorroso, inconfortable, incómodo, chinchoso, endemoniado, duro, cansino, cansado, fastidioso, gravoso, oneroso, trabajoso, premioso, pijotero, fuñique, grave, jodido, dichoso, prolijo, matador, mortal, hediondo, guerrero, infernal, pajolero, odioso, insoportable</i>
	INDUCCIÓN	32,72506461	<i>aguijonazo, aguijón, anzuelo, espolada, espuela, pinchazo</i>
	SUFRIMIENTO	32,53837553	<i>pensión, desplacer, punzada, patatús, martelo, arestín, disgusto, puñalada, varapalo, sinsabor, desazón, tribulación, pesadumbre</i>
	ENOJO	31,69693666	<i>enfado, enojo, entripado, picazón, pique, escozor, ofensa, atufó, mosca, mohína, indignación</i>
	DESAGRADO	31,60851849	<i>repelencia, horror, repulsión, hámagó, náusea, asco, repugnancia, grima, dentera</i>
	SUFRIMIENTO	31,39363334	<i>sufrimiento, hiel, agrá, agraz, acíbar, amargura, aflicción, pena, dolor, angustia, pesar, congoja, duelo, desconsuelo, sentimiento, apuro, ahogo, quebranto, presura, pasión, agonía</i>

Tabla 3-27 Grupos de selección preferente según MAD para el verbo sentir

Predicado	Grupo	Asociación de selección	Argumentos
PALPITAR	VALENTÍA	95,52566745	<i>hígados, redaños, arrestos, agallas, pecho, corazón</i>
	AGITACIÓN	88,7261206	<i>aleteo, salto, sístole, diástole, contracción, pulsación, palpitación, golpe, latido, pulso</i>
	CENTRO	85,46242544	<i>yema, ombligo, vientre, riñón, corazón, entraña</i>
	BENEVOLENCIA	84,92569175	<i>benevolencia, entraña, corazón, ternura, filantropía, humanidad</i>
	SENTIMIENTO	74,3173042	<i>sentimiento, podredumbre, psicología, entretelas, corazón, afección, sentir</i>
	DEDO	63,9634519	<i>pulgar, dedo, meñique, auricular, anular, corazón, medio, cordial, índice, dátil</i>
	VOLUNTAD	63,3948113	<i>voluntad, corazón, sentido, intención, merced, elección, albedrío, opción, facultad, potestad, libre albedrío, libertad, arbitrio</i>
	CORAZÓN	62,2941331	<i>corazón, pericardio, endocardio, válvula tricúspide, válvula mitral, ventrículo, aurícula, miocardio</i>
	PECHO	61,7956846	<i>pigidio, esternón, balancín, escólex, telson, escutelo, busto, pecho, seno, mama, teta, aréola, pezón, tetilla, ubre, chiche, mamila, metatórax, mesotórax, pronoto, protórax, cefalotórax, tórax</i>
	VENA	60,3168244	<i>hacecillo, vena, capilar, vaso, arteria</i>

Tabla 3-28 Grupos de selección preferente según MAD para el verbo palpitar

En los resultados extraídos respecto a los adjetivos brusco y doloroso se observa que ambos adjetivos se aplican a los mismos grupos, pero cada uno de ellos selecciona de forma preferente un único grupo que cumpla la restricción de  $MAD \geq 4,5$  en el caso de brusco lo hace con: “*escaramuza, alborada, refriega, pugna, batalla, combate, liza, lucha, lid, choque, acción, fregado*” mientras que doloroso lo hace con: “*azote, flagelo, ramalazo, zurriagazo, tártago, rayo, desventura, infortunio, mal, desgracia, desdicha, desastre, tristezas, lacra, suceso, calamidad, fatalidad, malaventura, cataclismo, catástrofe, bancarrota, tragedia, drama, miseria, plaga, desmán, trago, tormenta, través, hecatombe*” (Tabla 3-29 y Tabla 3-30).

<b>Predicado</b>	<b>Grupo</b>	<b>Asociación de selección</b>	<b>Argumentos</b>	<b>MAD</b>
<b>BRUSCO</b>	GUERRA	0,00548033	<i>escaramuza, alborada, refriega, pugna, batalla, combate, liza, lucha, lid, choque, acción, fregado</i>	6,47377325
	MANO	0,00117275	<i>mano</i>	1,59933181
	TALLO	0,00073327	<i>tallo, carda, cladodio, pedicelo, estolón, bohordo, tronco, talo, lleta, guía, mástil, estipe, mata, pie</i>	1
	CORAZÓN	0,00044792	<i>corazón</i>	0,61084944
	REVÉS	0,00039014	<i>azote, flagelo, ramalazo, zurriagazo, tártago, rayo, desventura, infortunio, mal, desgracia, desdicha, desastre, tristezas, lacra, suceso, calamidad, fatalidad, malaventura, cataclismo, catástrofe, bancarrota, tragedia, drama, miseria, plaga, desmán, trago, tormenta, través, hecatombe</i>	0,53205471

Tabla 3-29 Grupos de selección preferente según MAD para el adjetivo brusco

<b>Predicado</b>	<b>Grupo</b>	<b>Asociación de selección</b>	<b>Argumentos</b>	<b>MAD</b>
<b>DOLOROSO</b>	REVÉS	0,00548033	<i>azote, flagelo, ramalazo, zurriagazo, tártago, rayo, desventura, infortunio, mal, desgracia, desdicha, desastre, tristezas, lacra, suceso, calamidad, fatalidad, malaventura, cataclismo, catástrofe, bancarrota, tragedia, drama, miseria, plaga, desmán, trago, tormenta, través, hecatombe</i>	4,46302398
	CORAZÓN	0,00117275	<i>corazón</i>	0,37723348
	GUERRA	0,00073327	<i>escaramuza, alborada, refriega, pugna, batalla, combate, liza, lucha, lid, choque, acción, fregado</i>	0,37723348
	TALLO	0,00044792	<i>tallo, carda, cladodio, pedicelo, estolón, bohordo, tronco, talo, lleta, guía, mástil, estipe, mata, pie</i>	1
	MANO	0,00039014	<i>mano</i>	1,1029913

Tabla 3-30 Grupos de selección preferente según MAD para el adjetivo doloroso

Predicado	Grupo	Asociación de selección	Argumentos
ABUNDANTEMENTE	SUDOR	482,4811368	<i>trasudar, resudar, sudar, transpirar</i>
	INDICIO	394,648837	<i>transpirar, transparentarse, adivinarse, revelar, repuntar, amanecer, dibujarse, perfilarse, indiciar, indicar, aflorar, asomar, salir, surgir, declararse, reflejarse</i>
		305,9141637	<i>alargar, proporcionar, dar, dotar, surtir, servir, proveer, suministrar</i>
	PREPARACIÓN	228,4413626	<i>aparejar, disponer, aprevenir, perdigar, mullir, poner, prevenir, preparar, armar, formatear, estructurar, organizar, proveer, reservar, apercibir, pertrechar, parar, conrear, aprestar, aviar, prevenir, preconcebir, sembrar, intentar, aprontar, avisar</i>
	FUERA	213,5840966	<i>trasvinar, revenir, instilar, rezumar, exudar, sudar, manar, efluir, fluir, brotar, destilar, resurgir, surgir, surtir, saltar, emanar</i>

Tabla 3-31 Grupos de selección preferente según MAD para el adverbio abundantemente

Predicado	Grupo	Asociación de selección	Argumentos
ENÉRGICAMENTE	COLOR	882,4697135	<i>sonrosar, purpurar, embijar, rubificar, embermejer, arrebolrar, colorear, rojear, enrojecer</i>
		701,3689467	<i>invaginar, rebotar, redoblar, volver, trabar, triscar, engarabitar, engarabatar, escarzar, enarcar, ahorquillar</i>
	SUJECCIÓN	700,4808225	<i>embrochalar, encofrar, entibar, apuntar, apuntalar, sujetar, atrancar, amacizar, lastrar, recibir, estacar, enramar, envigar, escorar, barretear, atarugar, enrabar, pinzar, enzarzar, encañar, ahorquillar, estaquear</i>
	NEGAR	170,062065	<i>repulsar, rehusar, denegar, desestimar, desaprobar, negar</i>
	REGOCIJO	88,92719291	<i>palmeaar, palmoteaar, loqueaar, tronchase de risa, desternillarse de risa, mondarse de risa, carcajear, sonreír, reír</i>

Tabla 3-32 Grupos de selección preferente según MAD para el adverbio enérgicamente





## 4. Aplicación desarrollada

### 4.1. Introducción

Se presenta en este capítulo la aplicación ColexWeb, recurso lexicográfico que hace accesible a través de la web el acceso a la información contenida en la base de datos que alberga toda la información combinatoria extraída del corpus. El diseño de ColexWeb permite consultar los datos respecto a las capacidades combinatorias de una determinada palabra y visualizarlos, tanto textual como gráficamente, obteniendo así una comprensión adecuada de los mismos. A través de ella también se pone a disposición de la comunidad un motor de extracción de colocaciones sobre cualquier texto que se proporcione, en el que se han implementado los procedimientos recomendados según la experiencia alcanzada. De esta manera se ofrece a través de internet el acceso a este caudal de información combinatoria del español, de especial utilidad en la enseñanza de la lengua, estudios de los lexicógrafos o tareas de eliminación de ambigüedad en el Procesamiento del Lenguaje Natural.

Se conjuga en la aplicación que aquí se presenta características de otras herramientas orientadas a la extracción de colocaciones como son *Collocate*, el *Tollkit UCS* o el diccionario *Dice*:

- *Collocate* es una aplicación de escritorio de pago que ha sido desarrollada para extraer colocaciones y términos a partir de un corpus textual, mediante el análisis estadístico de las mismas. Extrae las colocaciones de un corpus como un todo, permitiendo la selección del tamaño de ventana de análisis y la elección de los indicadores a calcular entre información mutua, t-score y la razón de verosimilitud.
- El *Tollkit UCS* desarrollado en Perl y su versión mejorada en el R-Project permiten el análisis estadístico de coocurrencia de datos. Extrae las combinaciones existentes en un texto, compuestas por dos palabras, asociadas a sus frecuencias de aparición, y almacenadas en ficheros de texto plano. Permite operaciones como la visualización e impresión de datos, manipulación, ordenación y asociación con un rango de medidas calculadas.
- *DiCE* es una aplicación web que representa un diccionario de colocaciones del español, con una nomenclatura limitada actualmente al campo semántico de los nombres asociados a sentimiento. Además de otra información como por ejemplo (cuasi)-sinónimos y (cuasi)-antónimos, se presentan para cada unidad léxica las colocaciones con información de las funciones léxicas.

ColexWeb abarca aspectos de todas ellas ya que permite la consulta de las capacidades combinatorias de las palabras del español (*DICE*), proporcionándola de forma extensional, enriqueciéndola con las posibilidades de ordenación según los valores de alguno de los indicadores (*Collocate*) y facilitando la interpretación y valoración de los datos a partir de técnicas de visualización (*Toolkit UCS*). Como novedades adicionales se incorporan los indicadores basados en la detección de Outliers además de información de las preferencias de selección entre un predicado y sus argumentos a través de los valores de la asociación de selección sobre el espacio semántico del D.I.V.

Los requisitos fundamentales en el desarrollo de esta herramienta radican en el funcionamiento a través de la web, el diseño de una interfaz de usuario amigable, la forma de acceso a la BDD a través de un servicio web dedicado, disponible para la integración en herramientas que desarrollen otros programadores y la optimización en términos de velocidad, intentando siempre obtener respuestas en un rango de tiempo adecuado.

## 4.2. Arquitectura del sistema

Los principales componentes del sistema se resumen en el diagrama de despliegue en la Figura 4-1

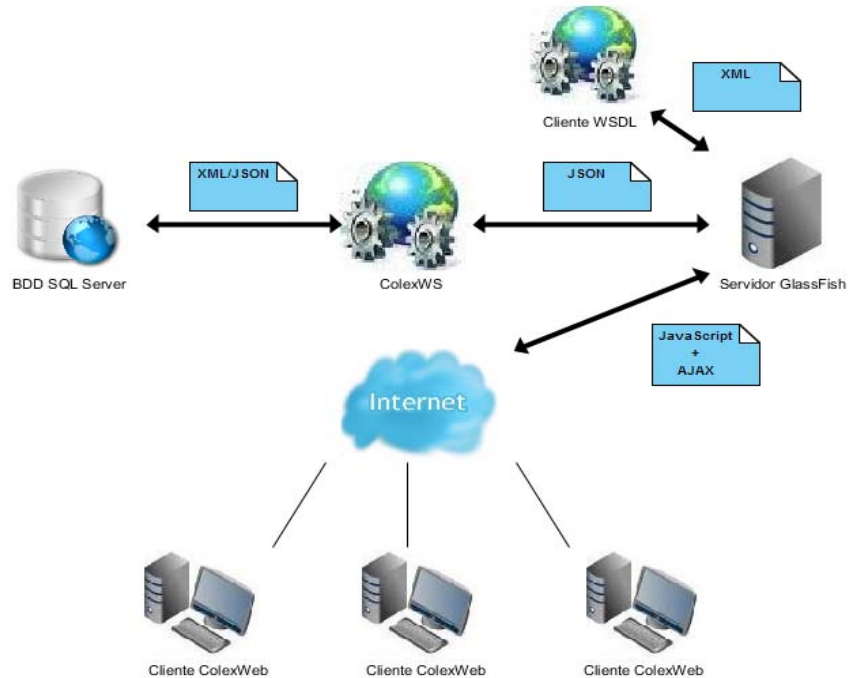


Figura 4-1 Diagrama de Despliegue

- **ColexWeb, aplicación web cliente.** Encargada de la interacción del usuario para realizar tanto el análisis de palabras sueltas como de un fichero de texto correctamente formateado.

- 
- **ColexWS, Web Service que enlaza base de datos con ColexWeb.** Permite el acceso a la base de datos desde la aplicación web cliente, realiza operaciones de consulta para obtener los datos de indicadores seleccionados por el usuario. Sigue el protocolo REST.
  - **Ciente Web Service Lematizador.** Servicio Web proporcionado por el Grupo de Estructura de Datos y Lingüística Computacional, GEDLC, de la Universidad de las Palmas de Gran Canaria. Se encarga de lematizar cualquier palabra del español identificando su forma canónica, categoría gramatical y la flexión o derivación que la produce, devolviendo las formas correspondientes. Es un web service de tipo SOAP, accedido a través del fichero WSDL, que a su vez genera las clases necesarias que pueden ser accedidas desde la aplicación web cliente.
  - **Servidor contenedor de aplicaciones.** Actúa como contenedor de la aplicación web cliente y el web service que interactúa con la base de datos, además de procesar las peticiones realizadas por el cliente. Las operaciones de solicitud que realiza son las siguientes:
    - *Carga del fichero de texto involucrado en el análisis.*
    - *Análisis de combinaciones en un fichero de texto.*
    - *Consulta de indicadores generales.*
    - *Consulta de indicadores propios.*
    - *Cálculo de selección léxica.*

Sobre el servidor se ejecutan dos servlets: el primero establece el proxy de comunicación asíncrona que es el encargado de establecer la conexión con la aplicación web cliente, mientras que el segundo es un servlet común destinado a la carga del fichero a analizar cuando se realiza esta operación. El motivo de este segundo servlet se debe a la necesidad de realizar el parseo traducir del fichero en el lado servidor, ya que el cliente GWT no dispone de herramientas para ello.

- **Base de Datos.** Recoge todas las palabras del corpus analizadas en su contexto y clasificadas por tipos de combinaciones, con los datos de categorías, frecuencias e indicadores, entre otros, correctamente formateados, sobre el Sistema Gestor de Base de Datos Microsoft SQL Server 2008 R2.

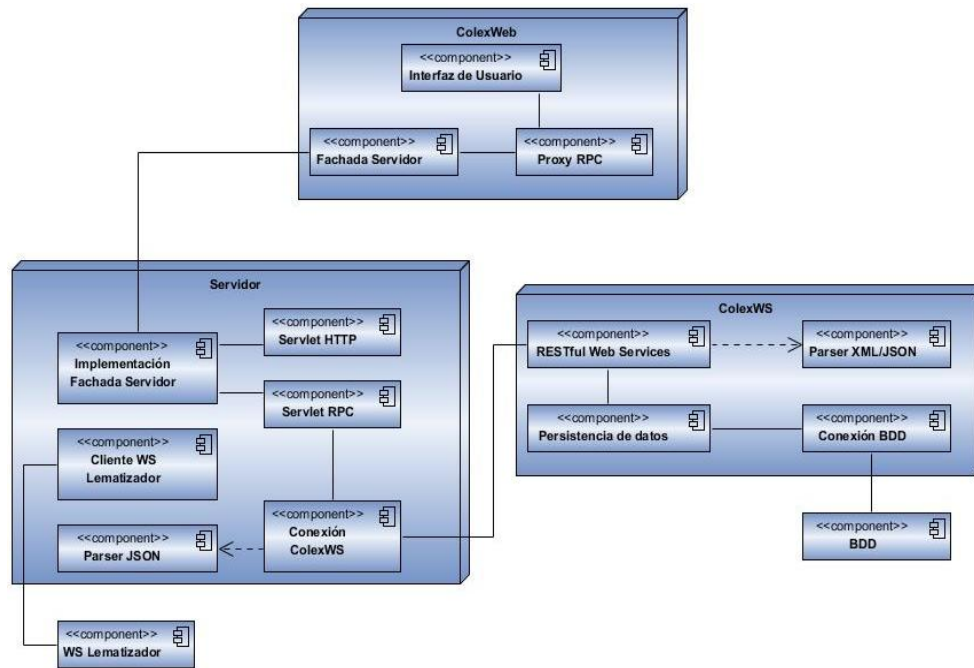


Figura 4-2 Diagrama de Componentes

El diagrama de componentes que aparece en la Figura 4-2 muestra la composición interna de cada uno de los elementos que intervienen en la arquitectura del sistema descrito:

### 4.3. Análisis de casos de uso

La aplicación web no precisa, hasta el momento, de un módulo dedicado a la identificación de usuarios para permitir su acceso ni restricción de procesos en base a perfiles. Por tanto, se considera un único actor que interactúa con la aplicación, denominado *lingüista*, el cual puede acceder a los procesos de análisis de palabras sueltas y procesado de un fichero de texto para extraer las combinaciones.

Teniendo en cuenta las consideraciones anteriores, se presenta el diagrama de casos de uso para la aplicación web cliente en la Figura 4-3.

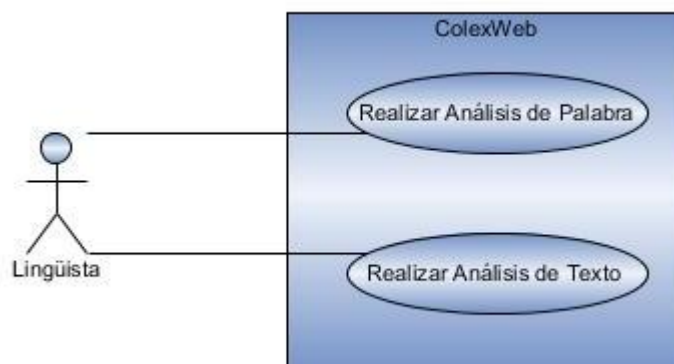


Figura 4-3 Diagrama de Casos de Uso

Se describen los casos de uso identificados:

- **Realizar análisis de palabra.** El lingüista introduce una determinada palabra del español a través de la aplicación, que solicita además seleccionar el indicador de consulta y el tipo de combinación a extraer de la base de datos de entre las permitidas por la palabra anterior. Se presenta así un conjunto de registros en la tabla con los datos solicitados, y una representación gráfica de los mismos, en el caso que se obtengan resultados en la base de datos.
- **Realizar análisis de un texto.** El lingüista indica el fichero que contiene el texto a analizar, además de especificar los indicadores que se desean calcular, el tipo de combinaciones a extraer y el tamaño de cinta para indicar la distancia máxima entre palabras. Tras realizar la carga del fichero en el servidor, se inicia el proceso de extracción. Se presenta así un conjunto de registros en la tabla con los datos solicitados, y una representación gráfica de los mismos, en el caso que se obtengan resultados con el tipo de combinación elegido.

#### 4.4. ColexWeb: Aplicación web cliente

Este apartado está dedicado a la descripción de ColexWeb. Como se ha mencionado, representa la aplicación web cliente a la que tiene acceso el usuario lingüista y con la cual interactúa para obtener los resultados que desea en función del proceso lanzado.

Su construcción se ha realizado utilizando el *framework Google Web Toolkit, GWT*. GWT impone una división cliente-servidor, de forma que en el cliente se realice el diseño de la interfaz de usuario y posibles validaciones, y en el servidor se realicen operaciones relacionadas

con la interacción con el modelo de datos de la aplicación. Por tanto, se considera la división de ColexWeb mostrada en la Figura 4-4.

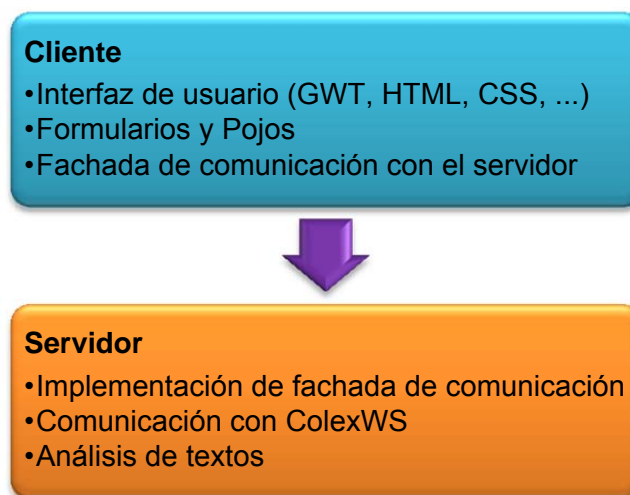


Figura 4-4 Estructura Cliente-Servidor

La interfaz de usuario que requiere ColexWeb no precisa de un gran despliegue, ya que requiere parámetros muy concretos para su funcionamiento en cualquiera de los procesos que ejecuta. Por ello, se busca simplicidad en el diseño para maximizar la fácil comprensión y usabilidad por parte del usuario. Dado que se ha orientado como una aplicación web en lugar de un conjunto de páginas, se dispone de una única pantalla principal sobre la que se realiza todo el trabajo que necesita el usuario, simulando así una aplicación de escritorio, pero compartida por la red. En la Figura 4-5 se observa pantalla inicial que puede visualizar un usuario.

Se identifican los elementos que se describen a continuación:

- Cabecera con título de la aplicación, situado en la parte superior del navegador.
- Menú desplegable ubicado en la zona central izquierda, justo debajo de la cabecera, con dos secciones diferenciadas para identificar los procesos: **Análisis de palabras** y **Análisis de un texto**. En la Figura 4-5 se muestra el correspondiente al análisis de palabras y en la Figura 4-6 al análisis del texto.
- Tabla que muestra los registros seleccionados, ubicada en la zona central derecha del navegador, justo debajo de la cabecera. Inicialmente, la tabla está vacía, se indica con el símbolo en gris la ausencia de resultados.

- Panel destinado a la representación de gráficas lineales para dibujar los datos presentados en la tabla anterior. Se ubica en la parte inferior de la pantalla, en blanco cuando ColexWeb ha sido iniciado.

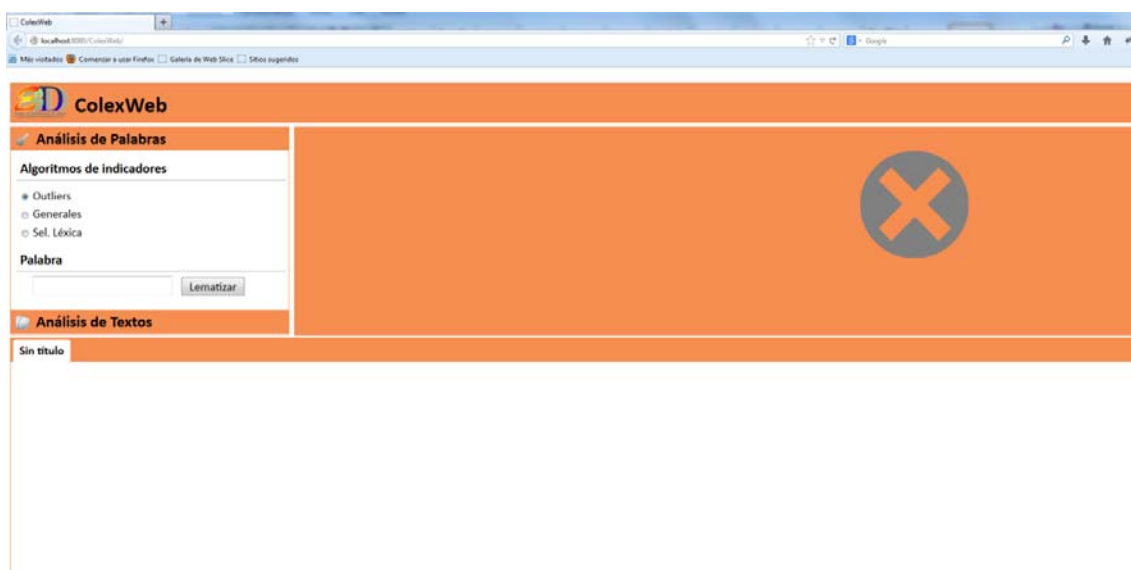


Figura 4-5 Pantalla inicial



Figura 4-6 Panel que se muestra para el análisis de un texto.

#### 4.4.1. Análisis de palabra

Se solicita la introducción de la palabra y el indicador o indicadores de interés. Se muestran los resultados de la lematización, especificando cada forma canónica resultante y su correspondiente categoría gramatical, y las categorías compatibles para construir una colocación sobre las que trabaja ColexWeb. El usuario puede realizar consultas sobre combinaciones cuyas formas canónicas tienen categoría gramatical: *verbo + sustantivo*, *sustantivo + adjetivo* y *verbo + adverbio*. Se permite seleccionar entre:

- Uno o varios indicadores de uso general: Frecuencia Relativa respecto a la base o el colocativo, Información Mutua, t-score y Z-Score (Figura 4-7).

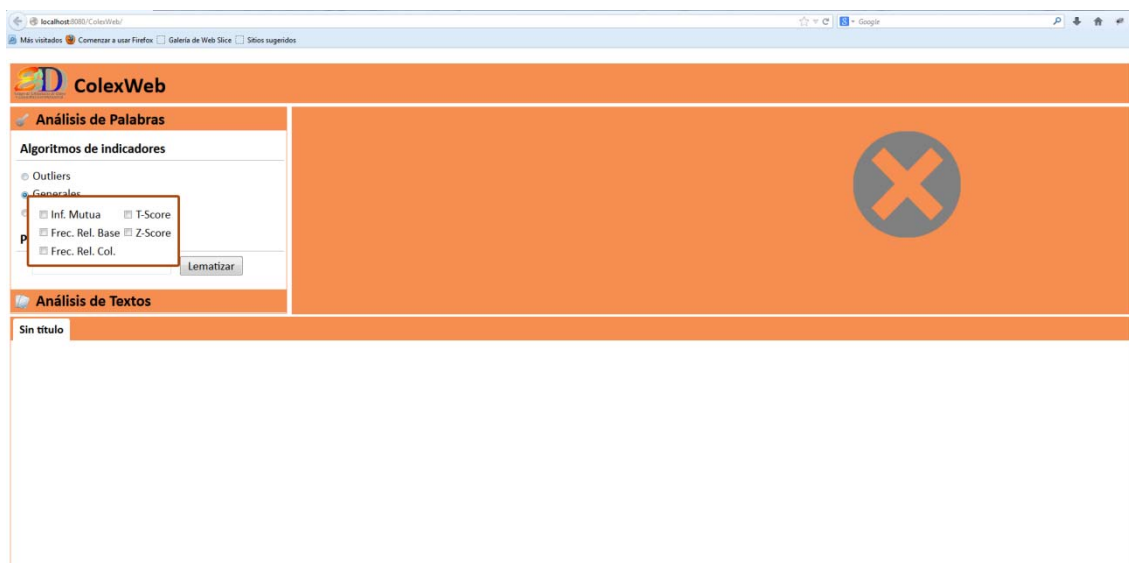


Figura 4-7 Análisis utilizando indicadores convencionales

- Indicadores basados en outliers: MAD y Box-Plot.
- Selección léxica.

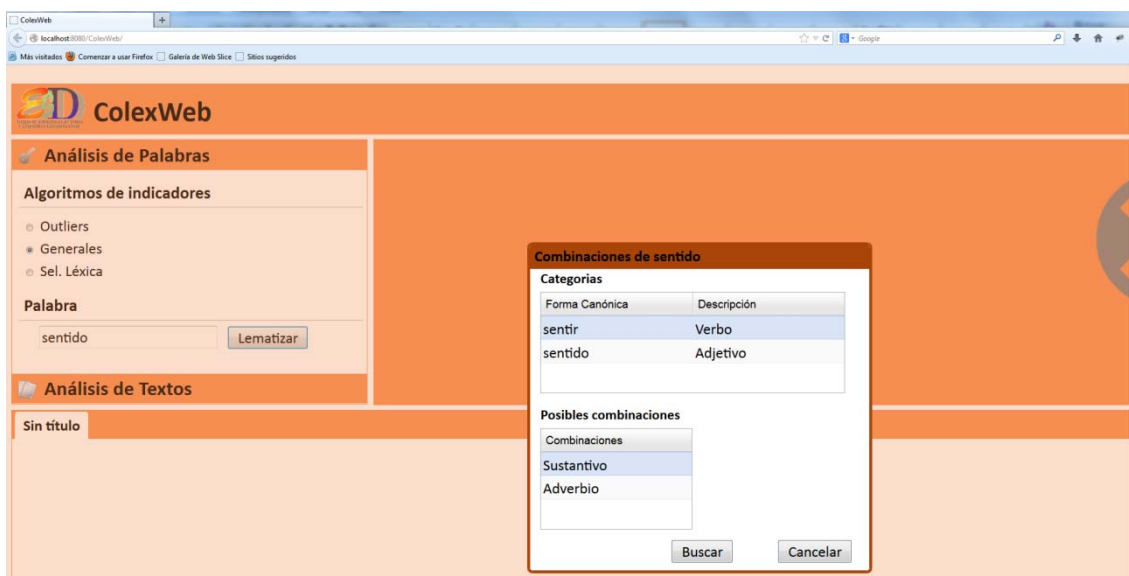


Figura 4-8 Lematización de la palabra y selección de combinaciones permitidas.

Una vez resuelta la consulta se muestra una tabla con las formas canónicas integrantes de la combinación y sus puntuaciones obtenidas del corpus según el indicador seleccionado. En



la parte inferior se muestra una gráfica con los resultados que facilita el análisis. El flujo se repite para la selección léxica y outliers, con pinceladas características de cada caso.

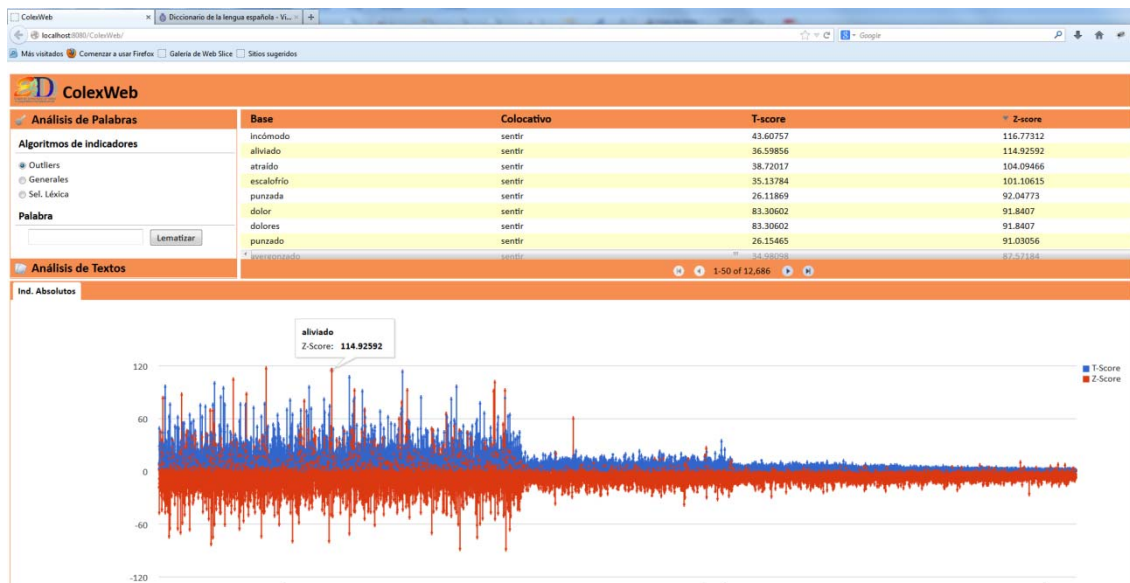


Figura 4-9 Resultados para la consulta sentido → sentir + sustantivo. Indicadores z-score y t-score

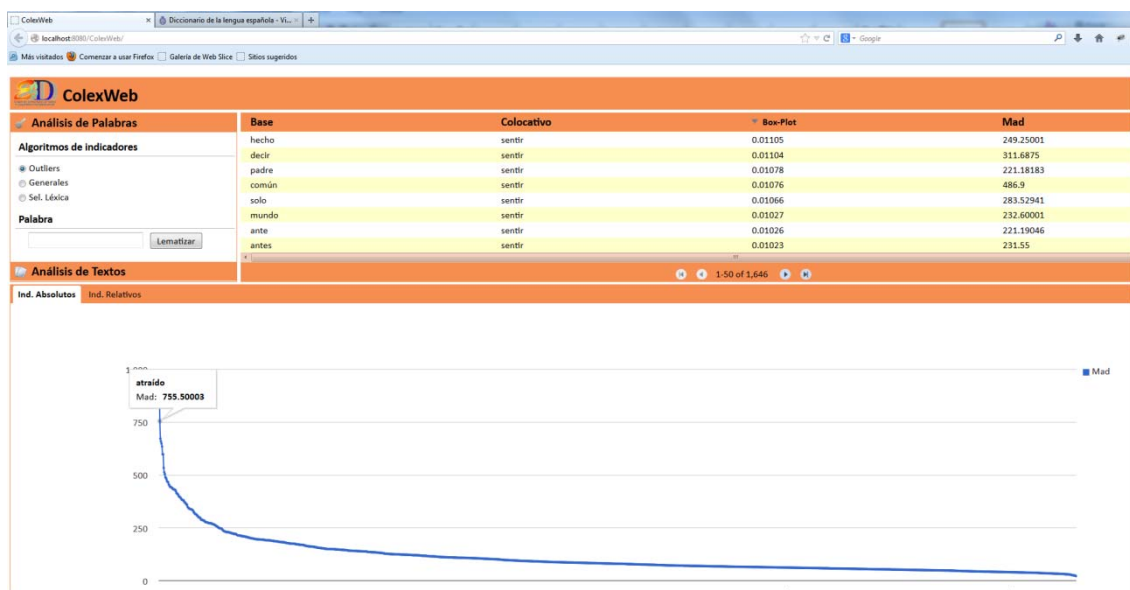


Figura 4-10 Resultados para la consulta sentido → sentir + sustantivo. Indicadores basados en outliers

Los resultados en la tabla se pueden mostrar ordenados por orden alfabético o por el valor numérico, según el usuario haga click sobre alguna de las columnas de formas canónicas o de indicadores. Si se sitúa el ratón sobre algún punto de la gráfica que se considere de interés se puede obtener la información de la combinación a la que hace referencia y su puntuación. En el caso de consultar sobre la Selección Léxica la información que se muestra en el globo es el conjunto de palabras que conforman el grupo semántico representado por ese punto.

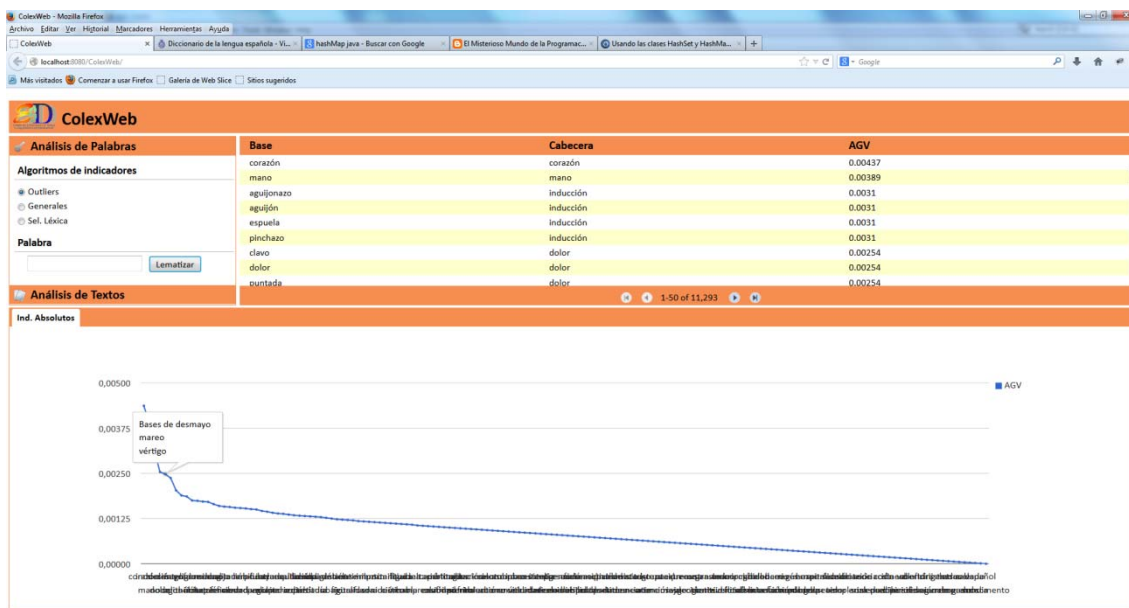


Figura 4-11 Ejemplo de consulta sobre selección léxica.

#### 4.4.2. Análisis de texto

En el caso de análisis del texto, los parámetros solicitados son el fichero de texto a procesar y los indicadores que se desean extraer; en este caso, se permiten los indicadores generales: *Frecuencias Relativas*, *Información Mutua*, *t-score* y *Z-Score* y también los basados en outliers. A continuación, se abre una ventana en la que poder indicar el tipo de combinaciones a extraer y el tamaño de la concordancia. Las concordancias se restringen al ámbito de las oraciones (Figura 4-12, Figura 4-13). Al igual que el proceso anterior, se presenta la tabla de datos con los indicadores seleccionados y una ventana previa que indica el tiempo de ejecución empleado, total de combinaciones extraídas, número de palabras analizadas en el texto y total de veces que se ha invocado el web service lematizador.

En cuanto a la creación de gráficas, el funcionamiento es distinto. La diferencia radica en que no se tiene una palabra principal sobre la que generar el conjunto de resultados, sino que puede ser cualquiera que forme una colocación. De esta manera, se permite generar las gráficas correspondientes a la selección del usuario sobre las celdas de las columnas para base y colocativo. Por ejemplo, si se elige la palabra *estrellar*, colocativo en la combinación de tipo Sustantivo-Verbo (*estrellar*, *presentado*), se obtienen todas aquellas combinaciones en las que *estrellar* es colocativo y genera el conjunto de gráficas asociadas en función de los indicadores elegidos. El flujo, asociado al ejemplo descrito, es el siguiente:

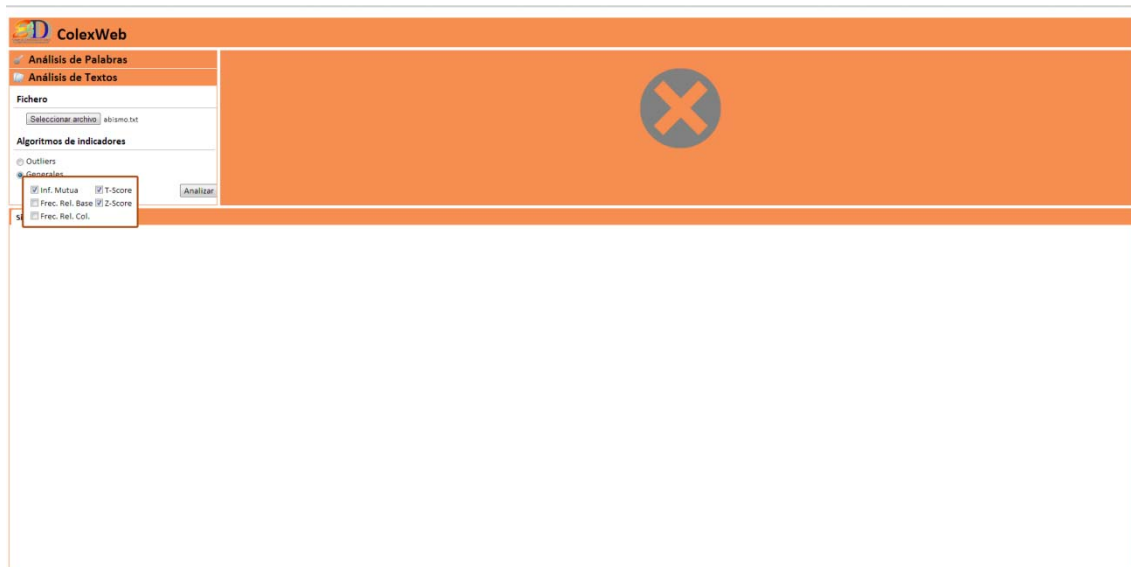


Figura 4-12 Inicio del análisis de textos.

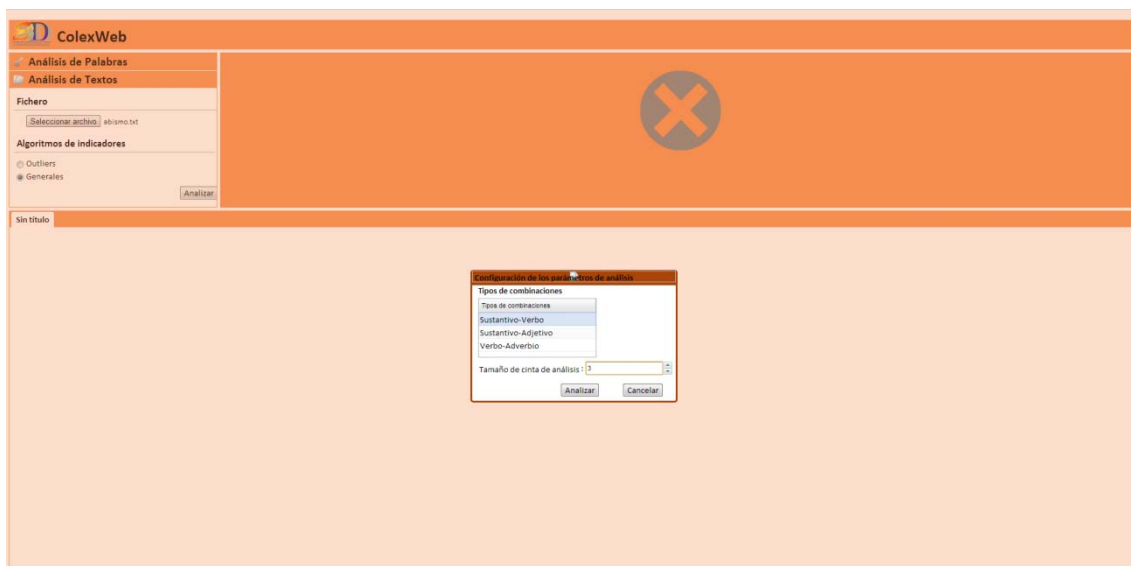


Figura 4-13 Configuración del análisis de textos.

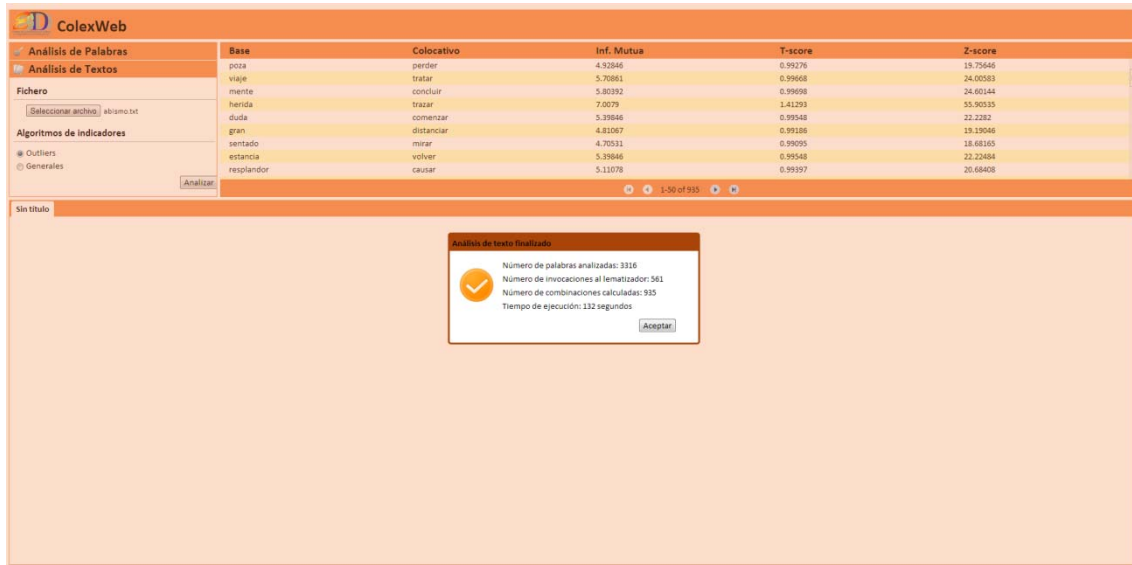


Figura 4-14 Resultados obtenidos del análisis de texto.

Durante la ejecución de los procesos, se presentan otras ventanas auxiliares. Dado que las operaciones de consulta y procesamiento de textos requieren un cierto tiempo de recuperación y presentación de datos, se visualizan diálogos de carga indicando el proceso realizado, que impide la interacción con la aplicación mientras se resuelve. También se realiza un control de errores sobre los datos introducidos, tales como la existencia de la palabra o del fichero de texto a analizar, se exige especificar como mínimo un indicador general si se solicita obtenerlos, se controla la imposibilidad de lematizar una palabra, etc.

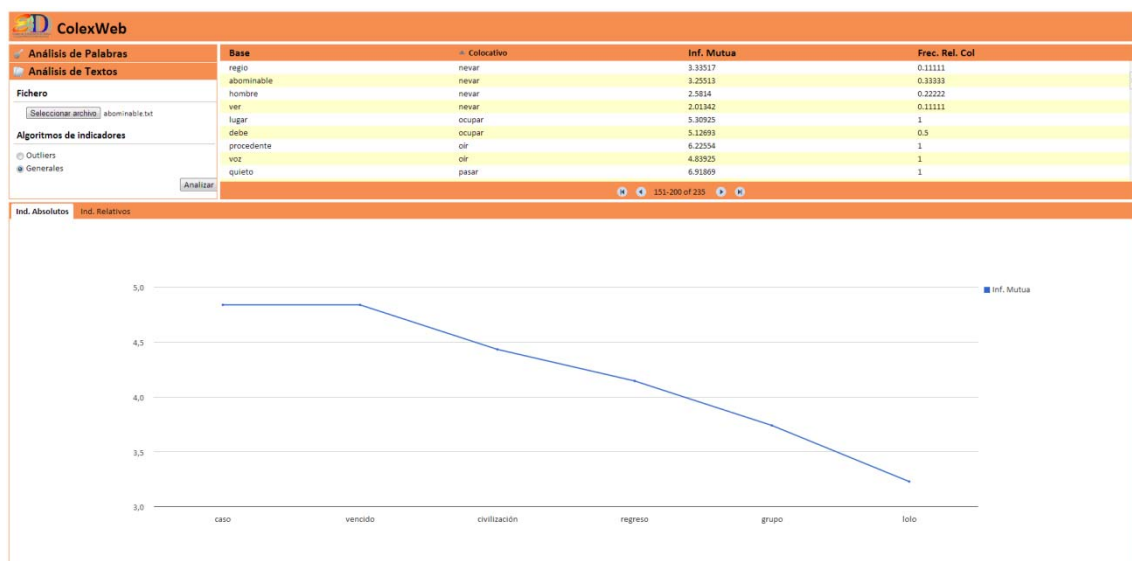


Figura 4-15 Resultados obtenidos del análisis de texto.

#### **4.5. Elementos para acelerar el análisis de textos.**

El análisis de los textos se lleva a cabo aplicando la misma metodología seguida para el procesamiento del corpus con el que se ha generado la BDD de combinaciones. Sin embargo, un factor crítico a tener en cuenta es el tiempo de respuesta en el procesamiento, por lo que además de unir estructuras con un acceso lo más rápido posible, se hace necesario construir un algoritmo que ofrezca una destacada eficiencia. A este fin, además de la lista de palabras vacías ya descrita en el capítulo 2, se ha generado un diccionario de palabras frecuentes.

Las palabras frecuentes consisten en un conjunto de palabras asociadas a una lista con todas las formas canónicas que produce su lematización, así como la categoría gramatical de la misma. Antes de lematizar una palabra se comprueba si ya existe en el diccionario, obteniendo los datos pertinentes de esta estructura local. El objetivo es evitar el coste computacional de ejecutar el lematizador cada vez que se necesite obtener las formas canónicas de una palabra. Debido al gran tamaño de este conjunto, su inicialización se delega al momento de inicio de la aplicación.

Este conjunto se ha extraído a partir de un corpus compuesto por 1.000.000 de palabras aproximadamente. Se han contabilizado las frecuencias de palabras no vacías utilizadas, extrayendo 62802 palabras diferentes, donde las 1000 más frecuentes representan el 50% del total empleado, y las 7000 más frecuentes cubren el 80%, como se ilustra en el gráfico de frecuencias acumuladas presentado en la Gráfico 4-1. Del mismo modo, a medida que se procesa un fichero de texto, se va actualizado la lista de palabras frecuentes con aquellas que aparecen en el corpus y que no forman parte del diccionario inicial.

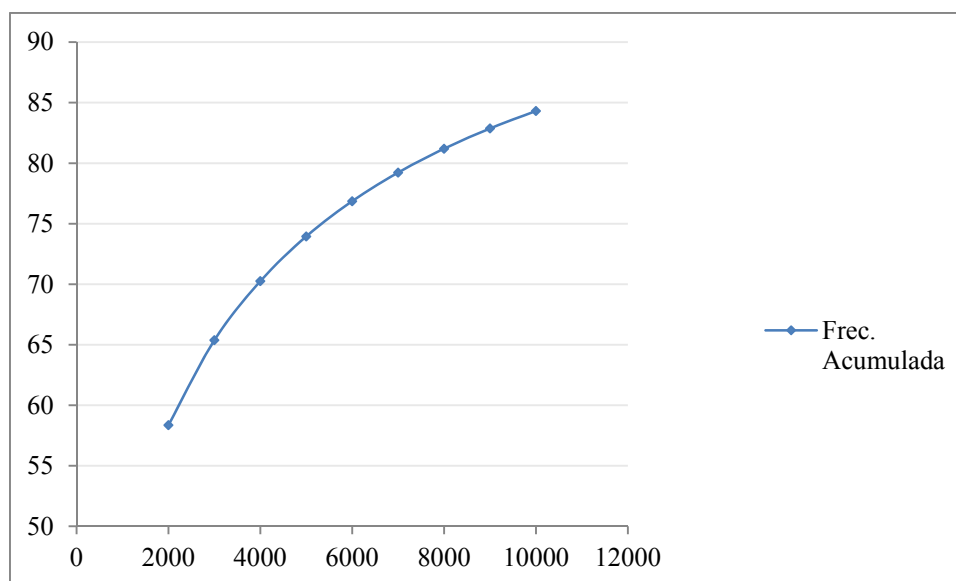


Gráfico 4-1 Frecuencias acumuladas para las palabras del corpus analizado

El Gráfico 4-2 y el Gráfico 4-3 muestran los tiempos invertidos en la ejecución del análisis de texto sobre varios ficheros de 6K, 12K, 113K, y 393K con una cantidad de 1011, 1955, 18724 y 64709 palabras respectivamente. Se observa que no necesariamente un mayor número de palabras conlleva más tiempo de procesamiento debido a la estrategia de incorporación de palabras frecuentes en tiempo real. Es decir, una vez que ha sido procesado un fichero, se mejora el rendimiento del programa, esto se debe por una parte, a que ya han sido cargados los recursos, establecido las conexiones, etc. Por otra parte, se disminuye la cantidad de invocaciones al lematizador, no solo por las palabras precargadas, sino por las que se van agregando en tiempo de ejecución. Finalmente, el tiempo de proceso lo consume la exploración de posibles combinaciones, ya que los tiempos de proceso no presentan diferencias extremas según la estructura de la colocación, pero hay una diferencia importante en la cantidad de combinaciones generadas según se trate de *Sustantivo + Verbo*, *Sustantivo + Adjetivo* o *Verbo + Adverbio*.

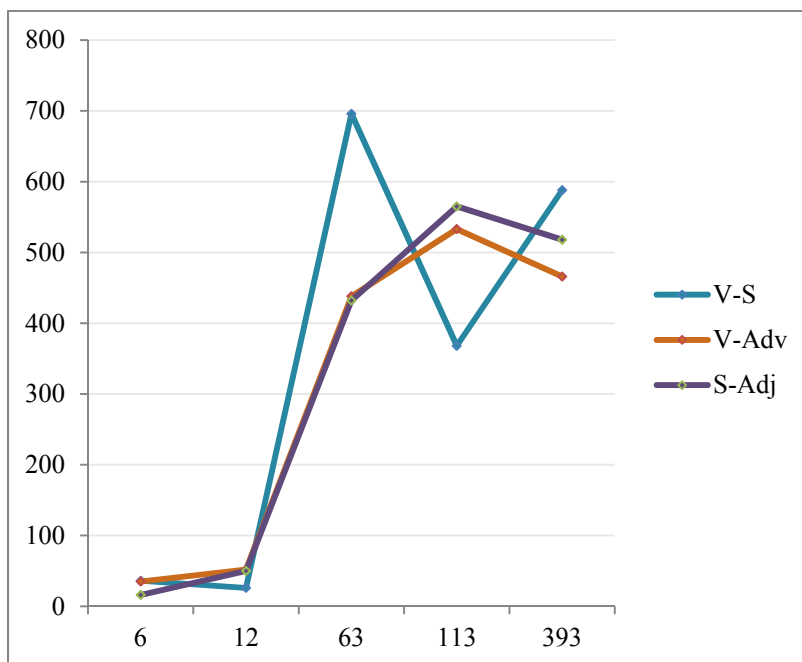


Gráfico 4-2 Tiempo de ejecución de ColexWeb para la Frecuencia Relativa sobre ficheros de distintos tamaños

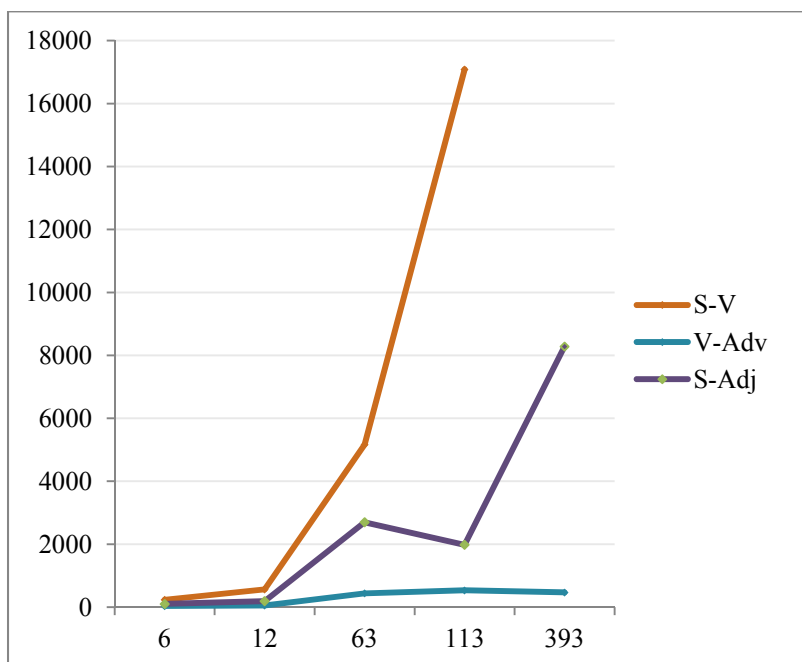


Gráfico 4-3 Combinaciones extraídas en ficheros de distintos tamaños.





## 5. Conclusiones

En esta tesis se ha abordado el estudio de la extracción automática de colocaciones léxicas del español a partir de corpus textuales extensos.

1. En primer lugar se ha realizado una recopilación de características lingüísticas de este tipo de combinaciones con objeto de utilizar aquellas que pudiesen repercutir en el tratamiento computacional del problema. En esta fase se concluyó que las palabras objetivo debían ser: nombres, adjetivos, verbos y adverbios, en las estructuras contenidas en la clasificación de colocaciones simples aportada por Koike (2002). Por otra parte, se considera esencial la flexibilidad formal, por lo que se opta por trabajar en el espacio de las formas canónicas en lugar de hacerlo en el espacio de las palabras. La característica de la preferencia, es decir, que la colocación se elija mayoritariamente por los hablantes frente a otras combinaciones posibles, se traduce en el enfoque basado en las frecuencias que tradicionalmente se ha dado al problema en cuestión. Esta última constituye el principal instrumento para el procesamiento de corpus textuales en busca de colocaciones. También destacan los numerosos casos en los que la asociación se produce entre una palabra y un grupo de palabras con rasgos semánticos comunes, característica directamente relacionada con la cualidad de los predicados de imponer restricciones sobre los posibles argumentos de los que predicen, hecho que explica la teoría lingüística de la selección léxica.
2. Respecto a los recursos que se utilizaron para el desarrollo del trabajo, se compiló un corpus de textos sin ningún tipo de información lingüística de un volumen tal que pueda considerarse una amplia muestra del uso del español (300000000 de palabras). También se dispuso de otro corpus de menor extensión, formado por 2000000 de palabras que fue utilizado para realizar estudios comparativos. La transformación de ambos al espacio de las formas canónicas se realizó mediante el lematizador del GEDLC. Como datos de control, que permitan verificar los resultados que se obtienen, se recopilaron todos los ejemplos de colocaciones presentados en artículos lingüísticos sobre colocaciones en español, además de las combinaciones en el DCECR. Otro recurso relevante en este trabajo lo constituye el DIV, como base de conocimiento respecto a grupos semánticos.
3. El corpus fue procesado para obtener todas las combinaciones de formas canónicas que corresponden a alguna de las estructuras objeto de estudio mediante una aplicación en la que el contexto de una palabra se restringe a la frase limitando la línea de concordancia a una amplitud de 10. Se registraron las formas canónicas y sus frecuencias, las combinaciones de formas canónicas que al menos aparecen 3 veces y sus frecuencias

en un BDD, esta fue utilizada para las investigaciones de las técnicas computacionales y para el desarrollo de la aplicación.

4. La vertiente computacional del problema se orientó a la recopilación y evaluación de un amplio catálogo de indicadores basados en la frecuencia de uso, abundantemente utilizados en la literatura. Las medidas de asociación que se probaron incluyen la más simple, frecuencias relativas, la información mutua y las basadas en test estadísticos: z-score, t-score y test de Dunning. Las combinaciones analizadas pertenecen a alguno de los tres grupos que se pueden establecer a partir del lematizador del GEDLC y que se adaptan a las estructuras colocacionales presentes en el español: sustantivo + verbo, sustantivo + adjetivo y verbo + adverbio.
5. El análisis exhaustivo de los datos revela las dificultades de la automatización de este proceso tradicionalmente basado en rankings o en puntos de corte. Por un lado, se establece que las conclusiones son fiables en la medida en que se disponga de suficiente cantidad de muestras de las palabras en el corpus tratado. Por otro lado, la posición de una combinación en el ranking depende, en mayor o menor medida según el indicador que se fije, del uso más o menos elevado de la palabra, normalmente debido a su valor funcional o a la cantidad de acepciones de la misma o a la estructura gramatical de la colocación. Determinar la cantidad de ejemplares de una combinación que se exige que haya en el corpus se realiza manualmente, así como el punto de corte que delimita las combinaciones libres de las colocaciones.
6. De los experimentos analizados sobre el corpus de Galdós, se desprende la diferencia en los valores de frecuencias registrados y en las medidas de asociación alcanzadas por las combinaciones. Se resalta también que la disparidad entre las combinaciones en las que intervienen palabras de baja frecuencia de uso y las que son muy frecuentes se agranda a medida que crece el tamaño del corpus, mostrándose el test de Dunning como el más estable entre todos ellos. Hay que hacer notar la dificultad para encontrar muestras suficientes de una determinada colocación en un corpus por muy extenso que sea, ya que solo aproximadamente la mitad de las combinaciones en el DCECR aparecen en el corpus completo, reduciéndose drásticamente la cantidad de las que cumplen con el requisito del número mínimo de muestras en el de Galdós.
7. Se propone una metodología en la que la extracción se realiza de forma independiente para cada palabra, evitando las comparaciones entre formas canónicas diferentes desde el punto de vista de la frecuencia. Se aborda el problema desde el enfoque de la detección de valores atípicos en las muestras, comparando distintas técnicas utilizadas en estadística para resolver el problema de la detección de outliers. La propuesta que se hace tiene el objetivo de detectar, fijada una palabra, con qué otras se puede establecer que su uso se sale de lo común en su ámbito, por lo que se determinan qué valores

---

extremos alcanzan las frecuencias relativas de la forma canónica seleccionada. Se han valorado tres posibilidades que se denominan: ZChebyshev —se basa en la media aritmética y en la desviación típica, MAD —se basa en la mediana y Box-Plot —se basa en los cuartiles, se incluyen estos últimos por ser más robustos. El primer caso tiene menor cobertura; sin embargo, el estadístico MAD y Box-Plot se revelan como indicadores fiables para la extracción automática de colocaciones, ya que proporciona como principal ventaja la precisión de los resultados, filtrando una gran cantidad de combinaciones libres que se colaban en las mejores posiciones de los rankings. Además proporcionan un valor objetivo en el que establecer el corte, si bien en la frontera se introducen combinaciones libres, las frecuencias relativas coinciden con valores similares a los que registran combinaciones del DCECR.

8. En la inspección de los datos se dejan entrever grupos de palabras relacionadas gracias a cierta afinidad semántica, como es de esperar según la teoría lingüística. Se indaga sobre fórmulas computacionalmente viables para deducir automáticamente cuáles son los que corresponden a los predicados presentes en el corpus. Se valoran técnicas basadas en matrices de contexto para cada palabra: análisis clúster y semántica latente. En ambos casos, debido a la dimensión del corpus, la construcción de las matrices es muy costosa en tiempo y además requiere la interpretación de los resultados para establecer los grupos. Por estas razones se considera que estas soluciones no son apropiadas para una aplicación software en tiempo real como la que marcan los objetivos de este trabajo.
9. Se propone como solución viable y efectiva para captar los nexos motivados por la selección léxica, utilizar la medida de asociación entre verbos y sustantivos de Resnik (1999), sustituyendo WordNet por los grupos semánticos que aparecen en los diccionarios ideológicos, para los que se determina una puntuación denominada asociación de selección. A partir de este valor se pueden determinar rankings de los grupos semánticos de tales recursos léxicos con los que un predicado establece un nexo semántico. Se encuentra que la asociación de selección tiene capacidad de discriminar los grupos semánticos que combinan con un determinado predicado: la mayor parte de los grupos que contienen alguna de las palabras que aparecen junto al colocativo obtienen una puntuación casi nula y siempre se observa un conjunto de grupos del DIV que destaca por su alto valor de asociación de selección. Incluso cuando ninguno de los elementos del grupo tenga una marcada preferencia individual, el conjunto de combinaciones aporta la suficiente fuerza de selección para hacer destacar al grupo. También se valoró sobre los grupos semánticos en el DCECR, se observa que los que se presentan en dicha obra para el predicado en cuestión se posicionan entre las de más alta puntuación. En general, otros grupos con puntuaciones elevadas hacen referencias a

argumentos que comparten rasgos semánticos con los que se combinan con él según este diccionario aunque se presenten en relación con otras entradas.

10. Todo el conocimiento adquirido en el desarrollo de este trabajo se pone a disposición de la comunidad a través de la web mediante la herramienta ColexWeb para la consulta de las capacidades combinatorias de las palabras del español basada en la información extraída del corpus. Mediante una interfaz amigable el usuario puede diseñar la consulta, escogiendo la palabra, alguna de las estructuras admisibles para ella y qué medidas de asociación se quiere utilizar para construir el ranking. Los resultados se muestran por extensión y gráficamente, lo que facilita su comprensión de forma rápida y efectiva. Los indicadores disponibles son: *frecuencias relativas*, *información mutua*, *z-score*, *t-score*, *Box-Plot* y *MAD*. Del mismo modo se admiten peticiones respecto a los grupos semánticos que selecciona un predicado. También se incorpora en ColexWeb una opción para el procesamiento de un fichero de texto plano siguiendo la metodología llevada a cabo para el corpus del GEDLC, mostrándose los resultados de forma similar a como lo hace la opción de consultas.

## 6. Bibliografía

AGIRRE, E., ARREGI, X., DÍAZ DE HARRAZA, A. AND SARASOLA, K. (1994). "Conceptual Distance and Automatic Spelling Correction" in *Workshop on Speech recognition and handwriting*. Leeds, England.

AGIRRE, E., RIGAU, G. (1995). "A Proposal for Word Sense Disambiguation using conceptual Distance", *International Conference on Recent Advances in Natural Language Processing*, Tzgov Chark, Bulgaria.

AGIRRE, E. AND MARTÍNEZ, D. (2001). "Learning class-to-class selectional preferences". In *Proceedings of the 2001 workshop on Computational Natural Language Learning*, Vol. 7. Association for Computational Linguistics, Stroudsburg, PA, USA, Article 3, 8 pages.

ALMELA, M. (2011). "The case for verb-adjective collocations: corpus-based analysis and lexicographical treatment". *Revista de Lingüística y Lenguas Aplicadas*, Vol. 6, págs. 39-52.

ALONSO RAMOS, M. (1994-1995) "Hacia una definición del concepto de colocación: De J. R. Firth a I. A. Mel'čuk", en *Revista de Lexicografía* Vol. I, págs. 9-28.

ALONSO RAMOS, M. (2002). "Colocaciones y contorno de la definición lexicográfica", en *Lingüística Española Actual* XXIV/1, págs. 63-96.

ALONSO RAMOS, M. (2002). "Diccionario de Colocaciones del Español, DICE". [Internet] [Consultado en: 2-12-2013]. Disponible en: <http://www.dicesp.com/paginas>  
ISBN: 978-84-693-9869-2

AGGARWAL, C. (2013), *Outlier Analysis*, Springer Science & Business Media.

BARGALLÓ ESCRIVÁ, M., CARAMÉS DÍAZ, J., FERRANDO ARAMO, V., MORENO VILLANUEVA, J.A., (1997-1998) "El tratamiento de los elementos lexicalizados en la lexicografía española monolingüe", en *Revista de Lexicografía*, Vol. IV, págs. 49-65.

BLASCO MATEO, E. (2002) "La lexicalización y las colocaciones", en *Lingüística Española Actual* XXIV/1, págs. 63-96.

BOSQUE, I. (2001) "Sobre el concepto de colocación y sus límites", en *Lingüística Española Actual* XXIII/1, págs. 9-40.

- BOSQUE, I. (2004). *REDES Diccionario combinatorio del español contemporáneo*. Ediciones SM.
- BRAUN, E. (1996) "Caos, Fractales y Cosas Raras". *Fondo de Cultura Económica*. México.
- BUDANISTSKY, A. (1999), "Lexical semantic relatedness and its application in Natural Language Processing". *Technical Report CSRG-390*, Computer Systems Research Group, University of Toronto.
- CASARES, J., (1999) *Diccionario Ideológico de la Lengua Española*. Editorial Gustavo Gili, SA.
- CASTILLO CARBALLO, M.A. (1997-1998), "El concepto de Unidad Fraseológica" en *Revista de Lexicografía*, Vol. IV, págs. 67-79.
- CASTILLO CARBALLO, M. A. (1998), "El término 'colocación' en la lingüística actual" en *Lingüística Española Actual XX/1*, págs. 41-54.
- CASTILLO CARBALLO, M. A. (2001), "Colocaciones léxicas y variación lingüística: implicaciones didácticas", en *Lingüística Española Actual XXIII/1*, págs. 133-143.
- CHURCH, K. W. AND HANKS, P. (1990), "Word association norms, mutual information, and lexicography". *Computational. Linguistic*. 16, 1 págs. 22-29.
- CORPAS PASTOR, G. (1996). *Manual de Fraseología española*, Madrid, Gredos.
- CORPAS PASTOR, G. (2001), "Apuntes para el estudio de la colocación" en *Lingüística Española Actual XXIII/1*, págs. 41-56
- CORPAS PASTOR, G. (2003), "Diez años de investigación en fraseología: análisis sintáctico-semánticos, contrastivos y traductológicos". *Lingüística Iberoamericana*, Vol. 20, Iberoamericana – Vervuert.
- Diccionario de sinónimos y antónimos*, (2005) Espasa-Calpe. [Internet] [Consultado en: 2-12-2011] en: <http://www.wordreference.com/sinonimos/>
- ALVAR EZQUERRA, M.(DIR.) (1998) *Diccionario Ideológico de la Lengua Española*, 2ª ed. VOX-Bibliograf.
- EVERT, S. (2005). "The statistics of word cooccurrences. Word pairs and collocations".

---

Dissertation, Stuttgart University.

GARCÍA PLATERO, J. M. (2002) “Aspectos semánticos de las colocaciones”, en *Lingüística Española Actual XXIV/1* 2002, págs. 25-34.

GARCÍA-PAGE, M. (2001), “El adverbio colocacional”, en *Lingüística Española Actual XXIII/1* 2001, págs. 89-106.

JORGE-BOTANA, G. (2010). “La técnica del Análisis de la Semántica Latente (LSA/LSI) como modelo informático de la comprensión del texto y el discurso. Una aproximación distribuida al análisis semántico”. Tesis doctoral. Universidad Autónoma de Madrid.

JORGE-BOTANA, G., LEÓN, J. A., OLMOS, R., HASSAN-MONTERO, Y. (2010). “Visualizing polysemy using LSA and the predication algorithm”. *Journal of the American Society for Information Science and Technology*. Vol. 61(8) , pages1706-1724

LIU, H., SHAH, S. & JIANG, W. (2004). "On-line outlier detection and data cleaning". *Computers & chemical engineering*, Vol. 28(9), pages 1635-1647.

HIGUERAS GARCÍA, M. (2006), “Estudio de las colocaciones léxicas y su enseñanza en español como lengua extranjera” en *ASELE, Colección monografías*, nº 9, Málaga

KOIKE, K. (2001), *Colocaciones léxicas en español*, Universidad de Alcalá, Takushoku University.

KOZIMA, H. AND FURUGORI, T. (1993). “Similarity between words computed by spreading activation on an English dictionary”. In *Proceedings of the sixth conference on European chapter of the Association for Computational Linguistics*, pages 232-239.

LANCIA, F. (2007). “Word co-occurrence and similarity in meaning”. In Salvatore, S., Valsiner, J. (Eds.) *Mind as infinite dimensionality*. Edizioni Carlo Amore, Rome.

LANDAUER, T. K., FOLTZ, P.W., LAHAM, D. (1998) “An introduction to latent semantic analysis”. *Discourse Processes*, 25(2-3), pages 259-284.

LIN, D. (1998) “An information-theoretic definition of similarity”. In *Proceedings of the 15th international conference on Machine Learning*, Vol. 1, pages 296-304.

LEYS, C., LEY, C., KLEIN, O., VERNARD, P., LICATA, L. (2013) "Detecting outliers: Do not use standard deviation around the mean, use deviation around the median". *Journal of Experimental Social Psychology*, 49, pages. 764-766.

LIN, F. AND COHEN, W. W. (2010). A very fast method for clustering Big Text Datasets. In Helder Coelho, Rudi Studer, and Michael Wooldridge (Eds.) *Proceedings of the 2010 conference on ECAI 2010: 19th European Conference on Artificial Intelligence*. IOS Press, pages 303-308.

MANNING, C. D. AND SCHÜTZE, H. (Eds.) (1999) *Foundations of Statistical Natural Language Processing*, MIT press, pages 294-300.

MANOJ, K., SENTHAMARAI KANNAN, K. (2013) "Comparison of methods for detecting outliers". *International Journal of Scientific & Engineering Research*, Vol. 4, (9),

MCCALLUM, A., NIGAM, K. AND UNGAR L. H. (2013). "Efficient clustering of high-dimensional data sets with application to reference matching". In *Proceedings of the sixth ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, New York, NY, USA, pages 169-178.

MORENO SANDOVAL, A. (1998) *Lingüística computacional*. Colección "Textos de Apoyo", Editorial Síntesis.

MORRIS, J. AND HIRST, G. (1998). Lexical cohesion computed by thesaural relations as an indicator of the structure of text. *Computational linguistics*, 17(1), pages 21-48.

OLMOS, R., JORGE-BOTANA, G., LEÓN, J.A, AND ESCUDERO, I. (2013). "Giving an interpretation for the semantic dimensions in Latent Semantic Analysis". In *Proceedings of the twenty-third Annual Meeting of the Society for Text and Discourse* Valencia, pages 16-18.

OLMOS, R., JORGE-BOTANA, G., LEÓN, J.A, AND ESCUDERO, I. (2014) "Transforming Selected Concepts Into Dimensions in Latent Semantic Analysis". *Discourse Processes*, 51(5-6), pages 494-510.

PEARCE, D. (2002). "A comparative evaluation of collocation extraction techniques". In *Proceedings of The Third International Conference on Language Resources and Evaluation*, Las Palmas de GC, Canary Islands, pages 1530-1536.

PENADÉS MARTÍNEZ, I. (2001), "¿Colocaciones o locuciones verbales?", en *Lingüística*



---

*Española Actual* XXIII/1, págs. 57-88.

QUASTHOFF, U. AND WOLFF, C. (2002) "The Poisson Collocation Measure and its Applications". In *Second International Workshop on Computational Approaches to Collocations*, 22/23 7 2002, Wien.

R version 2.15.0 (2012) Copyright © 2012. The R foundation for Statistical Computing.

REAL ACADEMIA ESPAÑOLA, (1992) *Diccionario de la Lengua Española*, Vigésima primera edición, Espasa Calpe.

RESNIK, P. (1993), "Selection and information a class-base approach to lexical relationships", *University of Pennsylvania Institute for Research in Cognitive Science Technical Report No. IRCS-93-42*.

RESNIK, P. (1999), "Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language". In *Journal of Artificial Intelligence Research* 11, pages 95-130.

RODRÍGUEZ HONTORIA, H. (2003), "Similitud Semántica" en Martí Antonín, M.A., Fernández Montraveta, A., Vázquez García, G.(Editoras) *Lexicografía Computacional y Semántica*. Colección UB 64, Ediciones de la Universidad de Barcelona.

SECO, M., ANDRÉS, O., RAMOS, G., (1999). *Diccionario del Español Actual*, Aguilar.

SERRA SEPÚLVEDA, S., (2009). "Las restricciones de selección en los diccionarios generales de lengua española". *Boletín de Filología de la Universidad de Chile*, tomo XLIV/ 2, págs.. 187-213.

SCHÜTZE, H. (1993). "Part-of-speech induction from scratch". In *Proceedings of the 31st annual meeting on Association for Computational Linguistics (ACL '93)*. Association for Computational Linguistics, Stroudsburg, PA, USA, 251-258.

PHAM-GIA T. AND HUNG T.L., (2001) "The mean and median absolute deviations", *Mathematical and Computer Modelling*, 34(7), pages 921-936.

REAL ACADEMIA ESPAÑOLA: Banco de datos (CREA) [en línea]. *Corpus de referencia del español actual*. <<http://www.rae.es>> [7/06/2011]

VINCZE O., ALONSO RAMOS M. (2013). "Incorporating Frequency Information in a Collocation Dictionary: Establishing a Methodology". *Procedia-Social and Behavioral Sciences*, Vol. 95, pages 241-248.

VAREAL, F., KUBARTH, H. (1994), *Diccionario Fraseológico del Español Moderno*. Gredos, Madrid.

ZULUAGA, A. (2002) "Los «enlaces frecuentes» de María Moliner. Observaciones sobre las llamadas colocaciones". *Lingüística Española Actual XXIV/1*, págs. 97-114.

## Anexo

Forma Canónica	Cabecera	A(palpitar, G)
agallas, arrestos, <b>corazón</b> , hígados, pecho, redaños	valentía	0,014154
aleteo, <b>contracción</b> , diástole, <b>golpe</b> , <b>latido</b> , palpitación, <b>agitación</b> pulsación, <b>pulso</b> , salto, <b>sístole</b>		0,01233356
<b>benevolencia</b> , <b>corazón</b> , <b>entraña</b> , filantropía, humanidad, <b>benevolencia</b> ternura		0,01185943
anular, <b>auricular</b> , <b>corazón</b> , cordial, dátíl, <b>dedo</b> , <b>índice</b> , <b>medio</b> , <b>meñique</b> , <b>pulgar</b>	dedo	0,01134058
albedrío, arbitrio, <b>corazón</b> , <b>elección</b> , <b>facultad</b> , <b>intención</b> , <b>libertad</b> , libre albedrío, <b>merced</b> , opción, <b>potestad</b> , <b>sentido</b> , <b>voluntad</b>	voluntad	0,01117997
afección, <b>corazón</b> , entretelas, podredumbre, psicología, <b>sentimiento</b> , <b>sentir</b>	sentimiento	0,01093138
<b>corazón</b> , <b>entraña</b> , ombligo, riñón, <b>vientre</b> , <b>yema</b>	centro	0,01039362
alborozo, <b>alegría</b> , alegrón, buen humor, <b>dicha</b> , euforia, <b>felicidad</b> , flash, gozo, <b>humor</b> , júbilo, <b>placer</b> , regocijo	alegría	0,01022647
<b>canción</b> , canto, leitmotiv, <b>matarile</b> , <b>motivo</b> , ritornelo, sintonía, solfa, unísono	canción	0,00904735
buche, cogollo, coletó, <b>cuerpo</b> , <b>entraña</b> , <b>medula</b> , médula, meollo, molla, núcleo, <b>pecho</b> , tuétano	esencia	0,00901596
<b>ánimo</b> , designio, <b>disposición</b> , <b>idea</b> , <b>intención</b> , <b>intento</b> , lineamiento, <b>mente</b> , <b>plan</b> , predestinación, predisposición, <b>propósito</b> , <b>proyecto</b> , talante, <b>vena</b> , <b>vista</b>	intención	0,00854951
<b>entrecejo</b> , <b>frente</b> , <b>sien</b> , sobreceja, testera, testuz	cabeza	0,00843233
batolito, brecha, conglomerado, crestón, eluvión, estrato, <b>filón</b> , saca, salbanda, <b>vena</b> , veta	suelo	0,00811938
arranque, arrechucho, basca, capricho, corazonada, <b>ida</b> , <b>impulso</b> , llamada, llamarada, <b>pronto</b> , <b>rapto</b> , <b>vena</b> , venida, ventolera, viaraza	determinismo	0,00806875
caída, cuesta, declive, dirección, grada, gradiente, medianil, <b>pecho</b> , <b>pechuga</b> , pendiente, rampa, ranfla, rasante, <b>repecho</b> , reventón, ribazo, salmonera, subida, talud, varga, vertiente	inclinación	0,00766294
<b>banda</b> , cadeneta, <b>cinta</b> , faja, franja, fres, friso, <b>goma</b> , <b>lista</b> , <b>pieza</b> , serpentina, <b>tira</b> , <b>vena</b> , veta	anchura	0,00760636

<i>altura, apotema, base, calibre, cateto, cosecante, coseno, línea</i>		0,00757381
<i>cuerda, diagonal, diámetro, eje, flecha, hipotenusa, lado, luz,</i>		
<i>mediana, mediatriz, módulo, radio, sagita, secante, seno</i>		
<b>latido, tictac</b>	<i>sonido</i>	0,00753423
<i>agrado, amenidad, delectación, deleite, delicia, dulzura, placer</i>		0,00704114
<i>fruición, gloria, goce, gozo, gusto, placer, regalo, sabrosura,</i>		
<i>satisfacción</i>		
<b>arteria, capilar, hacecillo, vaso, vena</b>	<i>vena</i>	0,00659664

Tabla 6-1 Los 25 grupos con mayor Asociación de Selección para el verbo palpar.

<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(brusco, G)</b>
<i>arribismo, electoralismo, maquiavelismo, maroma, oportunismo, secesión, secesionismo, transfuguismo, viraje, virazón</i>	<i>deshonestidad</i>	0,0438858
<i>ahogadilla, chapuzón, inmersión, sumersión, viraje, zambullida</i>	<i>dentro</i>	0,0380421
<i>bolina, bordada, ceñida, guiñada, huida, revuelta, viraje</i>	<i>desviación</i>	0,0339251
<i>insolación, polavisión, revelado, veladura, viraje</i>	<i>cinematografía</i>	0,0316622
<i>deflación, depreciación, devaluación</i>	<i>rebaja</i>	0,0165424
<i>ademanes, cortesía, formas, maneras, modales, modos, perfiles</i>	<i>cortesía</i>	0,0149323
<i>abducción, acarreo, porte, transbordo, transferencia, transfusión, transición, tránsito, transporte, transposición, traspaso, trasplante</i>	<i>transporte</i>	0,0139917
<i>arreada, atraco, golpe, mangoneo, sirla, tirón</i>	<i>robo</i>	0,0127632
<i>conversión, evolución, mutación, paso, salto, transformación, transición, tránsito</i>	<i>cambio</i>	0,0113854
<i>caracol, corvadura, curva, dobladura, entrada, meandro, mocheta, recodo, recoveco, rincón, sucucho</i>	<i>ángulo</i>	0,0112814
<i>concusión, sacudida, vaivén</i>	<i>agitación</i>	0,0111638
<i>ademán, aspaviento, batimán, coco, coquito, déxis, esguince, gesto, guiño, guiñada, jeribeque, llamada, mohín, morra, mueca, muestra, musaraña, rictus</i>	<i>gesto</i>	0,0111396
<i>acometida, agresión, apretón, arrancada, asalto, ataque, atentado, carga, contraataque, embate, irrupción, viaje</i>	<i>ataque</i>	0,0104063
<i>agarrón, jalón, orejón, repelón, tirón</i>	<i>atracción</i>	0,0100942
<i>acometida, bisagra, costura, empalme, engranaje, enlace, entronque, escarpe, inglete, juntura, traba</i>	<i>unión</i>	0,0093489
<i>alteración, cambio, evolución, innovación, modificación, muda, mudanza, reformismo, revés, sensibilidad, variación, variante, variedad, vuelta</i>	<i>cambio</i>	0,0093489
<i>circulación, desplazamiento, jubileo, locomoción, movimiento, paso, telequinesia, tráfico, tránsito, traslación, traslado, trayecto</i>	<i>tránsito</i>	0,0089822
<i>alboroto, alzamiento, asonada, cataclismo, comunidades, conmoción, contrarrevolución, convulsión, cuartelazo, golpe de estado, inquietud, insurrección, levantamiento, militarada, motín, movimiento, pronunciamiento, rebelión, revolución, revuelta, secesión, sedición, terrorismo, traición, tumulto</i>	<i>desobediencia</i>	0,0087821
<i>deducción, ilación, inducción, transición</i>	<i>razonamiento</i>	0,0087006

<i>acción, anagnórisis, anticlímax, argumento, clímax, desenlace, mensaje, enredo, episodio, golpe, intriga, jornada, máquina, movimiento, nudo, peripecia, perístasis, prótasis, solución, trama</i>	<i>0,0086801</i>
<i>balotada, batido, batuda, botepronto, brinco, cabriola, salto, corcovo, jumping, pedico, rebote, saltadero, salto, salto trenzado</i>	<i>0,0082866</i>
<i>aerobic, atletismo, calistenia, cros, culturismo, decatón, deporte, footing, gimnasia, gimnasia rítmica, halterofilia, joggin, lanzamiento, olimpismo, pedestrismo, pentatlón, salto, tiro, yoga</i>	<i>0,008161</i>
<i>evolución, marcha, meneo, moción, motilidad, móvil, movimiento, movilidad, movilización, movimiento, transmisión</i>	<i>0,0071675</i>
<i>atolondramiento, aturdimiento, azoro, confusión, conmoción, intranquilidad, corte, movimiento, tupición, turbación</i>	<i>0,0069941</i>
<i>alocución, arenga, declamación, manifestación, mitin, difusión, movimiento, proclama</i>	<i>0,0066185</i>

Tabla 6-2 Grupos de sustantivos con los valores máximos de Asociación de Selección para el adjetivo brusco.

<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(doloroso, G)</b>
<i>aguijonada, alfilerazo, clavadura, cornada, espoleadura, estacada, filazo, mojada, picadura, pinchazo, piquete, puntada, puntazo, puntura, punzada, repique</i>	<i>herida</i>	0,0195348
<i>ay, cojijo, gemido, guaya, lamentación, lamento, lástima, queja, quejido, querella, suspiro</i>	<i>lamento</i>	0,0182051
<i>acibar, aflicción, agonía, agrá, agraz, ahogo, amargura, angustia, apuro, congoja, desconsuelo, dolor, duelo, hiel, pasión, pena, pesar, presura, quebranto, sentimiento, sufrimiento</i>	<i>sufrimiento</i>	0,0161039
<i>acritud, analgesia, clavo, dolor, puntada, punzada, ramalazo, rayo, sinalgia, tormento, yaya</i>	<i>dolor</i>	0,0131621
<i>arestín, desazón, desplacer, disgusto, martelo, patatús, pensión, pesadumbre, punzada, puñalada, sinsabor, tribulación, varapalo</i>	<i>sufrimiento</i>	0,0130223
<i>aleteo, contracción, diástole, golpe, latido, palpitación, pulsación, pulso, salto, sístole</i>	<i>agitación</i>	0,0092431
<i>aféresis, apócope, contracción, epéntesis, metaplasmo, metátesis, paragoge, prótesis, síncope</i>	<i>representación</i>	0,0092057
<i>afán, agonía, anhelo, ansia, apetito, concupiscencia, deseo, desiderátum, empeño, emulación, envidia, gana, golondro, hambre, hipo, pelota, pelusa, pío, pique, prurito, pujo, rabanillo, reconcomio, sed, sueño dorado, voto, yesca</i>	<i>deseo</i>	0,0090338
<i>cortadura, corte, fraile, gargantil, incisión, marceo, mediacaña, picadura, rafa, raja, rumbo, sajadura, siete, tajo</i>	<i>corie</i>	0,0083892
<i>adormecimiento, anquilosis, artritis, artritismo, artrosis, calambre, contractura, convulsión, espasmo, gota, miopatía, osteopatía, raquitismo, reuma, reúma, reumatismo, rigor, risa sardesca, risa sardonía, risa sardónica</i>	<i>reumatismo</i>	0,0080548
<i>cucaracha, macuba, manila, palillo, picadura, rapé, tabaco, tirulo, tripa</i>	<i>veneno</i>	0,00789266
<i>boquerón, boquete, brecha, calada, gotera, magaña, picadura, pinchazo, punto, rasa, silbato, tomate</i>	<i>agujero</i>	0,00773525
<i>astringencia, contracción, contractilidad</i>	<i>disminución</i>	0,00727213
<i>cariadura, caries, neguijón, picadura</i>	<i>degeneración</i>	0,00701377
<i>aclamación, alarido, baladro, berrido, bramido, chillido,</i>	<i>voz</i>	0,00684484

<i>clamor, exclamación, grito, ladrido, rugido, voz</i>	
<i>arañazo, carnicería, chispazo, colisión, corte, herida</i>	0,00677393
<i>dehiscencia, fístula, grieta, herida, perforación, pupa, quemadura, rasponazo, rozadura, tocadura, transfixión, transverberación, yaya,</i>	
<i>admiración, asombro, consternación, escándalo, espanto, asombro</i>	0,00663455
<i>estupor, extrañamiento, extrañeza, flash, golpe, marasmo, novedad, pasmo, repullo, sorpresa, suspensión</i>	
<i>arcada, contorsión, convulsión, espasmo, jeribeque, agitación</i>	0,00641441
<i>socollón, tic, torozón</i>	
<i>agitación, agobio, alarma, ansia, ansiedad, batalla, intranquilidad</i>	0,00627659
<i>combate, comecome, comezón, conflicto, desabrimiento, desasosiego, estrés, guerra, hormigueo, inquietud, intranquilidad, malestar, mar de fondo, maretta, pesadilla, polvareda, preocupación, rebato, repunta, revuelo, susidio, suspense, taco, tártago, torozón, zozobra</i>	

Tabla 6-3 Grupos de sustantivos con los valores máximos de Asociación de Selección para el adjetivo doloroso.



<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(abundantemente, G)</b>
<i>alargar, dar, dotar, proporcionar, proveer, servir, suministrar, surtir</i>	<i>suministro</i>	0,1102675
<i>brotar, destilar, esfluir, emanar, exudar, fluir, instilar, manar, resurgir, revenir, rezumar, saltar, sudar, surgir, surtir, trasvinar</i>	<i>fuera</i>	0,0998385
<i>aparejar, apercibir, aprestar, aprevenir, aprontar, armar, aviar, avisar, conrear, disponer, estructurar, formatear, intentar, mullir, organizar, parar, perdigar, pertrechar, poner, preconcebir, preparar, prevenir, prevenir, proveer, reservar, sembrar</i>	<i>preparación</i>	0,0961376
<i>asolear, azacananar, azacanear, batallar, bregar, insudar, luchar, navegar, pegar, pernear, pringar, remar, sudar, trabajar, velar, zarandearse</i>	<i>actividad</i>	0,0500253
<i>llorar, resinar, sangrar, sudar</i>	<i>savia</i>	0,0466727
<i>agarrar, atrapar, cachar, cazar, empuñar, ligar, pescar, sacar, sudar</i>	<i>éxito</i>	0,0396564
<i>asperjar, espurrear, espurriar, irrigar, nebulizar, perlar, petroleo, pulverizar, regar, rociar, salpicar, vaporizar</i>	<i>dispersión</i>	0,0390331
<i>resudar, sudar, transpirar, trasudar</i>	<i>sudor</i>	0,0383897
<i>ensalivar, sudar, trasudar</i>	<i>humedad</i>	0,0383897
<i>enristrar, entramar, interceptar, regar, trabar, tramar, tranzar, trenzar</i>	<i>intersección</i>	0,0334439

Tabla 6-4 Grupos con mayor asociación de Selección para el adverbio abundantemente.

<b>Forma Canónica</b>	<b>Cabecera</b>	<b>A(enérgicamente, G)</b>
<i>acariciar, amasar, arrastrar, barrer, cardar, emborrizar, emprimir, estregar, fregar, fricar, friccionar, frisar, frotar, ludir, malaxar, masajear, mesar, peinar, rasar, rascar, raspar, rastrear, refregar, restregar, rozar, sobar, tazar, tropezarse</i>	<i>tangencia</i>	0,0388546
<i>aventurar, hacer muestras, motivar, parir, protestar</i>	<i>comunicación</i>	0,0356015
<i>chillar, encerrarse, manifestar, piar, protestar</i>	<i>oposición</i>	0,0342641
<i>acusarse, condenarse, confesar, protestar, reconciliarse, reconocerse</i>	<i>culpabilidad</i>	0,0275147
<i>agitar, cimbrar, menear, sacudir, traquear, traquetear, ventilar, vibrar, zamarrear, zangolotear, zarandar, zarandear</i>	<i>agitación</i>	0,026274

<i>ahuyentar, deponer, descartar, desechar, despedir, rechazo</i>	0,0257924
<i>desterrar, negar, rechazar, recusar, reherir, repeler, repudiar, repugnar, resistir</i>	
<i>confutar, contradecir, desmentir, negar, rebatir, rechazar, refutación</i>	0,0247705
<i>rectificar, refutar, repeler, resistir, triturar</i>	
<i>amordazar, atajar, atarse, atascar, contrariar, dificultad</i>	0,0201045
<i>embarazar, embargar, embazar, empachar, empatar, empezar, entorpecer, estorbar, estrangular, impedir, implicar, imposibilitar, incapacitar, obstruir, oponer, perturbar, reaccionar, trabar, tronchar, tropezar, vedar, yugular</i>	
<i>afrontar, arrostrar, contrarrestar, contrastar, desafiar, ataque</i>	0,018998
<i>enfrentarse, rechazar, recibir, reparar, resistir, revolver</i>	
<i>acaudillar, acudrillar, adiestrar, administrar, conducir, control</i>	0,0189946
<i>controlar, dirigir, encabezar, encaminar, gobernar, guiar, intervenir, maestrear, manejar, menear, normar, orquestar, patronear, presidir, regentar, regir, seguir</i>	

Tabla 6-5 Grupos con mayor Asociación de Selección para el adverbio enérgicamente.