Taller sobre herramientas de análisis textual:

La herramienta Sketch Engine

Facultad de Filología, Universidad Complutense de Madrid

18/02/2019

Autores:

Iván Arias Rodríguez (iarias01@ucm.es)

Ana Fernández Pampillón Cesteros (apampi@filol.ucm.es)

Doaa Samy (doaa_samy@hotmail.com) (doaasamy@cu.edu.eg)

Jorge Arús Hita (jarus@filol.ucm.es)

Objetivo:

El objetivo es proporcionar una introducción al uso de la herramienta *Sketch Engine*. El tiempo de aprendizaje previsto es de 2-3 horas en una única sesión o dos sesiones de hora y media.

Este documento es una actualización de:

Arias Rodríguez, Iván y Samy, Doaa y Fernández-Pampillón Cesteros, Ana María y Arús Hita, Jorge (2017) *Taller sobre herramientas de análisis textual: La herramienta Sketch Engine*. <u>https://eprints.ucm.es/46295/</u>

Sami, D.; Fernández-Pampillón, A.; Arús, J. (2011) "Taller sobre herramientas de análisis textual: la herramienta *Sketch Engine*". Disponible en: <u>http://eprints.ucm.es/13796</u>



This work is licensed under a <u>Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License</u>.

Contenido

| 1 INTRODUCCIÓN | 3 |
|---|----|
| 1.1 LAS HERRAMIENTAS DE ANÁLISIS TEXTUAL: LOS CORPUS | |
| 1.1.1 - Criterios para recolectar los datos: resumen de tipología de corpus | |
| 1.2 ¿Qué es Sketch Engine? | 5 |
| 1.3 ¿CÓMO REGISTRARSE Y ACCEDER A LA CUENTA DE <i>SKETCH ENGINE</i> ? | 5 |
| 2 SELECCIÓN DE UN CORPUS | 9 |
| 3 CONSTRUIR UN CORPUS CON SKETCH ENGINE | 14 |
| 3.1 CREAR UN CORPUS A PARTIR DE ARCHIVOS DE TEXTO LOCALES | |
| 3.1.1 - Corpus monolingües | |
| 3.1.2 - Corpus multilingües | |
| 3.2 CREAR UN CORPUS AUTOMÁTICO A PARTIR DE TEXTO DESCARGADO DE INTERNET | |
| 3.2.1 - Uso de palabras semilla | |
| 3.2.2 - Uso de URLs concretas | |
| 3.2.3 - Uso de un Website | |
| 3.3 División en subcorpus | |
| 4 EXPLORAR EL CORPUS | 24 |
| 4.1 Word Sketch: Relaciones gramaticales | |
| 4.2 DIFERENCIA SKETCH: DIFERENCIAS DE USO ENTRE PALABRAS | |
| 4.3 Tesauro: Creación de un tesauro automático | |
| 4.4 CONCORDANCIA: EJEMPLOS DE USO EN CONTEXTO | |
| 4.4.1 - Tipos de consultas de concordancia | |
| 4.4.1.1 - Concordancia simple | |
| 4.4.1.2 - Concordancia de lema | |
| 4.4.1.3 - Concordancia de frase (sintagma) | |
| 4.4.1.4 - Concordancia de forma | |
| 4.4.1.5 - Concordancia de caracteres | |
| 4.4.1.6 - Concordancia CQL | |
| 4.5 <i>LISTAS DE PALABRAS</i> : FRECUENCIAS DE APARICION | |
| 4.6 <i>N-GRAMAS</i> : EXTRACCION DE EXPRESIONES MULTIPALABRA | |
| 4. / PALABRAS CLAVE: EXTRACCION DE PALABRAS CLAVE Y TERMINOS MULTIPALABRA | |
| 4.8 UNELLICK DICTIONARY: CREACIÓN DE UN DICCIONARIO AUTOMÁTICO | |
| 4.9 CORCONDANCIA PARALELA: EXPRESIONES EQUIVALENTES EN DOS LENGUAS | |
| 4.10 Tendencias: Variación en el uso de palabras a lo largo del tiempo | 45 |



1. - Introducción

1.1. - Las herramientas de análisis textual: los corpus

Las herramientas de análisis textual permiten el tratamiento automático de la información para apoyar el **estudio empírico de las lenguas**. Permiten la consulta rápida de una o varias **colecciones de textos** electrónicos, así como la preparación y el mantenimiento de bases de datos textuales.

Los textos pueden analizarse desde dos niveles de complejidad:

- El primer nivel, el nivel de datos, es el nivel más simple desde el punto de vista operacional. El objeto de análisis es el texto puro, entendido como un conjunto de caracteres. Las operaciones automáticas que pueden realizarse son, únicamente, aquellas basadas en la manipulación de los caracteres del texto. Por ejemplo, la localización de todas las palabras que comienzan por un prefijo determinado, las palabras que acompañan a otra dada (colocaciones) o la frecuencia de aparición de un término.
- 2. En un segundo nivel de complejidad, el nivel de información se corresponde con la interpretación de los textos. A este nivel corresponden las operaciones que necesitan disponer, además de los textos, de la interpretación de los mismos, como por ejemplo la localización de todos los verbos transitivos en una colección de textos o los textos que tratan sobre "la desaparición del atún". Este tipo de consultas más "inteligentes" requieren que el texto tenga marcado de alguna forma este tipo de información. Es importante tener en cuenta que el procesamiento automático de la información es más complicado que el procesamiento de datos y que solo puede realizarse si previamente se ha preparado (pre-procesamiento) el texto (datos) insertando la interpretación (semántica) de estos datos textuales. Uno de los mecanismos de pre-procesamiento es la inserción en el texto de marcas o etiquetas explícitas con la información asociada a cada elemento textual (proceso de etiquetado).

Todas las herramientas de análisis textual permiten el análisis de texto puro (datos), y sólo algunas, como *Sketch Engine*, ofrecen la posibilidad de analizar también el **texto marcado** (información). Además, *Sketch Engine* dispone de sus propias herramientas para el etiquetado automático del texto, con lo que es capaz de preprocesar el texto puro, convirtiéndolo en texto marcado.

Para trabajar con herramientas de análisis textual hace falta una **colección de textos** representativos de aquello que se desea analizar. A esta colección representativa de textos se le llama un **corpus**. Para un estudio empírico bien fundamentado y unos resultados significativos, es necesario que la colección sea lo suficientemente **grande y equilibrada** como para asegurar que abarca las máximas ocurrencias posibles del fenómeno a estudiar. En este sentido se pueden distinguir varias tipologías de corpus.

1.1.1 - Criterios para recolectar los datos: resumen de tipología de corpus

Se puede establecer una **tipología** de corpus¹ en función de los **criterios** utilizados para la clasificación. Estos criterios pueden ser de diferente índole:

- Según la **modalidad** de la lengua.
 - o Corpus de lengua escrita.
 - Corpus de lengua hablada.
 - Transcripciones ortográficas de grabaciones, utilizadas sobre todo en lingüística de corpus.
 - Grabaciones y transcripciones ortográficas y/o fonéticas, empleadas en fonética y tecnologías del habla.
 - o Corpus mixtos.
- Según la naturaleza física de los datos.
 - Corpus **textuales** o escritos: compuestos por una colección de textos.
 - o Corpus orales: compuestos por una colección de audios.
 - Textos **multimodales** (imagen/video/audio y texto).
- Según la cobertura y la **temática**.
 - Por períodos.
 - o Por géneros.
 - o Por temas.
 - o Por autores.
 - Por registros.
- Según el **número de lenguas**.
 - Corpus **monolingües**: están formados por textos de una sola lengua. Se recopilan con el objetivo de dar cuenta de **una lengua o variedad lingüística**.
 - Corpus bilingües o multilingües: están formados por textos de dos (bilingües) o más (multilingües) lenguas sin que, en principio, sean traducciones unos de otros y sin compartir criterios de selección.
 - Corpus comparables ("paired texts"): consisten en una selección de textos en más de una lengua o variedad lingüística parecidos en cuanto a sus características y que comparten criterios de selección. Se utilizan sobre todo para comparar variedades de la lengua en estudios contrastivos.
 - Corpus paralelos ("bi-texts"): recogen textos en más de una lengua (bilingües o multilingües) pero, a diferencia de los anteriores, se trata del mismo texto traducido a una o más lenguas.
- Según la disponibilidad y el modo de recopilación.
 - o Corpus disponibles como recursos en instituciones o en grupos de investigación.
 - Corpus construidos a partir de textos seleccionados y **recopilados manualmente** según criterios concretos y para estudiar unos aspectos concretos.

¹ La palabra *corpus* es invariante en cuanto al número (es decir, su plural es igualmente *corpus*). No debe confundirse con la palabra *corpora*, que es el plural de la palabra inglesa *corpus* (que mantiene el plural latino).



 Corpus recopilados automáticamente por herramientas que permiten una búsqueda y recopilación a través de internet mediante el uso de palabras semilla².

1.2. - ¿Qué es Sketch Engine?

Sketch Engine es una **herramienta de análisis textual en línea** que recibe como entrada un **corpus en cualquier idioma** con, posiblemente, un cierto nivel de **anotación lingüística** para su posterior análisis. *Sketch Engine* ofrece múltiples funciones para el análisis lingüístico:

- Análisis de colocaciones gramaticales y léxicas (*Word Sketch*): Busca las apariciones en el corpus de ciertas estructuras gramaticales que contienen a una palabra/lema dado, así como coapariciones de términos.
- 2) Además permite hacer un contraste entre los resultados para las colocaciones de dos palabras (*Diferencia Sketch*), para apreciar las diferencias entre ambas.
- Tesauro: Crea automáticamente un conjunto de palabras relacionadas con una palabra dada (campo semántico).
- 4) **Concordancia**: Permite buscar contextos en los que aparecen ciertas (combinaciones de) palabras (formas flexionadas) y/o lemas.
- 5) Es posible también realizar una *concordancia en paralelo* en corpus multilingües para hallar los equivalentes a palabras y expresiones entre una lengua y otra.
- 6) *Lista de palabras*: Realiza conteos de cantidad de apariciones y frecuencias de las palabras/lemas que aparecen en el corpus.
- 7) Obtención de los *n-gramas* (expresiones multipalabra) más característicos del corpus.
- 8) Extracción de **palabras clave** y términos multipalabra: Se extraen los términos más representativos del texto (en comparación con un corpus de referencia), así como los términos multipalabra más característicos del corpus.
- 9) **OneClick Dictionary**: Crea un esbozo de diccionario a partir de un corpus y las relaciones sintácticas de sus palabras.
- 10) Análisis de *tendencias* en el uso de palabras según la fecha de publicación.

Sketch Engine permite trabajar con los siguientes tipos de corpus:

- Corpus que vienen integrados en la herramienta.
- Corpus disponibles en **instituciones** o en **grupos de investigación** que hay que "subir" previamente a la herramienta.
- Corpus construidos a partir de **textos seleccionados** y recopilados manualmente.
- Corpus **construidos automáticamente** con la herramienta **WebBootCat**, que permite la búsqueda y recopilación automática de documentos en Internet.
- Corpus paralelos.
- Corpus **etiquetados**.

1.3. - ¿Cómo registrarse y acceder a la cuenta de Sketch Engine?

Cualquier persona puede obtener una licencia gratuita temporal con una duración de 30 días. Además, todo aquel que cuente con una dirección de email de la Universidad Complutense puede crearse una cuenta **sin coste alguno hasta marzo de 2022** gracias a la financiación del **proyecto ELEXIS**. Si no se

² Una palabra semilla es una palabra representativa de un tema.

posee una dirección de correo de la UCM, el procedimiento para obtener una cuenta gratuita durante los próximos 30 días es el siguiente:

 Abra la página inicial de Sketch Engine en <u>https://www.sketchengine.eu/</u> y clique en el enlace Register que aparece en la parte superior izquierda de la página tal y como se muestra en la Figura 1.



- 2) Como aparece en la Figura 2, elija en el menú la opción de *Free 30-day trial*. Tras ello, rellene el formulario (puede utilizar cualquier dirección de correo) y clique en *SIGN UP*.
- 3) Al hacer esto, se le enviará la clave a la cuenta de correo que haya utilizado. Además se le redirige automáticamente a la página de *Log in*, que se muestra en la Figura 3. Rellene el campo de la cuenta de correo utilizada para darse de alta (o el nombre de usuario que haya escogido) junto con la clave que se le envió y presione el botón de *LOG IN*. Para futuros usos, puede acceder a la página de *Log in* clicando en el botón azul con el texto *LOG IN* que aparece en la página inicial de *Sketch Engine* (Figura 1).

| SKETCH ENGINE | |
|--|--|
| ign up | Sign up Free 30-day trial subscription |
| | User Name * Nombre_de_usuario_(nickname) |
| The complete functionality, 220+ corpora, 90 languages. May contain advertising. | E-mail * correo@cualquiera.cualquiera |
| • | First Name Nombre |
| Individual user account An academic or commercial use conducted by a single person. | Last Name Apellido |
| | I Will Use Sketch Engine Mainly For Inguistics, social science or humanities 🔻 |
| * | I Work Or Study At university, college, school or similar • |
| Multi-user account An academic or commercial use conducted by an institution. | Country * Spain Q |
| | Organisation Name Universidad Complutense de Madrid |
| Join a multiuser account An access code is required. | ✓ I agree to the Terms of use ☐ Inform me about new corpora and functions ✓ I agree to the processing of personal data |
| BACK | BACK SIGN UP |
| Problem with registration? Contact us at support@sketchengine.co.uk . | Problem with registration? Contact us at support@sketchengine.co.uk. |

| Log in User name or email Password LOG IN Forgot password? | |
|--|--------------------------------|
| Don't have an account? Sign up Or try open corpora. | Back to the original interface |
| Figura 3: página de Log in | |

7

٦

En caso de tener una dirección de correo de la UCM (o de alguna otra universidad incluida en el projecto ELEXIS) se puede registrar como usuario para obtener una **licencia institucional** que será válida hasta marzo de 2022. En este caso, proceda así:

- Abra la página inicial de Sketch Engine en <u>https://www.sketchengine.eu/</u> y clique en el botón azul de LOG IN que puede verse en la Figura 1.
- 2) En la página de *Log in* (Figura 3) clique el botón superior de la derecha que contiene el texto *Institutional login*.
- 3) Se le abrirá una ventana emergente como la que se muestra en la Figura 4. Simplemente escriba el texto *Complutense* en el campo de entrada de texto, escoja la institución UCM Universidad Complutense de Madrid y clique el botón (que cambia de Search a Continue tras seleccionar la UCM). Tras ello, deberá utilizar sus datos de correo de UCM y clave que utilice normalmente y clicar en Iniciar sesión.
- 4) En la siguiente página se le preguntará si quiere reactivar una cuenta antigua que se haya desactivado o si quiere crear una cuenta nueva. Si nunca ha tenido una cuenta en Sketch Engine o si no está seguro pero no le importa que su cuenta se reinicie (si en el pasado hubiera creado su propio corpus, se borraría), clique en NEW ACCOUNT. En caso contrario, clique en REACTIVATE ACCOUNT y siga las instrucciones que le aparezcan en pantalla.
- 5) Seguidamente le aparecerá un formulario como el de la Figura 2 (pero con menos campos) que deberá rellenar de modo similar a como se muestra en dicha figura, para después clicar en *SIGN UP*.
- 6) Tras esto, clique en GO TO SKETCH ENGINE en la ventana emergente para acceder a su cuenta. Además, se le enviará su clave de acceso a su correo institucional de la UCM que deberá utilizar en sucesivos accesos a Sketch Engine (clicando el botón azul de LOG IN en la página inicial y rellenando sus datos de cuenta de correo y clave de acceso).

| | | COMPLUTENSE |
|---|---------------------|--|
| SKETCH | | Acceso Web Unificado |
| ENGINE | | identificarse correctamente en esta página le habilitará la entrada en la mayoría de las aplicaciones y en los servicios en la nube @UCM. |
| | | Dirección de correo UCM |
| e Sketch Engine is a Corpus Query System allowing you to research | h how words behave. | usuario@ucm.es |
| Which organisation would you like to si | an in with? | Contraseña |
| which organisation would you like to si | gir in with: | |
| Complutense UCM - Universidad Complutense de Madrid Sign In | Search | Lolvidó la contraseña? Más información |
| Need help logging in2 | - | |
| The UK Access Management Federation Accessibility statement Privacy and Cookies Policy | | |
| | | |

Figura 4: introducción del nombre de la institución y verificación de usuario de la UCM



2. - Selección de un corpus

Para poder empezar a trabajar con *Sketch Engine*, lo primero que tenemos que hacer es seleccionar un corpus sobre el que trabajar. Así pues, la primera tarea es la de seleccionar alguno de los muchos corpus disponibles o crear un corpus propio. Es importante comprender que, una vez **elegido un corpus**, todas las **operaciones** (búsquedas, concordancias...) que se hagan **se realizarán sobre ese corpus**.

| 59 | BASIC ADVANCED | D MY CORPORA | SUBSCRIBE 🖘 🕜 🖆 |
|--|---|------------------------------------|------------------------------|
| ** | LANGUAGES Select a language and we w | will pick the best corpus for you. | QUICK START TUTORIAL |
| •••••••••••••••••••••••••••••••••••••• | ENGLISH ARABIC GERMAN | ITALIAN RUSSIAN PORTUGUESE | How to start |
| | JAPANESE More languages | | |
| ↓≡ NE | | | in 2 minutes |
| =0 الر | | | Back to the original interfa |

En el primer acceso a *Sketch Engine*, no habrá ningún corpus seleccionado, con lo que se nos redigirá automáticamente a la **página de selección de corpus**, como se muestra en la Figura 5. Como puede observarse en esta figura, la página muestra once iconos en la parte izquierda de la pantalla, que se corresponden con once herramientas disponibles. No obstante, al no haber escogido ni creado aún ningún corpus, solo las dos primeras opciones están disponibles, mientras las nueve siguientes aparecen atenuadas y no se pueden escoger. La primera de dichas herramientas es el **panel de control** (*Dashboard*). Si clicamos en este primer icono, veremos que no nos permite hacer nada más que seleccionar o crear un corpus, como puede observarse en la Figura 6.

Antes de empezar a usar *Sketch Engine*, hay una serie de acciones que podemos llevar a cabo clicando en los cuatro iconos que aparecen en la parte superior derecha de la pantalla. Es importante hacer notar que el botón rojo de *SUBSCRIBE* solo aparece cuando estamos usando una cuenta gratuita válida por 30 días. En caso contrario aparecerá *Get more space (+)*. De derecha a izquierda, dichos iconos son:

| Ø | DASHBOARD | type to search Q | SUBSCRIBE 🖘 🕐 💶 🗄 |
|----------|----------------|-----------------------------------|---|
| | | | RECENTLY USED CORPORA NEW CORPUS Nothing here |
| \odot | 1 | no corpus selected | |
| 00 | | Please, select corpus to start. | |
| *≣ | | SELECT CORPUS | |
| E+E | | | |
| | | | |
| ¥Ξ NΞ | | | Lexicom 2019 ~ |
| δ≣ | | | |
| 13 | RECENT RESULTS | FAVOURITE RESULTS | |
| | Nothing here | | |
| | | | Back to the original interface |
| | | Figura 6: nanel de control de Ske | etch Engine (nrimer acceso) |

- **Configuración**: Tiene tres opciones.
 - a. *Mi cuenta*: Da información sobre nuestro usuario y el tipo de cuenta que tenemos, así como sobre la cantidad de memoria utilizada en nuestros corpus. Por defecto se permite utilizer hasta 1.000.000 palabras, pero dicho espacio puede aumentarse.
 - b. Configuración: Permite cambiar la lengua en la que se muestra Sketch Engine, así como ajustar la densidad de información (si el texto aparece con más o menos espacio blanco alrededor). En adelante, se utilizará Sketch Engine con alta densidad de texto y en español (si bien no toda la herramienta está traducida y seguiremos viendo algunos textos en inglés).
 - c. *Cierre de sesión*: Permite cerrar la sesión abierta en el ordenador. La sesión debe cerrarse al terminar de usar *Sketch Engine* si estamos usándolo en un ordenador compartido.
- **Comentarios**: Permite hacer consultas a *Sketch Engine* cuando tengamos algún problema con su uso.
- **Guía de usuario**: Sketch Engine dispone de varias guías de usuario, tanto en modo texto como video, que permiten a un recién llegado familiarizarse con el uso de la herramienta.
- **Enlaces cortos**: Crea enlaces (URLs) cortos a un cierto corpus o página de *Sketch Engine* para enviárselo más comodamente a otro usuario.
- Obtener más espacio (+): Al clicar aquí se nos da la posibilidad de hacer una petición a Sketch Engine para obtener más espacio en nuestra cuenta y así disponer de más de 1.000.000 palabras de almacenamiento para nuestros corpus. Esta opción solo está disponible para usuarios con una cuenta institucional o de pago. Si de verdad nos hace falta el espacio extra de almacenamiento, podemos utilizar esta opción y disponer de más espacio sin coste extra.

Así pues, lo primero que debemos hacer para empezar a utilizar *Sketch Engine* es seleccionar un corpus sobre el que trabajar. Para ello, debemos ir a la página de selección de corpus, en cuya parte superior podemos ver tres pestañas:

Si nos fijamos en la página de selección de corpus, vemos que arriba muestra tres pestañas:

BÁSICO: Nos permite hacer una selección de un corpus basándose únicamente en la lengua. Si queremos escoger un corpus en una lengua que no sean las siete que aparecen en los botones de la Figura 5, debemos escribir el nombre de la **lengua** (en inglés) en el campo de texto que aparece en la parte inferior izquierda de la pantalla. Si escogemos el español (*Spanish*), se escogerá por defecto el corpus esTenTen³, un corpus etiquetado que contiene casi 10.000.000.000 palabras (es decir, 10¹⁰ palabras, de ahí su nombre). Este corpus está etiquetado automáticamente (sin revisión humana posterior) con el etiquetador morfosintáctico FreeLing⁴, que tiene una precisión de alrededor del 97% de acierto y que asigna etiquetas basadas en las recomendaciones EAGLES⁵.

Una vez hayamos seleccionado un corpus de trabajo, se nos abren ocho de las diez opciones que hasta ahora estaban desactivadas, tal y como se ve en la Figura 7, que muestra la pantalla de panel de control tras haber seleccionado el corpus español por defecto. Estas mismas herramientas aparecen tanto en el menú de iconos de la izquierda, como en el panel principal del panel de control (donde además se muestra la herramienta para la creación de diccionarios automáticos). También aparece un panel de los corpus recientemente escogidos a la derecha, y un panel de resultados en la parte inferior que se verá más adelante.



³ Jakubíček, M., Kilgarriff, A., Kovář, V., Rychlý, P. y Suchomel, V. (Eds.) 2013. The tenten corpus family.

⁴ Padró, L. y Stanilovsky, E. (Eds.) 2012. Freeling 3.0: Towards wider multilinguality.

⁵ Leech, G. y Wilson, A. (1996). EAGLES recommendations for the morphosyntactic annotation of corpora.

AVANZADO: Muestra **todos los corpus** que hay disponibles en *Sketch Engine*. Actualmente existen 492 corpus en 94 lenguas, como puede verse en la Figura 8. La página permite **filtrar** los corpus según la lengua y un gran número de características del contenido del corpus. Los resultados que aparecen en el listado se pueden ordenar en orden creciente o decreciente según la lengua del texto, el nombre del corpus o el número de palabras que contiene.

El listado, además de mostrar la lengua, el nombre y el tamaño de cada corpus, tiene un botón a la derecha (marcado con el símbolo ...) que permite obtener ciertos datos y realizar una serie de acciones sobre el corpus en cuestión (salvo la última, las demás solo son aplicables a corpus privados creados tal y como se verá en la siguiente sección):

- a. *Eliminar corpus*: Borra el corpus y deja de estar disponible en *Sketch Engine*.
- b. *Editar metadatos de corpus*: Permite modificar los valores sobre la descripción y contenidos del corpus.
- c. *Aumentar*: Añade contenido a un cierto corpus.
- d. Gestionar corpus: Permite crear subcorpus.
- e. *Compartir*: Para enviar a alguien un enlace a ciertos corpus.
- f. Descargar: Descarga una copia del corpus en el disco duro local.
- g. Ver detalles del corpus: Da información detallada del contenido del corpus, así como de su etiquetado y de su composición. En la Figura 9 se muestran los detalles del corpus español por defecto esTenTen11.

Además, en esta pestaña tenemos la posibilidad de crear nuestro propio corpus.

 MIS CORPUS: Se listan los corpus que haya creado el usuario, además de permitir crear nuevos corpus (Figura 10). Mientras no hayamos creado algún corpus propio, esta página permanecerá vacía de contenido.

| SELECT CO | ORPUS Spanish Web 2011 (esTenTen11, Eu + Am) | | Get more space 🕣 🕝 🦻 🗾 |
|-----------|--|--------------------|---|
| BÁSICO AV | MI CORPUS | | |
| | Q Con Word Sketch (492 corpus) (94 idiomas) Cualquier idioma | | NEW CORPU |
| Language | Name | ↓ Words | COPPUS CATEGORY |
| English | Timestamped JSI web corpus 2014-2018 English | 35 911 543 638 🚥 | ALL CA 492 |
| English | English Web 2013 (enTenTen13) | 19.685.733.337 *** | RECIENTES () 1 |
| English | Timestamped JSI web corpus 2014-2016 English | 18.315.071.361 *** | MI CORPUS () 0 COMPARTIO CONMIGO () 2 |
| Spanish | Spanish Web 2018 (esTenTen18) | 17.553.075.259 *** | DESTACADOS () 17 |
| German | German Web 2013 (deTenTen13) | 16 526 335 416 *** | USO GENERAL () 296 WEB () 242 |
| English | English Web 2015 (enTenTen15) | 15.703.895.409 *** | NON-WEB () 250 |
| Russian | Russian Web 2011 (ruTenTen11) | 14.553.856.113 *** | SPOKEN () 50 |
| English | English Web 2012 (enTenTen12) | 11.191.860.036 *** | SPECIALIZED () 187 |
| Czech | Czech Web 2017 (csTenTen17) | 10 502 222 474 *** | MULTIMEDIA @ 1 |
| French | French Web 2012 (frTenTen12) | 9.889.689.889 *** | CORPUS DE APRENDICES ① 4 ERROR ANNOTATED ① 2 |
| Spanish | Spanish Web 2011 (esTenTen11, Eu + Am) | 9 497 213 009 *** | DESARROLLO DE GRAMÁTICA () 0 |
| Japanese | Japanese Web 2011 (iaTenTen11) | 8 432 256 578 *** | Part is the second |

Figura 8: página de selección de corpus de Sketch Engine (Avanzado)

| INFORMACIÓN GE | NERAL | CUENTA | | LEYENDA DE E | TIQUETAS | SUFIJOS | LEMPOS | | TAMAÑOS D | DE LÉXICO | |
|-----------------------------------|---------------------------------------|--|---|---|--|--------------|--|--|--|------------|--|
| Language | Spanish | Tokens | 10.985.547.573 | noun | N.* | noun | | -0 | word | 23.693.613 | |
| Descripción del corpus | Description | words | 9.497.213.009 | verb | V.* | verb | | -V | tag | 385 | |
| Etiquetas | Description | Oraciónes | 407 322 228 | adjective | A* | adjective | | -1 | lempos | 20.740.962 | |
| Gramática de word sketch | Description | Párrafos | 213.364.889 | adverb | R.* | adverb | | -1 | gender_lemma | 19.405.974 | |
| | | Documentos | 22.287.384 | pronoun | P.* | pronoun | | -p | tags | 712 | |
| | | | | conjunction | C.* | conjunction | | -0 | morphemes | 19.593.136 | |
| | | | | preposition | S.* | preposition | í. | -4 | lc | 20.401.771 | |
| | | | | determiner | D.* | numeral | | -m | lemma | 19.469.130 | |
| | | | | | | | | | | | |
| | | | | interjection | 1.* | | | | shorttag | 13 | |
| | | | | Interjection numeral More information | 1* Z* | | | | shorttag | 13 | |
| | ATRIBUT(| OS ESTAD | | Interjection numeral More information L SUBCORPUS | 1* Z* | | Tokens | word | shorttag | 13 | |
| ESTRUCTURAS Y | ATRIBUT(22.287.384 | Subcorpu American | ÍSTICAS DE 15 _\$panish_domain | Interjection numeral More information L SUBCORPUS sar,bo,el,co,cr,cu,do,ec, | L* Z.* | .py,pe,uy,ve | Tokens 8.614.262.666 | word 7.447 | shorttag s % 192.496 78,415 | 13 | |
| ESTRUCTURAS Y doc (6) p (1) | ATRIBUT(22.287.384 213.364.889 | OS ESTAD Subcorpu American European | ÍSTICAS DE 15 _\$panish_domain 1_\$panish_domain | Interjection numeral More information L SUBCORPUS s_ar,bo,cl.co,cr,cu,do,ec, es | L* Z.* | .py.pe.uy.ve | Tokens 8.614.262.666 2.345.025.922 | word 7.447 2.027 | shorttag 5 % 192.496 78.415 319.124 21.346 | 13 | |
| ESTRUCTURAS Y doe (6) p (1) | ATRIBUT(22.287.384 213.364.889 | OS ESTAD Subcorpu American European Wikipedia | ÍSTICAS DE 15 _\$panish_domain 1_\$panish_domain | Interjection numeral (*) More information L SUBCORPUS s_ar,bo,cl,co,cr,cu,do,ec, es | L* Z* | .py.pe.uy.ve | Tokens 8.614.262.665 2.345.025.922 210.971 | word 7.447 2.027 | shorttag 5 % 192.496 78.415 319.124 21.346 182.388 0.002 | 13 | |
| ESTRUCTURAS Y doc (6) p (1) | ATRIBUT(22.287.384 213.364.889 | OS ESTAD Subcorpu American European Wikipedia boe.es_i | ÍSTICAS DE spanish_domain spanish_domain official_newspape | Interjection numeral () More information L SUBCORPUS s_ar.bo.cl.co.cr.cu.do.ec. es r_of_the_Government_of | L* Z* sv.gt.hn,mx,ni,pa | .py.pe,uy,ve | Tokens 8.614.262.666 2.345.025.922 210.971 103.576.368 | word 7.447 2.027 89 | shorttag 5 % 192.496 78.415 319.124 21.346 182.388 0.002 543.723 0.943 | 13 | |
| doc (6) p (1) | ATRIBUT(22.287.384 213.364.889 | Subcorpu American European Wikipedia boe.es_r | İSTICAS DE ss spanish_domain sspanish_domain official_newspape r.es_supermarke | Interjection numeral More Information L SUBCORPUS s_ar.bo.cl.co.cr.cu.do.ec, _es r_of_the_Government_of t_magazine | L* Z* sv.gt.hn,mx,ni,pa t_Spain | .py.pe,uy.ve | Tokens 8.614.262.666 2.345.025.922 210.971 103.576.368 14.102.624 | word 7.447 2.027 89 12 | shorttag s % 102.0406 78.415 319.124 21.346 182.368 0.002 543.723 0.043 191.984 0.128 | 13 | |
| ESTRUCTURAS Y dec (6) p (1) | ATRIBUT(22.287.384 213.364.889 | Subcorpu American European Wikipedia boe.es consume elmundo. | ÍSTICAS DE spanish_domain spanish_domain miticial_newspape r.es_supermarke s_newspaper | Interjection numeral More information L SUBCORPUS s_ar,bo,cl.co,cr.cu.do,ec, _es r_of_the_Government_of (megazine | L* Z* sv.gt.hn,mx,ni,pa | .py.pe,uy.ve | Tokens 8.614.262.666 2.345.025.922 210.971 103.576.368 14.102.624 25.155.342 | word 7.447 2.027 89 12 21 | shorttag s % 192.406 78.415 319.124 21.346 182.386 0.002 182.386 0.024 191.494 0.128 747.267 0.229 | 13 | |

Figura 9: detalles del corpus español por defecto esTenTen11

| | ECT CORPUS | S Spanish Web 2011 (esTenTen1 | 11, Eu + Am) 🤇 🕧 | Get more space 🕣 | ය () | 10 : |
|------|------------------------------|-------------------------------|------------------|------------------|--------|-------------|
| BÁSI | CO AVANZADO | MI CORPUS | | | | |
| | | | | | NEW CO | RPUS |
| Lan | guage | | Name 🛧 | | Wor | ds |
| Ning | gún corpus cumple los criter | ios seleccionados. | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| • | | | | | | |
| Ξ. | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |

Cuando queramos cambiar el corpus sobre el que trabajamos, podremos acudir a la página de selección de corpus para utilizar sus filtros. No obstante, en todas las páginas de *Sketch Engine* podemos ver, en la **parte superior**, un campo de texto que muestra el **corpus actualmente seleccionado**, y que permite escoger entre los corpus disponibles filtrando por su nombre. Además, la página del panel de control da la posibilidad de cambiar rápidamente entre los últimos corpus seleccionados, como ya se comentó previamente.

@ 0 8 0

3. - Construcción de un corpus

Para poder utilizar *Sketch Engine*, necesitamos seleccionar uno de los corpus existentes. Además de los casi 500 corpus actualmente disponibles, *Sketch Engine* permite crearnos nuestros propios corpus utilizando distintas estrategias:

- A partir de archivos de texto locales.
 - De forma automática, descargando el texto desde Internet:
 - Utilizando una colección de palabras semilla para descargar el corpus desde Internet.
 - Descargando archivos desde un sitio web (indicando simplemente el sitio web, o especificando las páginas concretas).

3.1. - Crear un corpus a partir de archivos de texto locales

Lo primero que debemos hacer es clicar en **NEW CORPUS**. Este botón está accesible desde varios puntos:

- En el panel de control, un botón azul en la parte superior derecha (Figura 7).
- En la página de selección de corpus, un botón rojo en la parte superior derecha tanto en la pestaña AVANZADO (Figura 8) como en la de MIS CORPUS (Figura 10).
- A través de la página de gestión de corpus (accesible al clicar el botón *GESTIONAR CORPUS* en el panel de control).

Sketch Engine permite crear corpus monolingües pero también corpus paralelos en varias lenguas. En las dos siguientes subsecciones se explicará cómo crear ambos tipos de corpus.

3.1.1 - Corpus monolingües

Al clicar en *NEW CORPUS*, aparece la pantalla que se muestra en la Figura 11, donde se debe asignar un nombre al corpus que se va a crear, especificar la lengua del texto que se incluirá y añadir una descripción sobre el corpus.

| 1. CREATE CORPUS > 2. ADD TEXTS > 0. | 3. COMPILE | | | | |
|--------------------------------------|-------------------------------|--|---------|-----------------|-----|
| | Build your own private corpus | from texts on the web or from your own documents. | | | |
| | Name | El ingeniero hidalgo | | | |
| | Language | Spanish Q | | | |
| | | MULTILINGUAL | | | |
| | Description | Las dos partes de El Quijote: El ingenioso hidalgo Don Quijote de la Mancha, y El | | | |
| | | ingenioso caballero Don Quijote de la Mancha | | | |
| | Storag | e used: 0 de 1,000,000 words (0%) | | | |
| | | | | | |
| | | Available features - | | | |
| | | ATRÁS | | | |
| | | | Back to | the original in | her |



Si clicamos en **Available features** se nos muestran las características que tendrá el corpus que vamos a crear. Estas características las produce directamente *Sketch Engine* cuando procese el corpus, y son las que se muestran en la Figura 12.



Figura 12: características que tiene un corpus privado creado con Sketch Engine

| 1. CREATE CORPUS > 2. ADD TEXT | S > 3. COMPILE | | |
|---|---|---|--|
| | Encuéntrame textos en el internet La descargamés textos de páginas web | Tengo mis propios textos Sube aus propios archivos (.bd., .pd) | |
| CORPUS CONTENT | | | |
| | Nada a |] quí | |
| | Corpus is empty Add texts o | using the options above. | |
| Storage used: 0 de 1,000,000 words (0%) | ATRÁS | KINENTE | |

Tras ello, se debe clicar en *SIGUIENTE*, con lo que se llega a la pantalla mostrada en la Figura 13. En este punto, caben dos opciones para elegir el texto que se debe incluir en el corpus:

- **Encuéntrame textos en el internet**: Esta opción debe utilizarse para crear un corpus automático compuesto por texto descargado de Internet. Su uso se explicará en la siguiente sección.
- Tengo mis propios textos: Se debe escoger esta opción para poder utilizar texto almacenado localmente.

Tras clicar en la segunda opción, se nos pedirá que subamos los archivos que contengan el texto (hasta un máximo de 100). Tras subirlos, *Sketch Engine* procesará el texto, dividiéndolo en *tokens*⁶, una tarea que hace a ritmo de unos 5.000.000 de palabras al minuto. Finalmente, la parte inferior se actualizará mostrando la localización y la cantidad de *tokens* que tiene el corpus hasta el momento (Figura 14).

| older | Words | |
|--------|---------|--|
| upload | 181 596 | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |

Figura 14: estado actual del corpus que se está creando

Sketch Engine puede manejar muchos tipos distintos de archivos: .doc, .docx, .htm, .html, .ods, .pdf, .tei, .tmx, .txt, .vert, .xlf, .xliff, .xls, .xlsx y .xml. Además, es posible también seleccionar **archivos comprimidos** de tipo .zip, .tar.bz2, .tar.gz o .tgz que contengan más de un archivo en su interior (todos aquellos archivos comprimidos que no tengan alguna de las extensiones descritas anteriormente, se ignorarán).

En vez de subir un archivo, también es posible copiar y pegar un texto directamente. Para ello basta con clicar en el enlace *or paste text* que aparece en la parte superior derecha de la pantalla. Al clicar se nos abre un cuadro de texto donde podremos pegar el texto deseado.

| CREATE CORF | PUS El ingeniero hidalgo 🔍 🔺 | Get more space 🕢 | C | 0 | | |
|------------------|--|------------------|------|-----------|-------------|----|
| CORPUS: El in | geniero hidalgo (Spanish) | | | | | |
| 1. CREATE CORPUS | 2. ADD TEXTS >> 3. COMPILE | | | | | |
| | | | | | | |
| | | | | | | |
| | Listo | | | | | |
| | | | | | | |
| | El corpus está lisito para comptación. | | | | | |
| | | | | | | |
| | ADD MORE TEXTS COMPILE | | | | | |
| | Evant entitione - 100 - | | | | | |
| | Expertsenings + Log + | | | | | |
| | | | | | | |
| | | | Back | to the or | riginal int | er |

Cuando se haya subido todo el texto, se debe clicar en el botón rojo **SIGUIENTE** situado en la parte inferior de la pantalla lo que nos lleva a la pantalla mostrada en la Figura 15. En caso de que nos

⁶ Token es un término utilizado en Lingüística Computacional y Procesamiento del Lenguaje Natural para referirse a un elemento de texto (habitualmente separadas por espacios) que ha de procesarse de manera independiente. Un token puede ser una palabra, pero también lo son los signos de puntuación, números, símbolos, etc.



| EXPERT SETTINGS | |
|---|---|
| Duplicated content ⑦ | Remove duplicated content |
| Structures and attributes to keep ⑦ | todos [ninguno |
| Sketch grammar ⑦ | ✓ g (27.115) ✓ s (5.087) ● Spanish (Freeling tagset) 2.0 recommended () ○ Universal (no tags) () |
| Term grammar (?) | O None (no word sketches) Snanish (Freeling) for term extraction 1.4. |
| | O None (no term extraction) |
| Structure name for documents ⑦ Structure for document ⑦ | doc Same as name for documents |
| | |
| CANCEL | AR SAVE AND COMPILE |
| Figura 16: ajustes previos | a la compilación de un corpus privado |

Tras haber añadido todo el texto al corpus, este se debe compilar clicando en el botón **COMPILE**. El proceso de compilación es el que aporta la información extra al corpus recién creado, principalmente gracias al etiquetado morfosintáctico, al lematizado y al análisis sintáctico, gracias a los cuales se realiza también un análisis acerca de patrones de aparición de lemas y su posición y función en la oración. Es necesario realizar la compilación antes de poder utilizar este corpus, aunque no es necesario hacerlo inmediatamente después de haber creado el corpus: su compilación puede quedar como pendiente (en caso de que aún no esté creado completamente) y realizarse más adelante.

El proceso de compilado puede llevar segundos o minutos, según sea su tamaño. Pasado este tiempo, el corpus ya estará **disponible** para su uso, con lo que aparecerá como un corpus más, tanto en el listado de corpus propios del usuario, *MIS CORPUS*, o en el listado de todos los corpus, *AVANZADO* (mostrados respectivamente en la Figura 10 y la Figura 8).

3.1.2 - Corpus multilingües

Para crear un corpus multilingüe, deberemos clicar en el botón *MULTILINGUAL* que puede verse en la Figura 11, tras lo cual se muestra el menú que aparece en la Figura 18. Lo único que hace falta es ponerle un nombre al corpus y adjuntar un archivo con el texto paralelo. Dicho archivo puede ser de formato .txt .tmx, .xls, .xlsx o .xliff y debe contener en cada línea el texto alineado (pueden ser frases o párrafos) en las dos o más lenguas que contenga. La primera línea debe contener el nombre en inglés de la lengua en cuestión. Tras escoger el archivo se debe pulsar en el botón *SIGUIENTE* y se nos

080

mostrará el menú que aparece en la Figura 18. Al crear un corpus multilingüe se crearán tantos corpus monolingües como lenguas haya en el archivo que hayamos utilizado, con lo que cada uno tendrá su propio nombre y tendrá especificada su propia lengua.

| 1. UPLOAD CORPORA > 2. CONFIG | GURACIÓN > 3. COMPILE | | | | | |
|-------------------------------|---|--|------|------------|-----------|---|
| | Build your own private corpus from texts on the web or from your own docume | ents. | | | | |
| | Name Corpus multilingüe | | | | | |
| | Corpus file | | | | | |
| | ARCHVO your caur uppead_text_tree_abc_view_ Find aut more about file form | and ats | | | | |
| | SINGLE LANGUAGE | | | | | |
| | Storage used: 181,596 de 1.000,000 words (18%) | | | | | |
| | ATRÂS BIGUIENTE | | | | | |
| | | | | r to the c | | |
| Figura 1 | 7: ajustes previos a la compilación de un corp | pus multilingüe | Bace | (to the c | Algelat B | |
| Figura 1 | 7: ajustes previos a la compilación de un corp | pus multilingüe Get more space () | G | () () | | - |
| Figura 1. CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web C BURACIÓN >> 3. COMPILE | pus multilingüe Get more space () | 6 | 0 | P | - |
| Figura 1 CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web Corpus de monolingue Corpus de de corpus names and/or the automatically detected language | pus multilingüe Get more space () al alow 95 | (c) | 0 | | |
| Figura 12 CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web GURACIÓN > 3. COMPILE Each language in the source file will be processed into a separate monolingue corpus and aligned with the corresponding corpus in the other language(s). Be you can change the corpus names and/or the automatically detected language Corpus name (Spanish) Corpus multilingúe, Spanish | pus multilingüe Get more space (*) al elow es | 60 | 0 | | - |
| Figura 12 CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web GURACIÓN > 3. COMPILE Each language in the source file will be processed into a separate monolingue corpus and aligned with the corresponding corpus in the other language(s). Be you can change the corpus names and/or the automatically detected language Corpus name (Spanish) Corpus multilingúe, Spanish Corpus language (Spanish) Spanish | pus multilingüe Get more space () al lelow 195 | 63 | 0 | | |
| Figura 12 CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web GURACIÓN > 3. COMPILE Each language in the source file will be processed into a separate monolingue corpus and algined with the corresponding corpus in the other language(s) Be you can change the corpus names and/or the automatically detected language Corpus name (Spanish) Corpus multilingüe, Spanish Corpus language (Spanish) Spanish Corpus name (English) Corpus multilingüe, English | pus multilingüe Get more space (*) al lelow los | G | 0 | | - |
| Figura 1. CONFIGURACIÓN | 7: ajustes previos a la compilación de un corp Corpus de web C GURACIÓN > 3. COMPILE Each language in the source file will be processed into a separate monolingue corpus and aligned with the corresponding corpus in the other language(s). B you can change the corpus names and/or the automatically delocted language Corpus name (Spanish) <u>Corpus multilingüe, Spanish</u> Corpus language (Spanish) <u>Spanish</u> Corpus mame (English) <u>Corpus multilingüe, English</u> | pus multilingüe Get more space (*) al lelow leis | G | 0 | | - |
| Figura 12 CONFIGURACIÓN | Corpus de web C Corpus de web C GURACIÓN > 3.COMPILE Each languaga in the source file will be processed into a separate monolingue corpus and aligned with the corresponding corpus in the other language(s) is you can change the corpus names and/or the automatically detected language Corpus name (Spanish) Corpus multilingüe, Spanish Corpus name (Spanish) Spanish Corpus name (English) Corpus multilingüe, English Corpus language (English) English | pus multilingüe Get more space (*) al lelow los | 60 | () | | |

El último paso consiste en pulsar el botón *SIGUIENTE* para que *Sketch Engine* compile los distintos corpus monolingües que componen el corpus multilingüe, y estén disponibles en la lista de corpus.

3.2. - Crear un corpus automático a partir de texto descargado de internet

No es necesario tener localmente el contenido del corpus que se vaya a crear, sino que dicho contenido se puede extraer directamente de Internet. Si queremos añadir este tipo de contenido a un corpus, deberemos proceder tal y como se explicó en la sección anterior, hasta llegar al punto mostrado en la Figura 13, momento en el cual deberemos clicar en *Encuéntrame textos en el internet*. Es importante notar que no hay ningún problema en crear un corpus que contenga texto extraído de internet junto con texto subido de archivos locales o el portapapeles.

Tras clicar en dicha primera opción se nos muestra el menú de la Figura 19, donde se debe elegir una de las tres maneras en las que se añadirá texto al corpus automáticamente:

- 1) Utilizando palabras semilla.
- 2) Descargando el texto de una lista de direcciones web (URLs) que indican páginas concretas.



3) Eligiendo **una web concreta** y dejando que *Sketch Engine* sea quien encuentre y descargue las páginas que se deben utilizar.

| Input type | Web search ⁽²⁾ URLs ⁽²⁾ Website ⁽²⁾ e.g. "shopping centre" multiplex hypermarket | 0 |
|---------------|--|---|
| Folder name ⑦ | Hit ENTER after each word or phrase. web1 Web search settings • Black list settings • White list settings • Size restrictions • | - |
| | Compile when finished CANCELAR | |

Estas tres formas se explican en las siguientes subsecciones.

3.2.1 - Uso de palabras semilla

Si se escoge el primero de los tres métodos en la Figura 19 (*Web search*), se deberán incluir en el campo de texto las palabras semilla, separadas por comas, que determinarán los archivos que se descargarán de forma automática. Se aconseja utilizar entre 3 y 20: cuantas más palabras semilla se definan, mayor será el número de páginas web descargadas y más amplio será el corpus (será mayor en tamaño y contendrá texto más diverso), ya que se utilizarán **subgrupos de tres** de dichas palabras semilla y con ellas se buscarán contenidos en internet usando el motor de búsqueda Bing.

Hay una serie de ajustes en relación a la manera en la que se extraerá el corpus:

- Web search settings: Se puede escoger cuántas páginas se descargan en cada búsqueda realizada, además de la posibilidad de incluir una lista de sitios web que serán los únicos de los que se descargarán páginas.
- Black list settings: Evita descargar páginas que contengan más palabras semilla que un cierto límite, y también aquellas que tengan más de un cierto número de palabras semilla distintas. Además, se puede introducir una lista de términos de forma que no se descargará ninguna página que contenga alguno de dichos términos.
- White list settings: Al contrario que en el grupo anterior, aquí se especifica el número mínimo de palabras semilla (distintas o no) que debe tener una página para aceptarse, y puede incluirse una lista de términos que tienen que estar incluidos en la página para ser aceptada. También se puede indicar cuál debe ser el ratio mínimo entre el número de palabras totales en el documento y el número de palabras clave.

- **Size restrictions**: Se pueden imponer límites de tamaño superior e inferior de los documentos descargados antes o después de limpiarlos de todo contenido no estrictamente textual.

Además, hay una opción, activada por defecto, *Compile when finished* que, como indica, procederá a la **compilación automática del corpus** tras la descarga y preprocesado del texto.

Cuando estén fijadas las palabras clave y todos los ajustes, se debe clicar el botón *iYA!* y esperar a que *Sketch Engine* descargue y procese las páginas, y cree finalmente un corpus con todo ese texto. Según la cantidad de palabras semilla que se haya utilizado, este proceso puede durar un tiempo considerable (desde algunos minutos, hasta horas) durante el cual se va mostrando el progreso. Al acabar, el corpus queda disponible para su uso.

No es necesario esperar a que el corpus se haya descargado para poder hacer otras cosas a la vez: la descarga se realiza desde el servidor de *Sketch Engine* y no desde el propio ordenador, con lo que se puede cerrar la ventana o incluso hacer otras tareas en *Sketch Engine* mientras el corpus crece.

3.2.2 - Uso de URLs concretas

Este método es el **menos automático** de los utilizados para crear un corpus descargando texto de Internet: crea el corpus a partir de una serie de documentos cuyas direcciones web (URLs) se indican de forma explícita. Basta utilizar el segundo de los tres métodos que aparecen en la Figura 19 (*URLs*), con lo que se mostrará el menú que aparece en la Figura 20.

| ← TEXTS FROM WEB | |
|------------------|---|
| Input type | Web search ⁽²⁾ URLs ⁽²⁾ Website ⁽²⁾ https://es.wikipedia.org/wiki/Lingüística https://es.wikipedia.org |
| Folder name ② | https://es.wikipedia.org/wiki/Lingüística_computacional |
| | Black list settings 👻 |
| | White list settings |
| | \checkmark Compile when finished $^{\odot}$ |
| | CANCELAR (YA! |

Las direcciones web deben escribirse en el campo de texto, separadas por espacios. Y se tienen unas opciones similares a las explicadas en la subsección anterior. Tras rellenar todos los datos, basta con clicar el botón *iYA!* y esperar a que se descarguen y procesen las páginas webs indicadas para que el corpus aparezca en el listado y pueda comenzar a utilizarse.



3.2.3 - Uso de un Website

Otra posibilidad a la hora de crear un corpus automático, es la de **elegir un** *Website* y dejar que *Sketch Engine* **encuentre páginas alojadas en dicha web**, y las utilice para crear un corpus. Esto es interesante, por ejemplo, para crear un corpus con artículos de un periódico online.

Para ello, se escoge el tercero de los tres métodos en la Figura 19 (*Website*). Al hacerlo, se muestra un campo de texto (como se ve en la Figura 21), en el que se debe introducir la **dirección de la web** (por ejemplo, <u>www.elpais.com</u>). También es posible indicar escoger un *path* concreto dentro de una página web, como <u>www.elpais.com/internacional/</u>, de forma que *Sketch Engine* solo buscaría en este caso aquellas páginas de la web <u>www.elpais.com</u> que estén **bajo el subdirectorio** /internacional/ (es decir, solo páginas de noticias internacionales). Por lo demás, las opciones que se dan son las mismas que en la subsección anterior.

| ← TEXTS FROM WEB | | |
|------------------|--|---|
| Input type | Web search ⁽²⁾ URLs ⁽²⁾ Website ⁽²⁾ http://www.elpais.es | 0 |
| Folder name ⑦ | web1 Black list settings White list settings Size restrictions | |
| | Compile when finished ⁽²⁾ | |

Tras especificar la página web de la que se extraerá el texto y seleccionar los ajustes, se clica en *¡YA!*, y como en los casos anteriores, se debe esperar hasta que haya terminado el proceso. *Sketch Engine* procurará buscar en la web indicada y extraer de ahí hasta **2.000 páginas**, con lo que el proceso puede alargarse durante **horas**.

Es importante tener en cuenta que usualmente las páginas web no desean gastar tráfico de datos sirviendo páginas a procesos robot como el que utiliza *Sketch Engine*, que se dedica a descargar de forma automática el contenido de una web. Es por ello que el uso de este método **no funcione correctamente en muchas páginas**.

3.3. - División en subcorpus

Muchas veces es interesante obtener datos de una cierta parte de un corpus y no de su totalidad. Para hacer esto se puede dividir un cierto corpus en varios subcorpus, de forma que tengamos acceso únicamente a ciertos contenidos a la hora de analizarlo. Para ello, teniendo seleccionado el corpus que queremos subdividir, clicamos en el botón de *GESTIONAR CORPUS* que aparece en la parte superior

central de la página del panel de control (Figura 7). Tras hacer esto, se nos muestra la página de gestión de corpus que se ve en la Figura 22, donde deberemos clicar en *Subcorpora*.

| Las dos partes de El Quijote: El ingenioso | hidalgo Don Quijote de La Mancha, | y El ingenioso caballero Don Quijote de la | Mancha | |
|--|-----------------------------------|--|--------------------------|--|
| Browse View documents and folders | Add texts to corpus | Share Share corpus with other users | Download Download | Compile compile corpus or change compil settings |
| Delete Remove corpus permanently | Subcorpora manage subcorpora | Configure Change corpus configuration | Logs View corpus logs | Hew corpus Crear corpus nuevo |
| | 7 | BACK TO DASHBOARD | | |
| | | | | |
| | | | | |

Si el corpus escogido no había sido subdividido previamente, no contendrá ningún subcorpus, y se nos mostrará una pantalla como la que aparece en la Figura 23. Para crear un subcorpus debemos clicar en el botón *CREATE SUBCORPUS*, lo que hace aparecer el menú mostrado en la Figura 24, donde se ha clicado en *Expandir todo* para mostrar el contenido el corpus. Basta con escribir el nombre del subcorpus y seleccionar los documentos que queremos que contenga (clicando sus nombres o IDs).

| SUBCORPORA El ingen | iero hidalgo 🔍 🕜 | Get more space ④ | Θ | 0 | |
|----------------------|-------------------|------------------|-------|-----------|------------|
| CORPUS: El ingeniero | hidalgo (Spanish) | | | | |
| | | | | | |
| | | | CREAT | E SUBCO | RPUS |
| | Nada aqui | | | | |
| | | | | | |
| | ATRÁS | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | Back | to the on | oinal inte |

| | | | Expandir todo contraer to |
|--------------|----------|-------------------------------------|---------------------------|
| File ID | 55 | File name 55 El Ingenioso Hida × | |
| file10046157 | <u> </u> | tipo de húsqueda Q | |
| file10049096 | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

Se pueden crear tantos subcorpus como se desee, que pueden tener contenido solapado entre ellos. Cuando lo hayamos hecho, sus nombres y tamaños aparecen en la página de subcorpus (Figura 25) que antes se mostraba vacía.

| Name | | Tokens | Words | % | | CREATE | SUBCOR | PUS |
|---------|-------|---------|---------|----|----|--------|--------|-----|
| Parte 1 | | 208.711 | 180.458 | 48 | ŧ. | | | |
| Parte 2 | | 225.856 | 195.282 | 52 | 8 | | | |
| | ATRÁS | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |

El tener un corpus dividido en varios subcorpus nos capacita para que *Sketch Engine* realice una serie de funciones únicamente sobre estas partes del corpus y no sobre el total.

@ 0 8 0

4. - Explorar el corpus

La página de selección de corpus de *Sketch Engine* (mostrada en la Figura 8) muestra los corpus de los que se dispone. La forma de trabajar en *Sketch Engine* es la de **escoger el corpus** sobre el que se va a trabajar, y una vez elegido, **todas las operaciones se hacen sobre ese corpus** (hasta que se escoja un corpus distinto). En todo momento se nos muestra el corpus escogido en la parte superior de todas las páginas de *Sketch Engine*.

Sobre este corpus seleccionado, *Sketch Engine* nos muestra en la pantalla de panel de control (Figura 7) una serie de herramientas disponibles. De ellas, hay ocho que están disponibles para todos los corpus monolingües que hayan sido compilados con *Sketch Engine* y que son las siguientes:

- 1) Word Sketch: Busca patrones sintácticos estadísticamente relevantes de un lema.
- 2) Diferencia Sketch: Compara los patrones sintácticos de dos lemas.
- Tesauro: Construye el tesauro o vocabulario de palabras relacionadas semánticamente con un cierto lema dado.
- Concordancia: Busca concordancias –las palabras deseadas, junto con todas las citas de los lugares en las que se hallan– simples o avanzadas.
- 5) *Lista de palabras*: Obtiene la lista de palabras con sus **frecuencias** de aparición en el corpus.
- 6) *N-gramas*: Extrae las expresiones multipalabra más típicas del corpus.
- 7) *Palabras clave*: Obtiene la lista de **palabras clave** (y términos multipalabra) del corpus.
- 8) **OneClick Dictionary**: Crea un esbozo de diccionario a partir de un corpus.

Además, hay otras dos herramientas que necesitan tener seleccionado un corpus con unas ciertas características especiales:

- 9) **Concordancia paralela**: El corpus seleccionado debe ser un corpus paralelo en dos o más lenguas. Establece correspondencias entre palabras de dos lenguas.
- 10) **Tendencias**: Se necesita un corpus cuyos documentos estén etiquetados con la fecha de publicación. Permite obtener datos sobre la variación en el uso de ciertas palabras en la lengua del corpus.

Estas herramientas se explicarán a lo largo de las siguientes secciones.

4.1. - Word Sketch: Relaciones gramaticales

Esta opción, permite explorar los patrones sintácticos de un lema concreto proporcionando **información sobre la posición y función sintáctica** en que aparece el lema en cuestión. Al clicar el botón *Word Sketch* se muestra un menú con dos pestañas en la parte superior: BÁSICO y AVANZADO. Nos centraremos en esta segunda pestaña, que es la que se muestra en la Figura 26. Ahí se debe rellenar el campo *lema*, que puede rellenarse con uno o más lemas. También se tienen otras opciones para escoger, como la **categoría gramatical** que se busca (puede haber lemas coincidentes para más de una categoría gramatical, como "ante" o "bajo", que pueden ser preposiciones o un nombre/adjetivo), la frecuencia mínima de aparición o la puntuación mínima. La puntuación es el resultado de una operación matemática en la cual se mide cuántas veces aparece el lema en el corpus



teniendo una cierta relación gramatical con otra palabra, en relación a cuántas veces aparecen por separado. También se puede restringir la búsqueda a un cierto **subcorpus** del corpus seleccionado.

| BASICO | AVANZADO | |
|------------------|------------|-----------------------------|
| Búsqueda | | Subcorpus |
| lema | | ningún (todo el corpus) 🔹 🕇 |
| | | |
| Categoría (| gramatical | |
| auto | | |
| adjective | | |
| adverb | | |
| noun | | |
| verb | | |
| prepositio | n | |
| Frecuencia mínin | na | Puntación mínima |
| auto | | 0 |
| | | |
| | | ITA! |
| | | |

El resultado del *Word Sketch* para "caballero" (en un corpus que está formado por el texto de las dos partes de El Quijote), aparece en la Figura 27. En la imagen se puede observar como "andante" es el **modificador más usual** de "caballero". Se puede observar también muchas otras columnas⁷, en las que se muestran los verbos que tienen a "caballero" **más frecuentemente como objeto** (en este caso, "armar"), ciertas otras relaciones sintácticas, y las **combinaciones más usuales con preposiciones** y verbos copulativos, por ejemplo. Cada columna representa una relación gramatical.

El menú desplegable de la parte superior izquierda nos indica que los resultados que se muestran son los que se obtienen cuando el lema actúa como una cierta categoría gramatical, y cuántas apariciones tiene. Si clicamos en él podemos observar que en el corpus aparece "caballero" como adjetivo 162 veces. No obstante, ya que el etiquetado es automático, podemos comprobar cómo muchos de dichos resultados muestran al lema "caballero" actuando como nombre y no como adjetivo.

A la derecha de dicho menú desplegable se tiene un botón marcado con el icono ... que nos permite pasar directamente a las herramientas de Corcondancia y de Tesauro, que se explican más adelante en sus respectivas secciones. En la parte superior derecha se tienen seis iconos, cuya utilidad es la siguiente (de izquierda a derecha):

⁷ Debido a las limitaciones en cuanto al número mínimo de ocurrencias, que se fija en el menú de la Figura 26, aparecen relativamente pocas columnas. En corpus mayores se obtienen resultados mucho más detallados.



 Cambiar criterios: Permite modificar los criterios de búsqueda, haciendo aparecer el menú que se vio en la Figura 26.



- **Descargar**: Descarga los resultados que se están visualizando como un archivo de extensión .csv, .xls, .xlm o .pdf.
- Cambiar ajustes de visualización: Permite mostrar o no, para cada resultado que se muestra en las tablas, la frecuencia de aparición (en realidad, el número de apariciones), la puntuación y el ejemplo más común en el corpus. Además, permite ordenar los resultados de cada tabla según la puntuación o según la frecuencia, así como agrupar resultados semánticamente similares. Al habilitar esta opción se nos mostrará una barra de desplazamiento en la que podremos escoger cómo de parecidos deben ser los lemas para agruparse.
- Mostrar visualización: Al clicar en este icono se nos muestra una interesante visualización gráfica de los resultados de las tablas. El gráfico es interactivo y tiene además varios ajustes de configuración a la derecha que permiten añadir o eliminar características del gráfico. Para el resultado obtenido de "caballero", dicha representación gráfica se muestra en la Figura 28.
- Detalles de pantalla: Muestra los criterios de búsqueda que se están utilizando.
- Agregar a favoritos: Si clicamos en este icono, se guardará esta búsqueda en nuestros favoritos, y se hará accesible desde la pestaña *RESULTADOS FAVORITOS* que aparece en la parte inferior del panel de control (Figura 7).

Además, cada una de las tablas que se muestran como el resultado de aplicación de la herramienta *Word Sketch* tiene en su parte superior cuatro iconos:

- El que está más a la izquierda, permite cambiar a nuestro gusto la **posición** en la que aparecen las distintas tablas.
- El segundo icono permite ver las **colocaciones** de los resultados mostrados en dicha tabla.
- El tercer icono elimina todas las demás tablas de la visualización y deja visible solo esa tabla.
- El que está más a la derecha elimina dicha tabla de la visualización.



Además, para cada resultado mostrado en cada una de las tablas, tenemos un icono con el símbolo ... que permite acceder directamente a las herramientas de Concordancia y *Word Sketch* del lema buscado ("caballero" en este caso) junto con el lema que aparece en dicho resultado, además de poder acceder a la herramienta de Tesauro del lema que aparece en el resultado.

4.2. - Diferencia Sketch: Diferencias de uso entre palabras

Esta herramienta permite introducir **dos lemas** de búsqueda para **comparar sus patrones sintácticos y colocaciones léxicas** según ocurren en el corpus. También es posible realizar esta misma tarea usando **dos formas de un mismo lema** o un **mismo lema en dos subcorpus** de un mismo corpus.

| DASICO AVA | | | |
|--|---------------------|--|---|
| lemmasformas de palabSubcorpora | pras | | Ð |
| Lema 1 | | Lema 2 | |
| quijote | | sancho | |
| Categoría gramatic auto adjective adverb noun verb preposition | al | | |
| Frecuencia mínima auto | | Máximos items en bloques comúnes 12 🔄 | |
| Fiaur | ra 29: menú de la l | iYAI | h |

Al clicar en esta opción en el panel de control y escoger la pestaña AVANZADO, se nos abre el menú mostrado en la Figura 29. Esta es la forma por defecto, y busca las diferencias entre dos lemas distintos, que deben compartir categoría gramatical. Si escogemos la segunda opción en la selección de arriba, podremos introducir un lema y dos formas distintas para este lema, de forma que busque las diferencias entre dichas dos formas correspondientes al mismo lema. Con la tercera opción se nos permite introducir un lema y dos subcorpus del corpus seleccionado.

Por ejemplo, para saber qué patrones son comunes a los lemas "Quijote" y "Sancho" y qué patrones son más propios de un lema o del otro, se introducen los dos en el formulario de Diferencia Sketch. Se obtiene el resultado de la Figura 30.

IMPORTANTE: los **nombres propios** se lematizan en *Sketch Engine* **en minúscula**, con lo que si se escriben con mayúscula inicial, no se encontrarán resultados.

| quijo | ote 2.1 | 59× | 6.0 | 4.0 2 | 0 0 -2 | 0 -4.0 -6. | sanc | h0 2. | 167× | | | | হ | ŧ | 0 | i |
|----------------|---------|-----------|----------|----------|------------|----------------|---------------|----------|-----------|--------------|---------|--------|--------------|--------|----------|-------------|
| 0 | de "qui | jote/sanc | :ho" 🧿 | | | | | | | | | | | - | - | |
| ₽ | | | | | Ω× | ₽ | | | | Ŋ X | , ₽ | | | | | <u>(</u>) |
| | modi | fiers of | "quijo | te/sanch | o'' | | de "quij | ote/sa | ancho" | | | a " | quijot | e/sanc | :ho" | |
| valeros | D | 17 | 0 | 11.4 | | locura | 13 | 0 | 11.0 | | suceder | 13 | 3 (|) · | 10.4 | |
| famoso | | 14 | 0 | 10.9 | | humor | 6 | 0 | 10.0 | | dejar | ţ | 5 (|) | 9.0 | - • |
| ingenio | 50 | 5 | 0 | 10.3 | | historia | 7 | 0 | 10.0 | | avenir | 4 | 4 (|) | 8.9 | - • |
| verdade | oro | 4 | 0 | 9.4 | | aposento | 6 | 0 | 10.0 | | tratar | 4 | 4 (|) | 8.9 | - • |
| caballer | 0 | 4 | 0 | 9.1 | | hazaña | 4 | 0 | 9.5 | | ver | 11 | I 5 | 5 | 9.8 | 9.0 • |
| mismo | | 6 | 0 | 8.5 | | razón | 6 | 6 | 9.8 | 10.2 ••• | decir | 33 | 3 32 | 2 ' | 10.9 | 11.1 • |
| grande | | 11 | 9 | 9.0 | 8.7 ••• | carta | 3 | 5 | 8.9 | 10.1 ••• | duque | 4 | 4 4 | 1 | 8.8 | 9.3 • |
| bueno | | 4 | 31 | 7.1 | 10.1 ••• | azote | 0 | 3 | — | 9.6 ••• | volver | 3 | 3 11 | 1 | 7.9 | 10.1 • |
| pobre | | 0 | 3 | _ | 9.3 ••• | corazón | 0 | 4 | _ | 9.8 ••• | llamar | (| | 0 | _ | 9.4 • |
| present | e | 0 | 3 | _ | 9.7 ••• | jumento | 0 | 4 | _ | 10.0 ••• | quijote | (| | 5 | _ | 10.1 • |
| escude | 0 | 0 | 16 | _ | 10.3 ••• | gobierno | 0 | 0 | _ | 10.4 ••• | don | (|) 1() 1(| 9 2 | _ | 10.2 • |
| amyo | | 0 | 10 | _ | 12.0 ••• | mujer | 0 | , | _ | 10.5 ••• | manuai | (| | , | _ | 10.0 • |
| | | | ~ | | | | | ~ | | | | | ~ | | | |
| ⇔ | | | | | <u>Ø</u> × | ¢ | | | | (<u>)</u> × | ج | | | | | 0 |
| ve | rbs wit | h "quij | ote/sar | cho" as | object | verbs | with "quijote | e/sand | cho" as s | ubject | | "quij | ote/sa | ncho" | a | |
| ser | | 42 | 16 | 8.7 | 7.3 ••• | rogar | 3 | 0 | 9.3 | | panza | 3 | 0 | 10 |).1 | - •• |
| quedar | | 16 | 7 | 8.8 | 7.8 ••• | mirar | 3 | 0 | 9.1 | | sancho | 8 | 0 | 10 |).1 | - •• |
| pregunt | ar | 27 | 14 | 9.6 | 8.8 ••• | volver | 7 | 3 | 10.2 | 9.0 ••• | don | 3 | 0 | 8 | 3.0 | - • |
| ver | | 19 | 13 | 8.8 | 8.4 ••• | oír | 4 | 3 | 9.5 | 9.1 ••• | pie | 0 | 3 | | — | 9.5 •• |
| oír | | 10 | 7 | 8.1 | 7.8 ••• | decir | 28 | 22 | 10.8 | 10.4 ••• | duque | 0 | 3 | | — | 10.0 •• |
| replicar | | 71 | 58 | 10.8 | 10.7 ••• | ver | 10 | 10 | 10.4 | 10.4 ••• | amo | 0 | 6 | | - | 10.9 •• |
| decir | | 289 | 255 | 11.8 | 11.7 ••• | responder | 9 | 10 | 10.4 | 10.6 ••• | | | | | | |
| estar | | 15 | 13 | 8.5 | 8.5 ••• | pasar | 3 | 4 | 9.0 | 9.4 ••• | | | | | | |
| llegar | lor | 257 | 207 | 11.0 | 8.0 ••• | hacer | 1 | 10 | 9.5 | 10.0 ••• | | | | | | |
| volvor | lei | 231 | 307 | 7.4 | 82 | der | 3 | 6 | 8.4 | 9.5 | | | | | | |
| prosequ | ir | 4 | 7 | 6.9 | 7.9 ••• | llegar | 0 | 4 | | 9.4 ••• | | | | | | |
| | | | ~ | | | | - | ~ | | | | | | | | |
| _ | | | | | (X) X | | | | | <u>(</u>) × | - | | | | | Ö |
| ←* | " | iioto/or | nebo" | and/or | | € ⁷ | | ore of | "quijoto | ancho" | €7 | "auiia | tologr | ebo" | . | |
| sancho | qu | 33 | 0 | 12.7 | — ••• | nuestro | 10 | 0 | 9.5 | - ••• | sierra | 3 | 0 | 11. | 2 | |
| escude | 0 | 3 | 0 | 10.4 | | mi | 34 | 7 | 9.0 | 6.7 ••• | suelo | 0 | 3 | - | - | 10.3 •• |
| don | | 4 | 36 | 9.5 | 12.3 ••• | vuestro | 3 | 3 | 6.4 | 6.5 ••• | - | | | | | |
| primo | | 0 | 3 | - | 9.7 ••• | su | 14 | 17 | 6.7 | 7.0 ••• | | | | | | |
| quijote | | 0 | 33 | — | 12.7 ••• | | | | | | | | | | | |
| ¢ | | | | | <u>Ø</u> × | ¢ | | | | (C) × | ¢ | | | | | 0 |
| | "qu | ijote/s | ancho" | sobre | | | "quijote/s | ancho | " de | | | "quijo | te/san | icho" | de | |
| asno | | 0 | 3 ∡ | _ | 12.3 ••• | mancha | 112 | 0 | 13.1 | | mancha | | 112 | 0 | 13.1 | - • |
| , and a | | • | | | n X | | | | | ΩX | | | | | | Ő |
| ¢ , | | | | | 0,11 | ←' | Harris 6. 6 | | | | €* | | | 4-1 | | |
| VOT | | por "q | uijote/s | ancho" | | razón | quijote/sa | incho' | con | 11.5 | dar | con | quijo | te/sar | 1cho" | |
| VCI | 8 | 0 | | 12.0 | | 142011 | 0 | , | | 11.0 | uai | 4 | U | | 0.1 | |
| ÷ | | | | | | | | | | | | | | | | |
| | 6 | entre "o | quijote | sancho' | | | | | | | | | | | | |
| pasar | | 4 | 0 | 12.3 | | | | | | | | | | В | ack to t | he original |

Los colores **verde y rojo** en la Figura 30 corresponden a cada uno de los lemas introducidos. El verde es el que indica la palabra a la izquierda ("quijote") y el rojo se asocia con la palabra a la derecha ("sancho"). El grado de degradación del color se asocia con la probabilidad de compartir patrones. Cuando se degrada el color significa que la colocación es menos cercana o menos típica de la palabra

@ 0 8 0

en cuestión y cuando es **blanco** es que es **común** a ambas palabras. Es decir, cuanto más intenso es el color más propia y distintiva es esta colocación para la palabra en cuestión.

Los resultados de Diferencia Sketch, además de por colores, también se muestran en los valores de las tablas. Las tablas de los patrones comunes a los dos lemas presentan 4 cifras al lado de cada una de las colocaciones. Las primeras dos cifras indican la **frecuencia de coocurrencia** con el primer lema y el segundo lema respectivamente. Las últimas dos cifras indican lo **distintivo** que es la **colocación** (*salient score*) respecto a cada lema. Las colocaciones se organizan en función del máximo de los dos índices de distinción (*salient scores*) y la coloración refleja la diferencia entre los índices (*scores*).

Esta herramienta es muy parecida a la de *Word Sketch* en cuanto a que las distintas tablas tienen encima iconos similares. Además, cada elemento de cada tabla tiene un icono ... a la derecha que permite abrir el resultado en otras herramientas de *Sketch Engine* como son la herramienta de *Concordancia*, la de *Word Sketch* y la de *Tesauro*. Además, en la parte superior tenemos unos iconos que son los mismos que los de *Word Sketch* y tienen la misma utilidad. Entre ellas cabe destacar el icono de *Show visualization* que da una representación gráfica para las tablas de resultados, y que se muestra en la Figura 31 para la tabla de modificadores de los lemas "quijote" y "sancho".



4.3. - Tesauro: Creación de un tesauro automático

Esta función calcula las palabras o **lemas que suelen aparecer con las mismas colocaciones** que una palabra dada. Basándose en estos cálculos se genera automáticamente un "**tesauro distribucional**" (*distributional thesaurus*) que recoge las palabras que aparecen en contextos similares a los de la palabra seleccionada.

Al clicar en el botón de la herramienta Tesauro, nos aparece un menú como el mostrado en la Figura 32 (como es habitual, nos centramos en la pestaña *AVANZADO*). Como puede verse, es un menú muy simple en el que únicamente se puede introducir el lema que se quiera buscar, además de escoger, si se desea, la **categoría gramatical**. Permite también decidir sobre si se agrupan los ítems similares y los



valores mínimos de puntuación necesarios para aparecer en el teasuro, así como el máximo de elementos que se mostrarán.

| Lema Categoría gramatical auto adjective adverb noun verb preposition Purtuación minima tesauro 0 € | BÁSICO A | AVANZADO | | - |
|--|--|----------|--|---|
| auto Número máximo de ítems adjective 100 🔄 adverb 100 🔄 noun Purtuación mínima tesauro verb 0 😒 | lema Catagoría gram | | Agrupar items similares | Ð |
| | adjective adverb noun verb preposition | | Número máximo de ítems <u>100</u> Purtuación minima tesauro <u>0</u> $\widehat{\bullet}$ | |
| jYA! | | | YA! | |

Por ejemplo, si se introduce la palabra "rocinante" y se genera el tesauro, como en la Figura 33, el resultado consiste en una lista de lemas (*caballo, asno, rucio...*) que tienen un comportamiento similar a la palabra "rocinante" con respecto a los patrones gramaticales y a las colocaciones.

Por lo demás, aparecen unos iconos que tienen un funcionamiento idéntico a los ya explicados para la herramienta de *Word Sketch*. Cabe resaltar el botón de *Mostrar visualización*, que crea muestra los resultados en una gráfica como la mostrada en la Figura 34.

| | cinante as noun 205× | | | | | | | | | | | | |
|------------|----------------------|-------|--------------------|----------------------|-----|---------------------|--------------------------|-----|---|------|----------|-------------|--------|
| 5 | | | | | | | | | হ | ± 0 | | (j) | ☆ |
| | | | Word | Listas de frecuencia | s ↓ | Similarity A | Agrupar | | | | | | |
| = | | | caballo | 22 | 2 | 0,14 | | | | | | | |
| | | | 2 asno | 10 | 0 | 0,118 | jumento 84 | | | | | | |
| 0 | | | ³ rucio | 12 | 0 | <mark>0,11</mark> 5 | | ••• | | | | | |
| •= | | | letra | 8 | 1 | 0,111 | | | | | | | |
| 5-5 | | | alcornoque | 1 | 7 | 0,108 | estribo 16 peña 32 | | | | | | |
| =•= =•= | | | 6 tabla | 2 | 4 | 0,091 | | | | | | | |
| t≡ | | | 7 montaña | 4 | 4 | 0,084 | | | | | | | |
| NE | | | alabio | 2 | 4 | 0,08 | | | | | | | |
| ۵≡ | | | 9 mulo | 6 | 0 | 0,08 | | | | | | | |
| 13 | | | 10 hombro | 4 | 1 | 0,079 | | | | | | | |
| | Líneas por pág | gina: | 10 ▼ 1–10 de 4 | 48 I< < | 1 | > >I | | | | Back | o the or | iginal inte | erface |



4.4. - Concordancia: Ejemplos de uso en contexto

Al clicar en la opción *Concordancia* del menú de la izquierda, se entra en la herramienta de búsqueda de concordancias. La versión más básica de esta herramienta (pestaña *BÁSICO*) tan solo acepta un parámetro: la forma que se quiere encontrar en el texto del corpus. La versión de la pestaña *AVANZADO* es la que se muestra en la Figura 35.

| BÁSICO AVANZADO | | |
|-------------------------|------------|---|
| Tipo de consulta | simple | Ð |
| simple | abc | |
| lema | | |
| frase | | |
| word | | |
| character | | |
| CQL | | |
| | | |
| Subcorpus | | |
| ningún (todo el corpus) | • + | |
| Filter context 🗸 | | |
| Tipos de texto 🗸 | | |
| | IAYI | |
| | | |



El menú de *Filter context* permite filtrar por un contexto basado en lemas o basado en categorías gramaticales. En ambos casos, la idea es que se pueden rechazar los resultados que contengan, o no, un cierto lema o una forma de una cierta categoría gramatical en un entorno prefijado de tokens alrededor de la expresión que queramos buscar.

Debajo, clicando en *Tipos de texto*, *Sketch Engine* permite la **elección de una parte del corpus** para hacer la búsqueda. Aparte de los subcorpus que hemos visto, hay algunos corpus que contienen **metadatos** en los archivos de texto que conforman su contenido. Uno de los metadatos más típicos es la **fecha de publicación**, pero existen otros como la **web** de la que se ha extraído el texto, o el **dominio de país** de la página web (como .es para páginas españolas, o .mx para páginas mexicanas). Sin embargo, esto no ocurre para todos los corpus, con lo que las opciones que se muestran al desplegar el menú de tipos de texto no son siempre las mismas.

4.4.1 - Tipos de consultas de concordancia

4.4.1.1 - Concordancia simple

La Figura 35 muestra la búsqueda de concordancias por defecto, que es la búsqueda simple. Tal y como se ve en dicha figura, basta con introducir una o más formas o lemas⁸ en el campo *abc*, y clicar en *¡YA*!.



⁸ En principio se pueden introducir uno o más lemas o formas flexionadas. El funcionamiento más simple es cuando se introduce una única forma flexionada: la concordancia se hará únicamente con dichas formas concretas. Pero si se introduce un lema, la búsqueda incluye a todas las posibles formas derivadas del lema. Igualmente, si se introduce más de una palabra (lema o forma), la búsqueda tratará de encontrar esa combinación de palabras, tratándolas como formas, y limitándose a ellas (aunque coincidan con los lemas).

El resultado, mostrado en la Figura 36, es una pantalla con todas las **apariciones de esta palabra en el corpus**. Por ejemplo, si se introduce el lema "molino", se obtienen todas las **líneas del corpus** (en este caso, 23) en las que aparece "molino" o "molinos". Cada línea que aparece en los resultados contiene la siguiente información:

- El **orden** de aparición **del lema/forma** (o combinación de varias formas) que se muestra. Así, el primer resultado tendrá la posición 1, el siguiente la 2, y así sucesivamente.
- Un botón que muestra un menú de información sobre los datos de esta aparición en concreto.
- El orden del token mostrado dentro de todo el corpus.
- El contexto izquierdo (palabras previas a la coincidencia encontrada).
- La propia **forma** que coincide con el patrón de búsqueda (destacada **en rojo**). Esta palabra es clicable, y al clicar en ella, muestra un contexto más amplio que el mostrado en el resultado.
- El contexto derecho (palabras posteriores a la coincidencia encontrada).

| tipo de búsqueda | tipo de búsqueda |
|------------------|---|
| ✓ s | 🗸 word (1) |
| | tag lempos gender_lemma tags morphemes word (lowercase) lemma shorttag |
| | |

Se puede observar que, como ya hemos visto en las otras herramientas, hay una serie de iconos en la parte superior derecha que dan acceso a una serie de ajustes con respecto a los resultados mostrados. Varios de ellos son ya conocidos, pero existen unos cuantos que son propios de esta herramienta o que tienen menús distintos a los ya vistos:

- **Deshacer última acción**: Como veremos más adelante, varias de las opciones disponibles cambiarán el resultado mostrado. Con este botón podremos deshacer el último cambio realizado.
- Opciones de visualización: Ya se ha visto previamente este botón, y esta vez muestra un menú como el que aparece en la Figura 37. En la Figura 36, cada una de las palabras del corpus aparecía sin ninguna información extra. Pero se puede mostrar cada palabra acompañada de ciertos datos como:
 - La etiqueta sintáctica individual con la que se etiqueta cada una de las formas.

- El **lema-categoría** del que deriva. Por ejemplo, la forma "amaneceres" deriva del lemacategoría *amanecer-Nombre*, mientras que "amanecerá" deriva de *amanecer-Verbo*.
- El lema manteniendo el género de la forma. Para palabras con variación en género y que aparecen en género femenino, se muestra el lema en femenino sin pasarlo al masculino como suele ser habitual. Así, "hermosas" se lematizará como hermosa, mientras que "hermosos" se lematizará como hermoso.
- Las etiquetas sintácticas múltiples que pueda tener la forma. Ocurre en algunas formas que pueden etiquetarse con más de una etiqueta (en español esto ocurre en las formas verbales con clíticos y las contracciones). Así, la forma "decírselo", tiene una etiqueta individual de VMN0000, pero una etiqueta multiple en la que cada clítico tiene su propia etiqueta: VMN0000, PP3CN00, PP3MSA0.
- Los morfemas. En realidad, esto se aplica al mismo tipo de formas que en el punto anterior, de manera que para la forma "del" muestra su división en lo que Sketch Engine denomina morfemas: de, el.
- La propia **palabra en minúscula**.
- o El lema del que deriva la forma.
- La etiqueta corta. Esta etiqueta consiste simplemente en el primer carácter de la etiqueta completa (que puede tener hasta un máximo de 7 caracteres para verbos y nombres en español), que es el que identifica la categoría gramatical.

También se puede decidir si esta información se muestra únicamente para la palabra para la que hemos solicitado la concordancia, o para todas las palabras incluyendo ambos contextos. Además, dicha información puede mostrarse directamente en pantalla bajo cada una de las formas, o permanecer oculta hasta que ponemos el puntero del ratón sobre la forma. Por último, se puede ocultar el número de orden que aparece a la izquierda de cada línea.

- **Tomar una muestra aleatoria**: Muestra una cantidad limitada de resultados tomándolos aleatoriamente de entre todos los resultados existentes.
- Ordenar: El orden por defecto en el que se muestran los resultados es el de aparición en el corpus. Con esta opción se permite ordenar (alfabéticamente) no solo según la palabra central sino también según palabras del contexto situadas hasta tres posiciones antes o después.
- *Filtrar*: Con esta opción se permite crear reglas de filtrado basadas en la aparición de una forma, lema, carácter... situado en el contexto de hasta cinco palabras previas o posteriores.
- Buenos ejemplos para diccionario: Se muestran primero aquellos ejemplos de uso que se consideran buenos ejemplos para diccionario⁹. Se considera que una aparición de una forma es un buen ejemplo de uso de dicha forma, si aparece en una frase de 10-25 palabras en la que no aparecen otras palabras raramente usadas, a ser posible sin anáforas, si la forma aparece en una colocación típica...
- Frecuencia: Se muestran los datos de frecuencia de aparición (apariciones por millón de formas) de la palabra clave buscada, así como de las palabras que la rodean en el corpus. Se puede calcular la frecuencia según la etiqueta morfosintáctica, la forma o los lemas. En el

⁹ Jakubíček, M., Kilgarriff, A., Kovář, V., Rychlý, P. y Suchomel, V. (Eds.) 2013. The tenten corpus family. Kilgarriff, A., Husák, M., McAdam, K., Rundell, M., & Rychlý, P. (2008). GDEX: Automatically finding good dictionary examples in a corpus.

ejemplo mostrado en la Figura 38 se muestran las diez etiquetas con mayor frecuencia para el lema "comer".

| Ø | CONCORDANCIA | El ingeniero hidalgo 🔍 🕡 | Get more | space 🕂 🕒 📀 | P : |
|----------|--------------------------------|--------------------------|------------------------|-----------------|----------------------|
| | | ର୍ 🛨 🗠 💿 🖑 | 🖗 🗶 📻 🖶 🛃 🖪 | •••• 11. KWIC - | + ☆ |
| 88 | Listas de frecuencias | CAMBIAR CRITERIOS VOLVER | A CONCORDANCIA | | |
| : | 1 | | | | |
| \odot | Mostrar frecuencias por millón | | | | |
| 00 | Tag | ↓ Listas de frecuencias | Freq. / M | | |
| •= | 1 VMN0000 | 100 | 230,11 | | ••• |
| 202 | 2 VMP00SM | 20 | 46,02 | | •••• |
| | 3 VMIP3S0 | 15 | 34,52 | | •••• |
| 1 | 4 VMIP3P0 | 10 | 23,01 | | •••• |
| ţ≡ | 5 VMG0000 | 10 | 23,01 | | |
| NE | 6 VMSP3S0 | 7 | 16,11 💻 | | |
| ≣۶ | 7 VMIS3P0 | 7 | 16,11 💻 | | |
| | 8 VMII3S0 | 6 | 13,81 💻 | | |
| ~ | 9 VMSI3S0 | 5 | 11,51 💻 | | |
| | 10 VMSP3P0 | 4 | 9,20 🗖 | | |
| | | Líneas por página: | <u>10</u> ▼ 1–10 de 32 | K K <u>1</u> | > > |
| | | | | Back to the | e original interface |

Figura 38: frecuencias de aparición de las distintas etiquetas morfosintácticas para el lema "comer"

| lema molno 23 (sz.s | h por M) | | | | | | | Q 1 10 0 | 4 × = = | 맘 🖬 🚥 | | KWIC • + | |
|---------------------|---------------|------------|---------|-------|-----------|-----|---------------|---------------|--------------------|---------|------|-----------|---|
| Colocacio | CAMBIAR CRI | | | ANCIA | | | | | | | | | |
| 001000010 | | | | | | | | | | | | | |
| Word | Coocurrencias | Candidatos | T-score | MI | 4 LogDice | | Word | Coocurrencias | Candidatos | T-score | MI | + LogDice | |
| 1 viento | 10 | 46 | 3,16 | 12,00 | 12,21 | | ti la | 7 | 10.052 | 2,44 | 3,72 | 4,51 | • |
| 🔋 rueda | 4 | 11 | 2,00 | 12,75 | 11,91 | *** | 12 y | 10 | 16.934 | 2,88 | 3,48 | 4,27 | • |
| 3 gigantes | 3 | 30 | 1,73 | 10,88 | 10,86 | | 13 que | 11 | 20.226 | 2,99 | 3,36 | 4,15 | • |
| 4 eran | 3 | 186 | 1,73 | 8,25 | 8,88 | | 14. no | 3 | 5.700 | 1,56 | 3,31 | 4,10 | |
| 5 cuando | 3 | 710 | 1,71 | 6,32 | 7,07 | | ts el | 4 | 7.826 | 1,79 | 3,27 | 4,06 | |
| una | 5 | 1.303 | 2,21 | 6,18 | 6,95 | | 16 . · | 14 | 35.011 | 3,25 | 2,92 | 3,71 | , |
| 7 los | 9 | 4.603 | 2,92 | 5,21 | 5,99 | | | | | | | | |
| • ; | 3 | 2.711 | 1,65 | 4,39 | 5,17 | | | | | | | | |
| 1 de | 18 | 18.578 | 4,01 | 4,19 | 4,99 | | | | | | | | |
| 👳 en | 6 | 7.730 | 2,28 | 3,87 | 4,66 | | | | | | | | |
| | | | | | | | | | Líneas por página: | 20 . | | 6 1 | |

Figura 39: colocaciones para la palabra "molino"

 Colocaciones: Se ofrece un menú en el que se puede escoger el entorno alrededor de la palabra clave en el que se buscan otros lemas que puedan formar potencialmente colocaciones con la palabra dada. Dado que se busca palabras que coaparecen frecuentemente con la palabra clave, se intenta eliminar las palabras muy comunes (como preposiciones o determinantes) penalizándolas. Para ello se dan varias fórmulas distintas para dicha penalización, utilizando distintos modelos estadísticos. El resultado para "molinos" se muestra en la Figura 39, utilizando los modelos por defecto (T-score, MI y logDice) y un entorno de ±5 palabras. Puede apreciarse como la colocación más habitual es "viento" usando el modelo logDice, mientras que es "rueda" si se usa MI y un sorprendente "de" cuando se usa el modelo T-score.

Distribución de ocurrencias en el corpus: Con este botón del menú, se puede obtener una "radiografía" de la aparición del término buscado en el corpus. Por ejemplo, podemos ver en qué partes del corpus aparece más a menudo el lema "comer". El resultado, mostrado en la Figura 40, indica que el lema está muy repartido a lo largo del corpus, si bien existen unas ciertas partes donde este lema es especialmente abundante. Puede cambiarse la granularidad: en el ejemplo dado, cada barra incluye un 0,5% del corpus, pero dicho valor puede oscilar entre el 0,1% y el 10%. Además, cada barra es clicable, y al hacer click se nos muestran los resultados de concordancias dentro de esa parte del corpus.



- KWIC: Es el acrónimo de KeyWord In Context. Con este menú podemos decidir si el entorno que se nos muestra en los resultados es el correspondiente a un cierto número de palabras previas y posteriores (cuya longitud se fija en las opciones de visualización explicadas más arriba), o si el entorno es simplemente la frase completa en la que aparece.
- **Crea subcorpus**: Con esta opción, se puede crear un subcorpus formado por las frases o documentos que contengan resultados de concordancia.

4.4.1.2 - Concordancia de lema

Para utilizar este tipo de consulta de concordancia se debe seleccionar *lema* en el menú de la Figura 35. Al hacerlo, el menú nos permite introducir un lema y también una **categoría gramatical** (cuando tiene más de uno). Así, la palabra "sobre" es un lema, que puede hacer referencia a una preposición o

0

a un nombre¹⁰, o "anochecer" puede ser un verbo o un nombre. Si no se especifica la categoría, se escoge la más frecuente (en caso de que en la consulta simple se introduzca un lema, dicha consulta sería equivalente a esta consulta de concordancia de lema).

Por lo demás, este tipo de búsqueda es muy similar a la explicada en la subsección anterior y el tipo de resultados y opciones de visualización son los mismos.

4.4.1.3 - Concordancia de frase (sintagma)

La opción *frase* de la Figura 35 nos permite hacer una búsqueda por **sintagmas**. Se pueden introducir varias palabras y se obtendrán las apariciones de exactamente dichas formas (se consideran formas flexionadas aunque coincidan con lemas), exactamente en el orden dado, y sin ninguna otra palabra entre medias. Al no dar ninguna información de categoría gramatical, *Sketch Engine* devuelve todas las apariciones de dicha combinación de palabras independientemente de su categoría gramatical. Si se busca "la vista", se devolverían resultados para la combinación de artículo + nombre ("ha perdido la vista"), pero también para la de pronombre + verbo ("necesito que la vista").

Este tipo de búsqueda puede conseguirse también si en la búsqueda simple se introducen los (varios) términos que se pretende encontrar. Por lo demás, la forma de mostrar los resultados es la misma que en los modos de búsqueda descritos anteriormente.

4.4.1.4 - Concordancia de forma

La opción *word* es similar a las anteriores, con la diferencia de que limita el número de resultados a la **forma concreta** dada (aunque coincida con un lema) pero a la vez, da la capacidad de elegir el **tipo de palabra** (categoría gramatical) que se desea. Así, se puede buscar el término "bajo" como adjetivo, en cuyo caso devolvería concordancias únicamente con dicha forma (solamente en masculino singular), o como preposición.

4.4.1.5 - Concordancia de caracteres

Para este tipo de búsqueda (usando la opción *character* en el menú de *Query type*) se devuelven los resultados al buscar la ocurrencia de la **cadena de caracteres** dada, aunque **no se trate de una palabra completa**. Para los tipos de búsquedas vistos hasta ahora, si se busca por "vista" solo se devolverán resultados para esta palabra en concreto. Si esta misma búsqueda se utiliza en el campo *character*, también se devolverán resultados para otras palabras como "revista" o "vistazo". Igualmente, si se busca por ejemplo el término "ábamos", se nos devolverán las terminaciones de las formas verbales que sean de primera persona del plural del pretérito imperfecto de indicativo de verbos de la primera conjugación. No obstante, existen formas mejores de conseguir este resultado.

4.4.1.6 - Concordancia CQL

Este es el tipo de búsqueda más completo, ya que permite definir de forma muy precisa los resultados que se buscan. Si se escoge la opción *CQL* (de *Corpus Query Language*), el texto a introducir en el campo *CQL* será una expresión compuesta por elementos del tipo **[atributo="valor"]**, donde el valor puede contener además expresiones regulares. Existen varios atributos, entre los que destacan *lemma* (lema), *word* (forma flexionada) y *tag* (etiqueta morfológica), pero en realidad existen tantos atributos

¹⁰ También podría considerarse como una forma flexionada del verbo "sobrar", pero en dicho caso no sería un lema, sino una forma.



como los que se han visto en las opciones de visualización de los resultados, tal y como puede verse en la Figura 41.

Si bien esta búsqueda es la que proporciona mayor libertad a la hora de escoger qué concordancias en concreto queremos encontrar, a cambio tiene cierta complejidad en su sintaxis ya que hace uso de expresiones regulares básicas y de etiquetas morfosintácticas. Debido a ello, *Sketch Engine* proporciona un motor de creación para estas expresiones, además de incluir un pequeño tutorial de vídeo acerca de su uso tal en el propio menú, tal y como puede observarse en la Figura 41. Por ejemplo, la expresión:

```
[word=".*[aeií]r" & tag="[^V].*" ]
```

devuelve todas las palabras que acabaran en –ar, -er, -ir, o –ír que no sean verbos (como "lugar", "mujer" o "Sir"). Igualmente, la expresión

[morphemes="meter,.*"] [shorttag="S" & word!="en"] [tag="D[ADP].*"] []? [lemma="asunto"]

devuelve como resultados expresiones que:

| 3 | CONCORDANCIA | El ingeniero hidalgo 🔍 | i | Get more space 🕀 | ප ? 📁 | : |
|-------------|---|--|-------------------------------------|--|--|--------|
| | BÁSICO AVANZADO | | | | | |
| • • • | Tipo de consulta simple lema frase word character CQL | caL [lemma + "libro"] [] (l Insert [] [] <> "" Atributo predeterminado Word lemma tag lempos gender_lemma tas | ,'3} [tag= "V_*"] € \ 1 ~ # TAGS | CQL 1: CON an introduc corpus duran 55 EK | plex c () () () () () () () () () | 0 |
| | none (the whole corpus) | word (lowercase) shorttag | <u> </u> | | | |
| | Tipos de texto 🗸 | | jYAI | | | |
| الآم ا | | | | (| Back to the original inte | erface |

- Comiencen con el lema "meter" (que solo puede ser verbo) conjugado en cualquiera de sus formas y estén seguidos de al menos un clítico.
- Tengan a continuación una preposición que no sea "en".
- Continúe con un determinante que sea un artículo, posesivo o demostrativo.
- Pueda llevar un token de cualquier tipo a continuación.
- Termine con una forma flexionada del lema "asunto".

Con esta expresión capturaríamos expresiones como "métase con el malditos asuntos" o "metiéndolo a ese asunto", que indicarían variaciones de la colocación típica "meterse en mis asuntos". Estas búsquedas pueden tardar varios minutos en realizarse si el corpus es muy grande.

4.5. - Listas de palabras: frecuencias de aparición

Para obtener la lista de palabras de un corpus junto con las frecuencias de cada una, se puede clicar en *Lista de palabras*, que aparece en la parte superior del menú izquierdo. En el menú que aparece en la Figura 42, puede apreciarse que hay diferentes opciones para obtener la lista de palabras y su frecuencia. La más básica se corresponde con *word*, que hace recuentos de apariciones de palabras diferenciando por sus formas. Sin embargo, existen multitud de variantes, comenzando por el recuento a nivel de lemas, o de etiquetas sintácticas, por ejemplo, además de las distintas categorías gramaticales.

Puede hacerse el recuento de los *tokens* que empiecen o terminen por una cierta cadena de caracteres, o que la contengan. También pueden usarse **expresiones regulares** para definir las ocurrencias que entrarán en el recuento, y puede incluirse también una lista de palabras concretas que no deben incluirse en el recuento aunque cumplan el resto de condiciones. Incluso pueden eliminarse del recuento los *tokens* que no sean palabras (como símbolos y signos de puntuación).

Por último, existe la posibilidad de seleccionar cómo se agrupan los recuentos. Por defecto se muestra una lista simple, en la que se tiene una única columna de lema, forma, etiqueta... y otra columna con las frecuencias. Sin embargo, si se usa *Mostrar como* es posible mostrar más de una columna, con lo que se puede mostrar, por ejemplo, un listado de formas, con su correspondiente etiqueta morfosintáctica y el lema, seguido de su frecuencia.





4.6. - N-gramas: extracción de expresiones multipalabra

Con esta herramienta se puede analizar el corpus para encontrar las expresiones multipalabra más típicas. Al clicar el botón de N-gramas¹¹, se muestra el menú que aparece en la Figura 43. En el menú se puede escoger el **Tamaño de N-gramas**, especificando un rango de longitud de las expresiones multipalabra que se buscarán (por defecto 2 o 3). Es posible buscar recuentos de formas, lemas, etiquetas... hasta los ocho tipos que ya se han visto en las opciones de visualización de la Figura 37. También se puede restringir la búsqueda de la expresión multipalabra indicando las letras por las que debe empezar, usando expresiones regulares, o indicando las palabras con las que debe empezar o acabar.

| BÁSICO AVANZADO | Ð |
|--|---|
| Tamaño N-gramas 2 3 4 Word Listas de frecuencias min. 5 0 | 5 6 Anidar n-gramas Incluir no palabras 1 A = a excluir siguientes palabras 1 A = a excluir siguientes palabras 1 Subcorpus • 1 • none (the whole corpus) • |
| todosstarting with lettersending with letterscontaining lettersstarting with wordcontaining wordending with wordmatching regular expressionde esta lista | |
| Figura 43: me | iYA! enú de la herramienta N-gramas |

También se da la opción de anidar los resultados obtenidos (de forma que se muestren juntas las expresiones que empiecen por las mismas palabras, independientemente de que por su frecuencia deberían aparecer más tarde), incluir *tokens* que no sean palabras, diferenciar o no entre mayúsculas o minúsculas, y excluir términos multipalabra que incluyan alguna forma de las incluidas en una lista.

Tras clicar en el botón *¡YA!* se nos muestra la lista de resultados. En la Figura 44 se muestra el listado de las 50 expresiones de cuatro lemas más frecuentes en el corpus de El Quijote.

¹¹ Un n-grama es simplemente un conjunto de n palabras que aparecen de forma consecutiva en el corpus. Cuando n es igual a 2 se habla de bigramas, y cuando es igual a 3, de trigramas. Para números de n superiores a 3 (cuyo uso no es muy común), se habla de 4-gramas, 5-gramas... y así sucesivamente.

| Word | ↓ Count | | Word | | | Word | | Word | + C | ount | |
|-------------------|--------------|------|----------------------------|--------------|---|-----------------------------|----|----------------------------|-----|------|-----|
| don quijote de | əl 127 | | 14 el que haber de | 36 | | 27 mi señor don quijote | 27 | 40 en todo el día | | 23 | |
| quijote de el m | ancha 118 | ••• | 15 caballero de el triste | 34 | | 28 el cura y el | 27 | el duque y el | | 23 | |
| dar a entender | que 47 | | 16 a el que responder | 34 | | 20 de el andante caballería | 26 | 42 que no parecer sino | | 22 | |
| de el triste figu | a 44 | | 1/ que ser el que | 32 | | 30 que ser el más | 25 | 43 que me parecer que | | 22 | ••• |
| señor dulcinea | de toboso 43 | | 18 don quijote y sancho | 32 | | 31 no parecer sino que | 25 | 44 que haber de ser | | 22 | |
| el caballero de | el 43 | | 19 no se haber de | 31 | | 😰 a don quijote y | 25 | 45 el puerta de el | | 22 | |
| ser uno de el | 42 | | 20 así ser el verdad | 31 | | 33 uno de el más | 24 | 4 de allí a poco | | 22 | |
| en el mitad de | 40 | | 21 que deber de ser | 30 | | 34 que a mí me | 24 | 47 ser el de el | | 21 | |
| decir a este sa | tón 40 | | 😂 que se haber de | 29 | | 36 mi señor el duque | 24 | 48 se haber de hacer | | 21 | |
| 10 haber en el mu | ndo 39 | | a no poder dejar de | 28 | | 36 el señor dulcinea de | 24 | 49 que haber de hacer | | 21 | |
| qué ser el que | 37 | | 21 de el caballero andante | 28 | | 37 de el que haber | 24 | 50 no haber más que | | 21 | |
| 🔹 que no haber o | e 36 | | 25 todo el día de | 27 | | 36 cura y el barbero | 24 | | | | |
| no haber para | qué 36 | | 75 no haber de ser | 27 | | 39 a el puerta de | 24 | | | | |
| | | Líne | eas por página: 50 💌 | 1-50 de 1.46 | 1 | IC C <u>1</u> > >I | | | | | |

4.7. - Palabras clave: extracción de palabras clave y términos multipalabra

Por último, se tiene la herramienta *Palabras clave* que hace un análisis a nivel del corpus completo (tiene que ser un corpus creado por el usuario) y extrae de él las **palabras clave** más características, así como los **términos multipalabra** más propios de dicho corpus. Esto lo hace analizando el texto del corpus del que se quieren extraer las palabras clave, y **comparando** dicho análisis con un **corpus equilibrado**, que representa múltiples temas y múltiples variedades de la lengua. Así, se identifican aquellas palabras o grupos de palabras que resultan frecuentes en el corpus que está siendo analizado, pero que no lo son tanto en un corpus mayor y con contenidos más variados.

| Foreign and the second s | |
|---|--|
| 1 Por lo menos uno alfanumérico | V Sólo alfanuméricos |
| Keywords (single-words) settings | erms (multi-words) settings |
| Reference corpus Ref Spanish Web 2011 (esTenTen11, Eu + Am) Span | eference corpus anish Web 2011 sample Q |
| Reference subcorpus Ref ningún (todo el corpus) Q nin | ference subcorpus ngún (todo el corpus) |
| Número máximo de items Attribute for keywords Núm 1000 Image: Iemma 100 | nero máximo de ítems 00 \$ |
| correspondiente a este regex corre | respondiente a este regex |

Al clicar en el botón de esta herramienta se nos muestra el menú que aparece en la Figura 45. Entre otras cosas, en él se puede escoger cuál es el corpus de referencia frente al que se compara nuestro corpus privado, y algunas características de los términos que se extraerán.

Tras pulsar en el botón *¡YA!* se obtienen los resultados, que se muestran en la Figura 46. Como puede apreciarse, al estar analizando un corpus que consiste en las dos partes de El Quijote, las palabras clave son, fundamentalmente, nombres propios usados en el universo de El Quijote, junto con conceptos típicos como "caballero andante" o "libro de caballería".

@080

| | . , | | OLAVE | | | · | | | | |
|---|-----|--------------|---------------------|------------------|-----|----|----------------------|----------------|--------------------|------------|
| | | | | _ | | | | | र् ± ⊙ (|) 12 |
| | | | SINGLE-WOF | RDS | | | MU | LTI-WORDS | 0 | |
| | | Word | Frecuencia 1 | Reference corpus | | | Word | Frecuencia 1 F | Reference corpus | |
| | 1 | sancho | 2.167 | 33.724 | | 1 | don quijote | 2.148 | 888 | |
| | 2 | quijote | 2.159 | 60.447 | | 2 | sancho panza | 278 | 95 | |
| | | dulcinea | 279 | 3.704 | | 3 | caballero andante | 222 | 82 | |
| = | 4 | andante | 351 | 8.089 | | 4 | señor don | 161 | 103 | |
| = | 5 | vuesa | 202 | 1.090 | | 5 | señor don quijote | 128 | 4 | |
| | 6 | rocinante | 205 | 1.867 | | 6 | señora dulcinea | 89 | 4 | |
| | 7 | toboso | 162 | 1.522 | | 7 | quijote de la mancha | 110 | 153 | |
| | 8 | lotario | 142 | 608 | | 8 | señor mío | 61 | 44 | |
| 1 | 9 | ventero | 148 | 1.301 | | 9 | don fernando | 131 | 288 | |
| | 1 | rucio | 131 | 1.714 | | 10 | señor caballero | 45 | 6 | |
| | 1 | luscinda | 99 | 199 | | 11 | libro de caballería | 40 | 117 | |
| | 1 | 2 cardenio | 102 | 635 | | 12 | señora mía | 36 | 9 | |
| | 1 | insula | 144 | 5.472 | | 13 | doña rodríguez | 35 | 0 | |
| | 1 | doncella | 217 | 14.678 | | 14 | sansón carrasco | 33 | 4 | |
| | 1 | o dorotea | 114 | 3.232 | | 15 | señor gobernador | 33 | 224 | |
| | 1 | e merced | 969 | 118.282 | | 16 | cide hamete | 35 | 16 | |
| | 1 | 7 panza | 330 | 33.587 | | 17 | don antonio | 62 | 311 | |
| | 1 | jumento | 84 | 864 | | 18 | sancho amigo | 26 | 0 | |
| | 1 | ∍ priesa | 72 | 437 | | 19 | andante caballero | 26 | 1 | |
| | 2 | barbero | 172 | 16.327 | | 20 | señora la duquesa | 24 | 0 | |
| | | Líneas por p | xágina: <u>20 ▼</u> | 1–20 de 1.000 | K K | 1 | > >I | | Back to the origin | nal interf |

4.8. - OneClick Dictionary: creación de un diccionario automático

Con esta herramienta, *Sketch Engine* es capaz de crear un esbozo de diccionario utilizando para ello las concordancias en el corpus. Esta herramienta, que necesita una cuenta en la web https://www.lexonomy.eu/, hace un análisis gramatical del corpus para extraer atributos de términos e inferir su significado.

Como su propio nombre indica, basta con clicar en la herramienta para que se genere el diccionario utilizando el corpus seleccionado.

4.9. - Corcondancia paralela: expresiones equivalentes en dos lenguas

Con esta herramienta se utiliza un corpus paralelo en varias lenguas para extraer las expresiones equivalentes en dos lenguas. El menú es similar al de la herramienta de *Concordancia*, con la diferencia de que se realiza en dos lenguas. Se debe escoger un término en una lengua origen, y una lengua

destino. Una vez escogida la expresión que se quiere traducir a otra lengua, se obtiene un resultado como el mostrado en la para la palabra española "toro" y su equivalente en búlgaro.

| CONCORDANCIA PARALELA EUR-Lex Spanish 2/2016 | Get more space 🕘 🕢 🔁 🧷 🛃 |
|--|--|
| simple toro 2.156 (2,66 por M) | 오 🛓 🗠 🐵 🖉 🗙 (align 🕶) |
| F F 🗄 | 🚎 \Xi EUR-Lex Bulgarian 2/2016 |
| Considerando que es necesario distinguir entre la admisión a la inseminación artificial de los teres de raza selecta y su semen que han sufrido todas las pruebas del examen oficial previsto para su raza en un Estado miembro y la admisión de los toros y su semen únicamente para fines de experimentación; | като има предлид, че е необходимо да се прави разлика между одобрение за изкуствено осеменяване на чистопородни бикове и тяхната сперма, които са преминали през официалните тестове за тяхната порода, предвидени в определена държава-иленка и одобрението на бикове и сперма от тях, единствено за целите на тестването; «З> |
| <s> Considerando que es necesario distinguir entre la admisión a la inseminación artificial de los toros de raza selecta y su semen que han sufrido todas las pruebas del examen oficial previsto para su raza en un Estado miembro y la admisión de los toros y su semen únicamente para fines de experimentación ;</s> | като има предвид, че е необходимо да се прави разлика между одобрение за изкуствено осеменяване на чистопородни бикове и тяхната сперма, които са преминали през официалните тестове за тяхната порода, предвидени в определена държава-ипенка и одобрението на бикове и сперма от тях, единствено за целите на тестването; «бо |
| <s> Considerando que es deseable que los toros de raza selecta y su semen sean identificados por el análisis del grupo sanguíneo de dichos toros o por cualquier otro método adecuado; </s> | «s> като има предвид, че е желателно чистопородните бикове и тяхната сперма да се идентифицират по кръвна група или други подходящи методи; |
| <c>> Considerando que es deseable que los toros de raza selecta y su semen sean identificados por el análisis del grupo sanguíneo de dichos toros o por cualquier otro método adecuado ; </c> | «so има предвид, че е желателно чистопородните бикове и тяхната сперма да се идентифицират по кръвна група или други подходящи методи ; |
| Sin perjuicio de las normas de policía sanitaria, los Estados miembros velarán por que no se prohiban, limiten u obstacuilcen la admisión para la reproducción de las hembras de bovino de raza selecta ni la admisión para la cubrición natural de los toros de raza selecta. | Държавите-членки гарантират, че без да се засягат ветеринарно-санитарните правила няма да установляят забрани, с орзинчения или прелятствия за одобрението на чистопородни женски живот от рода на одрия рогат добитък за целите на развъждането, както и за одобрението на чистопород бикове за остствено сосменияване. |
| -so- la admisión, para someterse a una prueba oficial de toros de raza selecta o la utilización de su semen en los límites cuantitativos necesarios para la realización de dichas pruebas oficiales por ormanisme o asociariones autorizados. «So | - одобрението за официално тестване на чистопородни бикове или употребата на тихната сперма в количества, необходими на одобрените организации или асоциации за извършването на такива официални тестове. |

Figura 47: corcondancia paralela entre la palabra española "toro" y su correspondiente en búlgaro

4.10. - Tendencias: variación en el uso de palabras a lo largo del tiempo

Para poder usar esta herramienta es necesario tener seleccionado un corpus cuyos documentos estén etiquetados según su fecha de creación. De esta manera, *Sketch Engine* es capaz de comprobar qué palabras se usan más o se usan menos en los documentos más recientes en relación con los más antiguos. El menú de la herramienta permite decidir si mostrar la tendencia para formas, lemas... o cualquiera de las ocho opciones que se muestran en la Figura 35. El resultado para el corpus EUR-Lex en español aparece en las figuras a continuación: la Figura 48 muestra las palabras que se usan menos en documentos más recientes, mientras que Figura 49 se muestran las que se usan más.

| TE | | CIAS | EUR-Lex Spanish 2/ | 2016 | Q () | | | | Get more space 🕀 | ٥ | Ð | ? | 1 |
|----|-------------|-----------------------------|--------------------|-------|----------|---|-----------------|------------------|------------------|--------|-----|---------------|------------|
| | | | | | | | | | | | র | ± 0 | () |
| | Lemma | \mathbf{v} | Listas de frecuer | ncias | Muestra | | Lemma | \checkmark | Listas de fr | ecuenc | ias | Muestra | |
| 1 | avis | $\mathbf{N}_{\mathrm{exc}}$ | 99 | 9.016 | | | 11 señor | \mathbf{N}_{i} | | 7.6 | 42 | ~ | |
| 2 | juridique | N | 99 | 9.461 | | | 12 comunidad | \mathbf{N}_{i} | | 687.4 | 15 | | |
| | important | N | 101 | 1.119 | | | 13 sobrepasar | \mathbf{N}_{i} | | 27.6 | 84 | ~ | |
| 4 | comunidades | N | 207 | 7.904 | \frown | | 14 cee | \mathbf{N}_{i} | | 380.9 | 42 | \frown | |
| | finés | N | 11 | 1.633 | ~ | | 15 referente | \mathbf{N}_{i} | | 57.2 | 36 | ~ | |
| 6 | auténtico | N | 38 | 3.039 | \sim | | 16 constitutivo | \mathbf{N}_{i} | | 104.2 | 12 | | |
| 7 | conveniente | N | 45 | 5.096 | \sim | | 17 convenir | \mathbf{N}_{i} | | 87.5 | 89 | ~ | |
| 8 | sueco | N | 25 | 5.705 | ~ | | 18 repartir | \mathbf{N}_{i} | | 13.5 | 79 | \sim | |
| 9 | kilogramo | \mathbf{N}_{i} | 19 | 9.618 | \sim | | 19 cantidad | \mathbf{N}_{i} | | 292.8 | 14 | ~ | |
| 10 | subordinar | $\mathbf{N}_{\mathrm{eff}}$ | 12 | 2.783 | \sim | | 20 modalidad | \mathbf{N}_{i} | | 67.8 | 49 | ~ | |
| | | Líneas po | r página: 20 🔻 | 1–2 | 0 de 254 | K | < 1 | > >I | | | Ва | ck to the ori | ginal in |

Figura 48: principales tendencias negativas (palabras en desuso)

| 5 | TENDEN | CIAS | EUR-Lex Spanish 2/2016 | Q (j | | | | Get more space 🕁 🛛 🧕 | G | • | |
|----------|------------------|---------------|------------------------|------------|------|-----------------|--------------|----------------------|-------|------------------|---------------|
| | | | | | | | | | ৎ | ± o | i ☆ |
| | Lemma | \mathbf{V} | Listas de frecuencias | Muestra | | Lemma | \checkmark | Listas de frecuer | ncias | Muestra | |
| | 1 documentar | 1 | 11.043 | | | 11 bz | 1 | | 459 | ~ | |
| - | 2 rechtsanwälte | 1 | 613 | | | 12 mengozzi | 1 | 1 | .092 | | |
| 0 | 3 eficiencia | 1 | 40.768 | ~ | | 13 trazabilidad | 1 | 6 | i.034 | | |
| ō | 4 previsibilidad | 1 | 2.047 | | | 14 fi | 1 | 20 | .196 | \sim | |
| = | 5 reevaluar | 1 | 863 | | | 15 supervisión | 1 | 91 | .173 | \sim | |
| •= | 6 liderazgo | 1 | 3.468 | | | 16 auditoría | 1 | 65 | .963 | ~ | |
| | 7 demostrable | 1 | 748 | | | 17 operativo | 1 | 86 | 6.767 | ~ | |
| ■ | 8 maximizar | 1 | 3.049 | | | 18 remoto | 1 | 3 | .297 | | |
| | 9 fráncfort | 1 | 1.283 | | | 19 validar | 1 | 10 | .026 | ~ | |
| | 10 patrulla | 1 | 804 | ~ | | 20 remitente | 1 | 46 | 6.347 | ~ | |
|) | | Líneas por pá | igina: <u>20 ▼</u> 1-2 | 0 de 1.000 | < | < 1 | > >1 | | E | Back to the orig | ginal interfa |
| | | Figur | a 49: principal | es tende | ncia | s negativas | (palak | oras de moda) | | | |