

SAR

Práctica 4: NLTK

Para uso de la herramienta NLTK se recomienda la lectura de <http://www.nltk.org/book>

Como ayuda para resolver los ejercicios propuestos se ha elaborado el documento [Guia_NLTK](#).

Se debe entregar un programa Python que resuelva todas las acciones propuestas en cada ejercicio. Además, se deben responder por escrito al subir la tarea a la pregunta 12 del ejercicio 1 y a la pregunta 13 del ejercicio 3.

EJERCICIO 1.

Escribe las instrucciones de Python adecuadas para realizar las acciones propuestas en cada apartado donde se adjunta el resultado correcto de su ejecución, si procede.

1	Acceder al corpus en castellano cess_esp
2	Mostrar el número de palabras que contiene este corpus 192685
3	Mostrar el número de frases que contiene 6030
4	Obtener las frecuencias de aparición de los ítems que componen el primer fichero del corpus anterior. Un ítem es un par (key, value) donde key es la palabra y value es la frecuencia de aparición de la palabra. Visualizar los 20 más frecuentes. [('de', 23), ('la', 12), ('en', 9), ('y', 8), ('.', 6), ('-Fpa-', 5), ('-Fpt-', 5), ('para', 5), ('una', 5), ('EDF', 5), ('millones', 4), ('como', 4), ('con', 4), ('que', 4), ('megavatios', 3), ('*0*', 3), ('gas', 3), ('EAA', 3), ('central', 3)]
5	Obtener el vocabulario del primer fichero del corpus (ordenado por frecuencia). ['de', 'la', 'en', 'y', '.', '-Fpa-', '-Fpt-', 'para', 'una', 'EDF', 'millones', 'como', 'con', 'que', 'megavatios', '*0*', 'gas', 'EAA', 'central', 'a', 'por', 'México', 'principal', 'Altamira_2', 'Saltillo', 'potencia', 'euros', 'Río Bravo', 'el', 'natural', 'construcción', '495', 'se', 'dólares', 'Mitsubishi', 'licencias', 'invertir', 'centrales', 'cada', '28', 'Una', 'compañía', 'duración', '186', '194', 'explicó', 'prevé', 'empresa', 'encargará', 'estatal', 'ciclo', 'proyecto', 'grupo', 'francesa', 'empezar', 'construir', 'japonés', 'prevista', 'explotarla', 'previsto', 'Mitsubishi_Corporation', 'creada', 'quiso', 'es', 'su', 'al', 'Comisión_Federal_de_Electricidad', 'funcionar', 'mexicana', 'combinado', 'mayoritaria', 'pública', 'dos', 'cuya', 'asistente', 'años', 'pasará', 'combustible', 'El', 'mayo_del_2002', 'energía', 'Electricité_de_France', '25', 'revelar', 'del', '134', 'un', 'norte', 'poner_en_marcha', 'debe', 'EFE', 'venta', 'licitación', 'utilización', 'funcionará', 'japonesa', 'quedaron', 'no', 'La', 'intervendrá', '247', 'participación', 'producida', 'Electricidad_Aguila_de_Altamira', ':', 'red', 'pagó', 'accionista', 'Tampico', 'Tuxpán', 'tiene', 'en_virtud_de', 'electricidad', 'posteriormente', 'portavoz', 'compra', '51_por_ciento', 'anunció', 'eléctrica', 'acuerdo', 'hoy', 'Altamira', '1998', 'participaron', 'cuánto', 'primera', 'jueves', 'CFE', 'eléctricas']
6	Obtener de forma ordenada las palabras del vocabulario de longitud mayor que 7 y que aparezcan más de 2 veces en el primer fichero del corpus. ['megavatios', 'millones']

EJERCICIO 2.

Escribe un programa en Python para calcular cuántas veces aparecen las palabras what, when, where, who y why en cada una de las categorías del Corpus Brown (['adventure', 'belles_lettres', 'editorial', 'fiction', 'government', 'hobbies', 'humor', 'learned', 'lore', 'mystery', 'news', 'religion', 'reviews', 'romance', 'science_fiction']) como una lista donde para cada una de las 5 palabras tengamos la lista de categorías del corpus y la frecuencia de aparición de la palabra en dicha categoría.

Resultado en forma de lista:

```
[ 'what', [ 'adventure', 110, 'belles_lettres', 244, 'editorial', 84, 'fiction', 128, 'government', 43, 'hobbies', 78, 'humor', 36, 'learned', 141, 'lore', 130, 'mystery', 109, 'news', 76, 'religion', 64, 'reviews', 44, 'romance', 121, 'science_fiction', 27], 'when', [ 'adventure', 126, 'belles_lettres', 252, 'editorial', 103, 'fiction', 133, 'government', 56, 'hobbies', 119, 'humor', 52, 'learned', 227, 'lore', 182, 'mystery', 114, 'news', 128, 'religion', 53, 'reviews', 54, 'romance', 126, 'science_fiction', 21], 'where', [ 'adventure', 53, 'belles_lettres', 107, 'editorial', 40, 'fiction', 76, 'government', 46, 'hobbies', 72, 'humor', 15, 'learned', 118, 'lore', 97, 'mystery', 59, 'news', 58, 'religion', 20, 'reviews', 25, 'romance', 54, 'science_fiction', 10], 'who', [ 'adventure', 91, 'belles_lettres', 452, 'editorial', 172, 'fiction', 103, 'government', 74, 'hobbies', 103, 'humor', 48, 'learned', 212, 'lore', 259, 'mystery', 80, 'news', 268, 'religion', 100, 'reviews', 128, 'romance', 89, 'science_fiction', 13], 'why', [ 'adventure', 13, 'belles_lettres', 36, 'editorial', 10, 'fiction', 18, 'government', 6, 'hobbies', 10, 'humor', 9, 'learned', 20, 'lore', 25, 'mystery', 25, 'news', 9, 'religion', 14, 'reviews', 9, 'romance', 34, 'science_fiction', 4]]
```

EJERCICIO 3.

Escribe las instrucciones de Python adecuadas para realizar las acciones propuestas en cada apartado donde se adjunta el resultado correcto de su ejecución, si procede.

1	Cargar el documento "quijote.txt" en una única cadena una cadena unicode UTF-8
2	Mostrar todos los símbolos del documento ordenados por orden alfabético ! " ' () , - . 0 1 2 3 4 5 6 7 : ; ? A B C D E F G H I J L M N O P Q R S T U V W X Y Z] a b c d e f g h i j l m n o p q r s t u v x y z ¡ « » ¿ Á É Í Ñ Ó Ú à á é í ñ ó ù ú ü
3	Eliminar del texto los símbolos siguientes: ¡ " ' () , - . : ; ¿ ?] « »
4	Mostrar todos los símbolos del documento filtrado ordenados por orden alfabético 0 1 2 3 4 5 6 7 A B C D E F G H I J L M N O P Q R S T U V W X Y Z a b c d e f g h i j l m n o p q r s t u v x y z Á É Í Ñ Ó Ú à á é í ñ ó ù ú ü
5	Obtener el número de palabras y el número de palabras distintas del texto filtrado. Mostrar la 10 primeras y las 10 últimas en orden alfabético 381212 24480 10 16 1604 1614 1615 17 23 A ABC ACADÉMICO última últimamente últimas último últimos única único únicos útil útiles
6	Obtener las frecuencias de aparición de los ítems que componen el documento filtrado. Un ítem es un par (key, value) donde key es la palabra y value es la frecuencia de aparición de la palabra. Visualizar los primeros 20 ítems. [('que', 20549), ('de', 17997), ('y', 17166), ('la', 10202), ('a', 9532), ('el', 7962), ('en', 7907), ('no', 5787), ('se', 4690), ('los', 4681), ('con', 4053), ('por', 3779), ('las', 3423), ('le', 3396), ('lo', 3393), ('su', 3320), ('don', 2538), ('del', 2465), ('me', 2345), ('como', 2244)]
7	Crear un nuevo documento eliminando las stopwords del texto filtrado.
8	Obtener el número de palabras y el número de palabras distintas del texto sin stopwords. Mostrar la 10 primeras y las 10 últimas en orden alfabético 183251 24066 10 16 1604 1614 1615 17 23 ABC ACADÉMICO ACADÉMICOS última últimamente últimas último últimos única único únicos útil útiles
9	Obtener las frecuencias de aparición de los ítems que componen el documento sin stopwords. Visualizar los primeros 20 ítems. [('don', 2538), ('Quijote', 2164), ('Sancho', 2145), ('si', 1798), ('dijo', 1789), ('tan', 1219), ('ser', 1056), ('respondió', 1053), ('bien', 964), ('señor', 948), ('así', 905), ('merced', 900), ('sino', 694), ('dos', 672), ('pues', 639), ('decir', 577), ('caballero', 573), ('hacer', 535), ('aunque', 525), ('Dios', 518)]
10	Crear un nuevo documento sustituyendo cada palabra del texto sin stopwords por su raíz. Para ello se utilizará el stemmer snowball.
11	Obtener el número de palabras y el número de palabras distintas del nuevo documento. Mostrar la 10 primeras y las 10 últimas en orden alfabético 183251 10134 10 16 1604 1614 1615 17 23 abad abadej abades zoroastr zorr zorrún zuec zulem zumb zurd zurrón zuz ñud
12	Obtener las frecuencias de aparición de los ítems que componen el nuevo documento. Visualizar los primeros 20 ítems. [('don', 2656), ('quijot', 2180), ('sanch', 2158), ('si', 1966), ('dij', 1882), ('señor', 1812), ('respond', 1277), ('tan', 1243), ('hac', 1158), ('buen', 1115), ('así', 1095), ('bien', 1069), ('ser', 1057), ('dec', 967), ('caballer', 955), ('merc', 900), ('pues', 865), ('parec', 833), ('algun', 811), ('cos', 805)]
13	Justificad los resultados obtenidos en los pasos 5, 8 y 11.