



SOBRE MI

- Ingeniería Informática y Matemáticas en la Universidad Autónoma de Madrid
- Máster en Big Data y Data Science en Finanzas en AFI (1ª promoción)
- Data Scientist y Product Owner en BBVA AI Factory
- Data Product Manager en Stuart

It's Your Turn!



Text Mining y NLP... ¿para qué?



Google



Microsoft

amazon

Text Mining y NLP... ¿para qué?

- Análisis de texto
- Extracción de información
- Clasificación de textos
- Búsqueda de información relevante
- *Sentiment analysis*
- Generación automática de texto
- ...

Indice

- 1. *BACK TO BASICS***
- 2. LENGUA, 3º ESO**
- 3. DISTANCIAS**
- 4. TODO FUNCIONA MEJOR CON NÚMEROS**
- 5. CASO DE USO: ANÁLISIS DE SENTIMIENTO**

1. BACK TO BASICS



Estructuras de datos

FICHERO/DOCUMENTO



Estructuras de datos

FICHERO/DOCUMENTO

↳ PÁRRAFO (separación por \n)

“Primera parte del Ingenioso hidalgo Don Quijote de la Mancha.”

“Capítulo Primero. Que trata de la condición y ejercicio del famoso hidalgo don Quijote de la Mancha.”

“En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no hace mucho que vivía un hidalgo de los de lanza en astillero, adarga antigua, rocín flaco y galgo corredor. Una olla de más vaca que carnero, salpicón todas las noches, (...), consumían las tres partes de su hacienda.”

Estructuras de datos

FICHERO/DOCUMENTO

↳ PÁRRAFO (separación por \n)

“Primera parte del Ingenioso hidalgo Don Quijote de la Mancha.\n”

“Capítulo Primero. Que trata de la condición y ejercicio del famoso hidalgo don Quijote de la Mancha.\n”

“En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no hace mucho que vivía un hidalgo de los de lanza en astillero, adarga antigua, rocín flaco y galgo corredor. Una olla de más vaca que carnero, salpicón todas las noches, (...), consumían las tres partes de su hacienda.”

Estructuras de datos

FICHERO/DOCUMENTO

↳ PÁRRAFO (separación por \n)

↳ FRASE (separación por .)

“Primera parte del Ingenioso hidalgo Don Quijote de la Mancha.”

“Capítulo Primero.”

“Que trata de la condición y ejercicio del famoso hidalgo don Quijote de la Mancha.”

“En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no hace mucho que vivía un hidalgo de los de lanza en astillero, adarga antigua, rocín flaco y galgo corredor. ”

Estructuras de datos

FICHERO/DOCUMENTO

↳ PÁRRAFO (separación por \n)

↳ FRASE (separación por .)

↳ PALABRA (separación por ' ')

“En”, “un”, “lugar”, “de”, “la”, “Mancha”, “de”, “cuyo”...

Estructuras de datos

FICHERO/DOCUMENTO

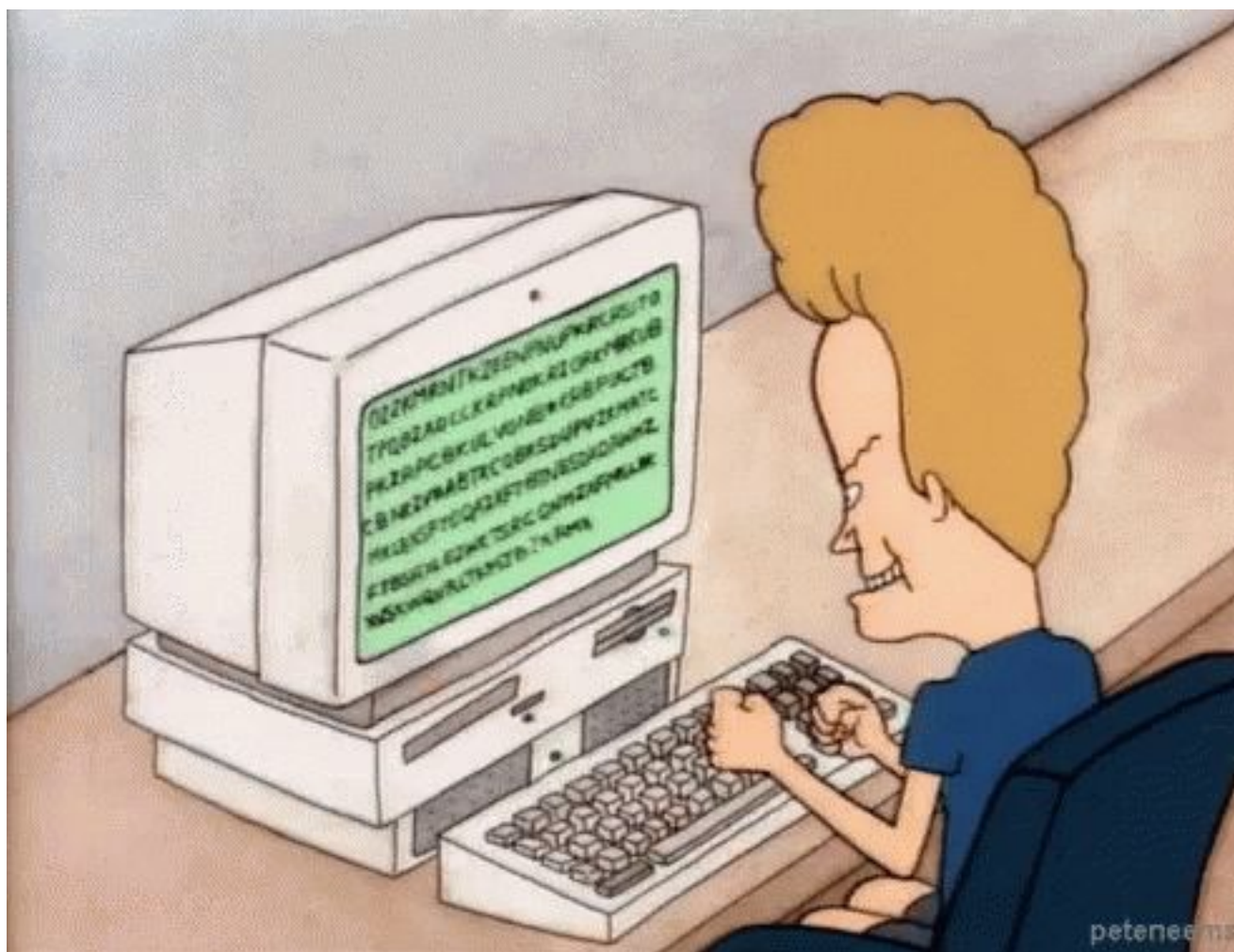
└→ PARRAFO (separación por \n)

└→ FRASE (separación por .)

└→ PALABRA (separación por ' ')

└→ CARÁCTER/LETRA

"E", "n", " ", "u", "n", " ", "l", "u", "g", "a", "r"...



Expresiones regulares

Las expresiones regulares, también conocidas como regex o regexp, son una herramienta común a todos los lenguajes de programación y sirven para realizar búsquedas no específicas pero regulares dentro de textos.

Son una de las herramientas más potentes y más utilizadas a la hora de extraer información, realizar comprobaciones o desarrollar productos de datos.

Expresiones regulares

Complementan las funciones vistas en el apartado anterior multiplicando su capacidad. Por ejemplo, podemos utilizar regex para buscar todas las palabras que empiezan por mayúscula sin tener que especificar todas las letras, buscar fechas dentro de textos o abreviaturas.

Expresiones regulares (*regex*) (I)

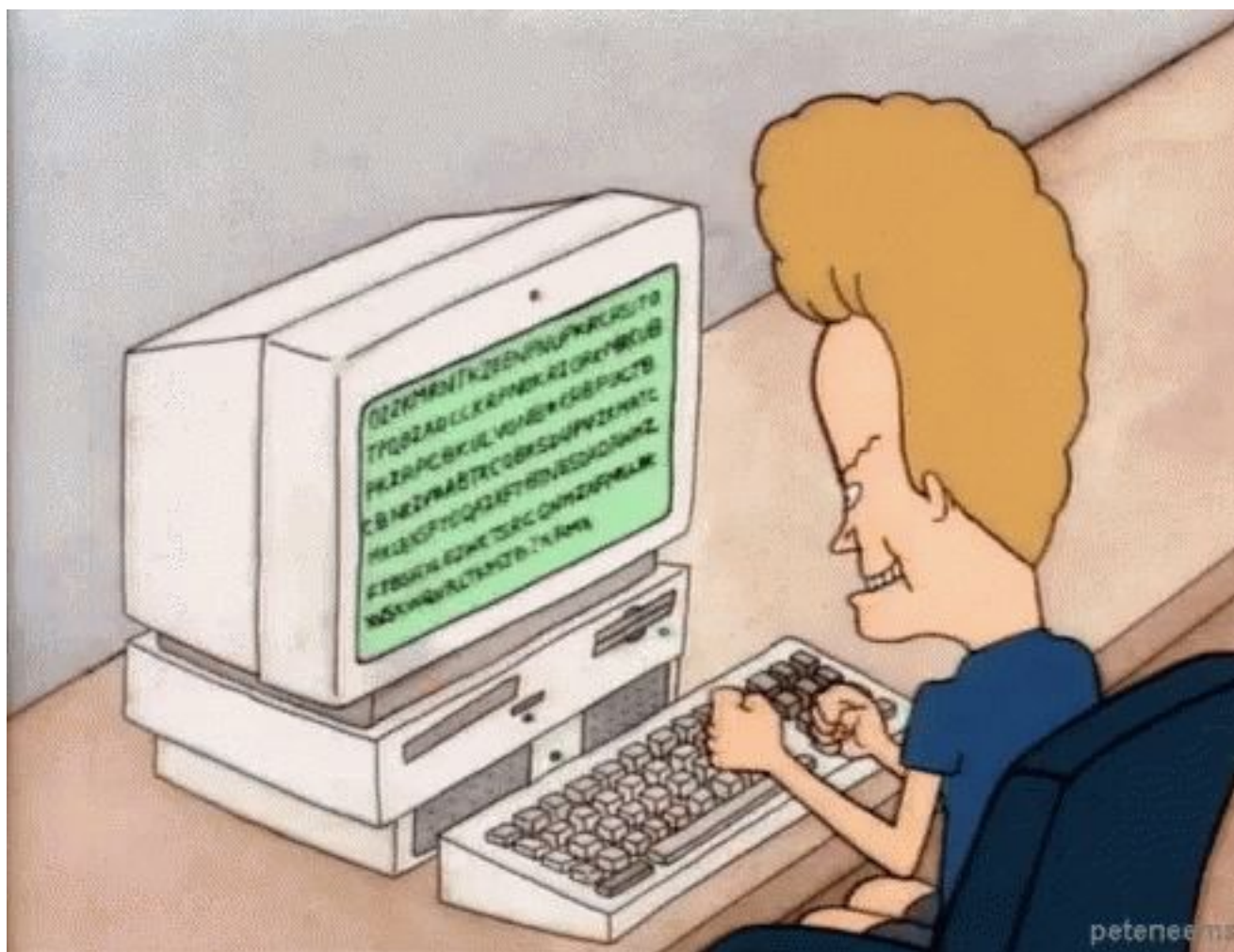
.	Un carácter cualquiera.
^	Inicio de string.
\$	Final de string.
[...]	Caracteres incluidos dentro.
a-z	Letras minúsculas.
A-Z	Letras mayúsculas.
0-9	Dígitos.
[^abc]	Cualquier carácter menos los que van después.
a b	O un carácter u otro.

Expresiones regulares (*regex*) (II)

<code>\</code>	Escape de caracteres.
<code>\b</code>	Limite de palabra.
<code>\d</code>	Dígito.
<code>\D</code>	Cualquier carácter menos un dígito.
<code>\s</code>	Cualquier blanco (‘ ‘, tabulador, salto de linea)
<code>\S</code>	Cualquier carácter menos los blancos.
<code>\w</code>	Cualquier valor alfanumérico.
<code>\W</code>	Cualquier carácter menos un alfanumérico.

Expresiones regulares (*regex*) (III)

*	0 o más veces.
+	1 o más veces.
?	1 vez o ninguna.
{n}	n veces.
{n,}	Al menos n veces.
{,n}	Como mucho n veces.
{n,m}	Al menos n veces y como mucho m veces.



petenee ms

2. LENGUA, 3º ESO



NLTK

- *Natural Language Tool Kit*
- Librería que incluye multitud de funciones para llevar a cabo operaciones sobre texto.
- Funciones **morfológicas** *e.g. Stemming/Lemmatization*
- Funciones **semánticas** *e.g. Sinónimos/Antónimos*
- Funciones **sintácticas** *e.g. Part Of Speech (POS)*
- Modelos pre-entrenados en inglés, capacidad de re-entrenar con corpus anotados de otros idiomas.

SpaCy

- Librería que incluye multitud de funciones para llevar a cabo operaciones sobre texto.
- Gran escalabilidad.
- Soporte directo a múltiples idiomas.
- Funciones **morfológicas** *e.g. Stemming/Lemmatization*
- Funciones **semánticas** *e.g. Sinónimos/Antónimos*
- Funciones **sintácticas** *e.g. Part Of Speech (POS)*

Tokenización

- Primer paso para poder aplicar muchas funciones.
- Dividir texto en estructuras más simples.
- *Sentence tokenizer*: división en frases. Función por defecto y posibilidad de adaptación a las estructuras de distintos idiomas.
- *Word tokenizer*: división en palabras.
- Se mantienen los signos de puntuación en los tokens.
- [*SpaCy*] Mantiene los saltos de línea también.

Part Of Speech

- Ambas librerías proporcionan una función que, dada una frase, asigna a cada una de las palabras una etiqueta de la función que realiza en la frase. Además incluyen visualizaciones de los resultados.
- [NLTK] Por defecto, tiene etiquetador entrenado en inglés. Existe la opción de entrenar el etiquetador con distintos corpus anotados para adaptar la herramienta a distintas sintaxis. Distintos idiomas, lenguajes de programación...
- [SpaCy] Etiquetador ya entrenado en castellano.

Abbreviation	Meaning
CC	coordinating conjunction
CD	cardinal digit
DT	determiner
EX	existential there
FW	foreign word
IN	preposition/subordinating conjunction
JJ	adjective (large)
JJR	adjective, comparative (larger)
JJS	adjective, superlative (largest)
LS	list marker
MD	modal (could, will)
NN	noun, singular (cat, tree)
NNS	noun plural (desks)
NNP	proper noun, singular (sarah)
NNPS	proper noun, plural (indians or americans)
PDT	predeterminer (all, both, half)

Press un

Stemming (sólo en NLTK)

- El proceso de stemming consiste en reducir las palabras a su raíz.
- El resultado no siempre (de hecho, pocas veces) es una palabra válida.
- Para llevar a cabo un proceso de stemming no se tiene en cuenta el contexto. No se tiene en cuenta el papel que juega la palabra en una frase.
- Existen varios algoritmos de stemming. El más popular es el de Porter.



The Porter stemming algorithm

Links to resources

- [The stemmer in Snowball](#)
- [Sample English vocabulary](#)
- [Its stemmed equivalent](#)
- [The 'official' home page of the Porter stemming algorithm](#)

Here is a case study on how to code up a stemming algorithm in Snowball. First, the definition of the Porter stemmer, as it appeared in *Program*, Vol 14 no. 3 pp 130-137, July 1980.

THE ALGORITHM

A *consonant* in a word is a letter other than A, E, I, O or U, and other than Y preceded by a consonant. (The fact that the term 'consonant' is defined to some extent in terms of itself does not make it ambiguous.) So in TOY the consonants are T and Y, and in SYZYGY they are S, Z and G. If a letter is not a consonant it is a *vowel*.

A consonant will be denoted by c, a vowel by v. A list ccc... of length greater than 0 will be denoted by C, and a list vvv... of length greater than 0 will be denoted by V. Any word, or part of a word, therefore has one of the four forms:

CVCV ... C

CVCV ... V

VCVC ... C

VCVC ... V

These may all be represented by the single form

[C]VCVC ... [V]

[C]VCVC ... [V]

<https://snowballstem.org/algorithms/porter/stemmer.html>

Lemmatizacion

- El objetivo es el mismo que el de aplicar *stemming*: encontrar la raíz de la palabra.
- Sin embargo, en lemmatizacion el resultado siempre es una palabra existente.
- Para llevar a cabo un proceso de lemmatizacion es útil tener en cuenta su contexto en la frase. Una misma palabra puede devolver distintos resultados dependiendo de esto.
- [SpaCy] Comete algunos errores por no tener en cuenta el contexto. [NLTK] Puede incluirse y lo utiliza.

Named Entities Recognition (NER)

- *NER* es el proceso de encontrar dentro del texto palabras o grupos de palabras que identifican a personas, lugares, empresas, etc.
- Más adelante se ha extendido a números, medidas, fechas y otro tipo de expresiones incluyendo cifras.
- Este tipo de referencias suelen ser muy informativas a la hora de analizar el contenido de un texto. Son de gran ayuda para clasificar textos.

Stopwords

- El término *stopwords* hace referencia a aquellas palabras que sirven de unión dentro de las frases pero que no tienen significado propio.
- *E.g.* El, a, aquellas, que, de, las, que.
- Es importante manejarlas correctamente a la hora de trabajar con los textos ya que pueden tener una gran influencia por su elevada frecuencia.
- Existen conjuntos predefinidos de *stopwords* en distintos idiomas pero a veces es necesario adaptarlos al problema particular.

Text @NLTK

- NLTK define una estructura llamada *Text* que implementa varias funciones para realizar búsquedas dentro del texto.
- Son funciones muy útiles para manejar textos grandes. Pueden servir incluso para analizar el estilo de la persona que ha escrito el texto.
- NLTK incluye varios ejemplos de textos dentro de la librería `books`.
- Funciones : `concordance`, `similar`, `common_contexts`, `dispersion_plot`, `collocation`.

Extra info @spaCy

- En el momento en que se lee un texto con SpaCy automáticamente se procesa y se llevan a cabo todos los análisis.
- Cada análisis resulta en un etiquetado de los tokens que puede consultarse.
- Un ejemplo sería el análisis de los conjuntos de nombres (noun_chunks) que analiza las relaciones de los nombre con el resto de elementos de la frase.
- Mantiene el soporte a múltiples idiomas.

Wordnet @NLTK

- Word Net es una librería de funcionalidades incluida dentro de NLTK.
- A partir de un corpus, Word Net proporciona funcionalidades relacionadas con la semántica de las palabras.
- Definiciones, sinónimos, antónimos, frases de ejemplo...
- Da soporte a múltiples idiomas pero siempre con el inglés como idioma central.

3. DISTANCIAS

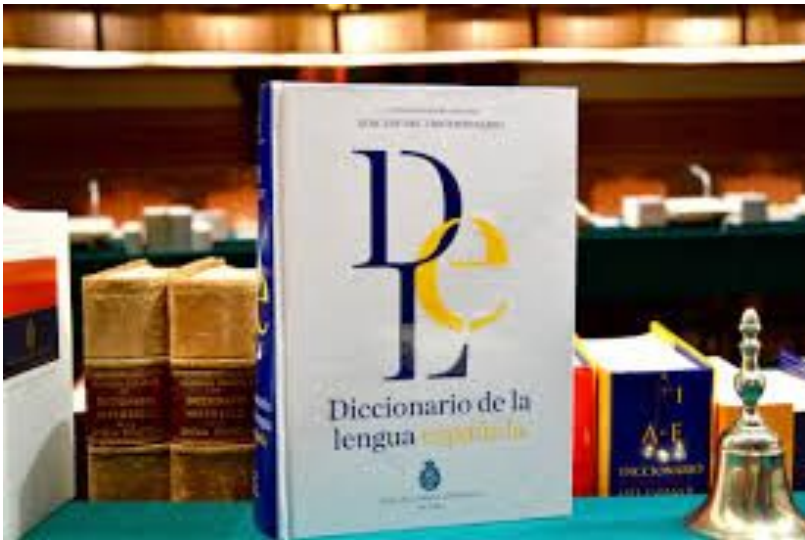




¿Cómo haríais un corrector de escritura en móvil en media hora?



¿Cómo haríais un corrector de escritura en móvil en media hora?



Función de
+ distancia entre
palabras

Distancia de Hamming

- Nº de posiciones en las que tenemos el mismo símbolo en ambas cadenas.
- Esta distancia solo sirve para cadenas de la misma longitud. Cuanto más alta más cercanas.
- $\text{Hamming}(\text{loca}, \text{lolo}) = 2$
- $\text{Hamming}(\text{loca}, \text{loac}) = 2$
- $\text{Hamming}(\text{loca}, \text{lonas}) = \text{NO APLICA}$

Distancia de Levenshtein

- Nº mínimo de inserciones, borrados o cambios de caracteres para convertir una cadena en otra.
Cuanto más baja, más cercanas son las palabras.
- $\text{Lev}(\text{loca}, \text{lolo}) = 2$
- $\text{Lev}(\text{loca}, \text{loac}) = 2$
- $\text{Lev}(\text{loca}, \text{lonas}) = 2$

Distancia de Damerau-Levenshtein

- Igual que la de Levenshtein pero una trasposición de símbolos adyacentes sólo cuenta como 1.
- $\text{Lev}(\text{loca}, \text{lolo}) = 2$
- $\text{Lev}(\text{loca}, \text{loac}) = 1$
- $\text{Lev}(\text{loca}, \text{lonas}) = 2$

Distancia de q-gramas

- Número total de q-gramas que aparecen en alguna cadena pero no en la otra.
- $2\text{-gramas}(\text{loca}, \text{lolo}) = 3$ (lo, oc, ca, lo, ol, lo)
- $2\text{-gramas}(\text{loca}, \text{looc}) = 2$ (lo, oc, ca, lo, oo, oc)
- $2\text{-gramas}(\text{loca}, \text{lonas}) = 5$ (lo, oc, ca, lo, on, na, as)

Otras distancias

- Distancia del coseno
- Distancia de Jaccard
- Distancia de Jaro-Winkler
- ...



4. TODO FUNCIONA MEJOR CON NÚMEROS

¿ $\text{dist}(\text{malo}, \text{bueno}) > \text{dist}(\text{malo}, \text{mago})$?



¿ $\text{dist}(\text{malo}, \text{bueno}) > \text{dist}(\text{malo}, \text{mago})$?

Las distancias del capítulo anterior sirven para corregir typos, intentar predecir qué palabra está escribiendo el usuario una vez ha empezado a escribirla o trabajar con entidades propias.

Pero, ¿sirven para otro tipo de problemas como la clasificación o la generación de textos?

¿En qué casos creéis que la respuesta a la pregunta de las distancias es sí y en cuales es no?

De textos a vectores

Todos los modelos de Machine Learning están basados en el manejo de vectores y matrices de números.

¿Qué hacemos cuando queremos trabajar con textos? Convertir los textos a vectores o matrices que los representen: dos textos que tienen cosas en común deben generar vectores que estén cercanos.

Vectorización aditiva

La mayoría de algoritmos de vectorización de textos son aditivos. Es decir, el vector que representa a un texto se genera a partir de la suma de los vectores que representan a sus elementos: palabras, n-gramas, parejas de palabras...

me = (rojo, azul, blanco);

gusta= (blanco, blanco, blanco);

mucho=(blanco, amarillo, blanco);



“me gusta mucho” =
(rojo, verde, blanco)

Preprocesado

- Limpieza de texto
- Tokenización
- Normalización
- *Lemmatización/Stemming*

Bag of words

- Es la manera más simple de vectorizar textos.
- Se generan vectores del tamaño del vocabulario del dataset.
- Cada palabra del vocabulario tiene asignada una posición y el vector correspondiente a esa palabra tiene un 1 en esa posición y 0s en el resto.
- A cada texto del dataset se le asigna el vector aditivo de los elementos que lo forman.

Bag of words

Dataset

T1: "Me ha gustado la película"

T2: "No me ha gustado"

T3: "La película es mala"

Bag of words

Dataset

T1: "Me ha gustado la película"

T2: "No me ha gustado"

T3: "La película es mala"

Vocabulario

me = 1; haber = 2; gustar= 3; la=4; película=5; no=6; ser=7; mala=8

Bag of words

Dataset

T1: "Me ha gustado la película"

T2: "No me ha gustado"

T3: "La película es mala"

Vocabulario

me = 1; haber = 2; gustar = 3; la = 4; película = 5; no = 6; ser = 7; mala = 8

Vectorización

T1 = [1,1,1,1,1,0,0,0]

T2 = [1,1,1,0,0,1,0,0]

T3 = [0,0,0,1,1,0,1,1]

TF-IDF

- *Term Frequency-Inverse Document Frequency.*
- Se generan vectores del tamaño del vocabulario del dataset.
- Al contrario que bag of words, TF-IDF le da distintos valores a cada palabra según su aporte de información para hacer un texto único. El valor depende de la frecuencia con la que aparece la palabra en el dataset.
- $TF-IDF = TF * IDF.$

TF

- $TF = N^{\circ}$ de veces que la palabra aparece en el texto.
- Aunque esta es la definición básica existen otras opciones para calcular el término TF.
 - Booleano: 1 si la palabra aparece, 0 si no.
 - Frecuencia normalizada: se divide por el tamaño del texto o por la frecuencia de la palabra que más aparece en el texto.
 - Escalado logarítmico: $\log(1 + TF)$

IDF

- $IDF :=$ cuánto ayuda esa palabra a diferenciar textos = $\log(\text{Nº de textos totales} / \text{Nº de textos que incluyen esa palabra})$
- A veces se modifica la fórmula básica para protegerse ante la aparición de nuevas palabras que podrían generar divisiones por 0.
 - Ajuste del denominador: $\log(\text{Nº de textos} / (1 + \text{Nº de textos que incluyen}))$
- Palabras raras \Rightarrow IDF alto,
- Palabras comunes \Rightarrow IDF ≈ 0

TF-IDF: ejemplo

TEXTO 1

Término	Frecuencia
gustar	7
película	8
buena	3

TEXTO 2

Término	Frecuencia
gustar	8
película	11
mala	3
no	8

TEXTO 3

Término	Frecuencia
película	4
Tarantino	1
Oscar	1

TF-IDF: ejemplo

TEXTO 1

Término	Frecuencia	Frec. normalizada
gustar	7	7/8
película	8	8/8
buena	3	3/8

TEXTO 2

Término	Frecuencia	Frec. normalizada
gustar	8	8/11
película	11	11/11
mala	3	3/11
no	8	8/11

TEXTO 3

Término	Frecuencia	Frec. normalizada
película	4	4/4
Tarantino	1	1/4
Oscar	1	1/4

TF-IDF: ejemplo

TEXTO 1

Término	Frecuencia	Frec. normalizada	IDF
gustar	7	7/8	$\log(3/2)$
película	8	8/8	$\log(1) = 0$
buena	3	3/8	$\log(3)$

TEXTO 2

Término	Frecuencia	Frec. normalizada	IDF
gustar	8	8/11	$\log(3/2)$
película	11	11/11	$\log(1)=0$
mala	3	3/11	$\log(3)$
no	8	8/11	$\log(3)$

TEXTO 3

Término	Frecuencia	Frec. normalizada	IDF
película	4	4/4	$\log(1)=0$
Tarantino	1	1/4	$\log(3)$
Oscar	1	1/4	$\log(3)$

TF-IDF: ejemplo

TEXTO 1

Término	Frecuencia	Frec. normalizada	IDF	TF-IDF
gustar	7	7/8	$\log(3/2)$	0.15
película	8	8/8	$\log(1) = 0$	0
buena	3	3/8	$\log(3)$	0.17

TEXTO 2

Término	Frecuencia	Frec. normalizada	IDF	TF-IDF
gustar	8	8/11	$\log(3/2)$	0.13
película	11	11/11	$\log(1)=0$	0
mala	3	3/11	$\log(3)$	0.13
no	8	8/11	$\log(3)$	0.34

TEXTO 3

Término	Frecuencia	Frec. normalizada	IDF	TF-IDF
película	4	4/4	$\log(1)=0$	0
Tarantino	1	1/4	$\log(3)$	0.12
Oscar	1	1/4	$\log(3)$	0.12

TF-IDF

Dataset

T1, T2, T3

Vocabulario

gustar = 1; película = 2; buena = 3; mala = 4; no = 5; Tarantino = 6;
Óscar = 7

Vectorización

T1 = [0.15,0,0.17,0,0,0,0]

T2 = [0.13,0,0,0.13,0.34,0,0]

T3 = [0,0,0,0,0,0.12, 0.12]

TF-IDF: ¡ojo con el split de los datos!

Uno de los errores más comunes que se cometen al aplicar TF-IDF es mezclar los conjuntos de entrenamiento y test.

Al igual que ocurre con el resto de datos, la información necesaria para aplicar TF-IDF debe provenir únicamente de los datos de entrenamiento.

El proceso a seguir es el siguiente: se calcula el término IDF para todas las palabras del conjunto de entrenamiento y se guardan para utilizarlas tanto en entrenamiento como en la fase de inferencia. Nunca se recalculan con el conjunto de test.

Word embeddings

- Tanto *bag of words* como *TF-IDF* tienen en cuenta las palabras incluidas en el vocabulario pero no tienen en cuenta el contexto en el que se utilizan.
- Existen una serie de algoritmos de vectorización conocidos como *word embeddings* que utilizando la potencia de las redes neuronales generan representaciones vectoriales de las palabras analizando el contexto en el que se utilizan.

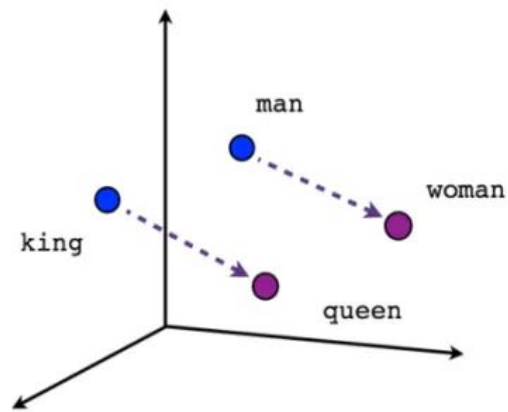
Word embeddings

- Los *word embeddings* aseguran que las representaciones vectoriales de dos palabras que se utilicen en contextos parecidos estarán cercanas en el espacio vectorial definido.
- Para poder entrenar estos algoritmos hacen falta datasets muy grandes que permitan generalizar y sacar conclusiones sobre los contextos de cada palabra.
- Pueden ser entrenados con datasets de distintos ámbitos.

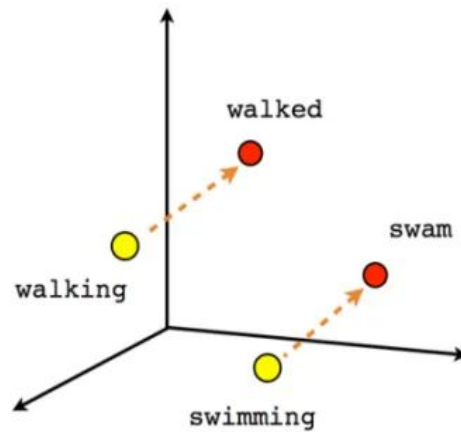
word2vec

- Creado en 2013 por un grupo de investigadores de Google.
- Se basa en una red neuronal de 2 capas y está entrenado con datos de Google.
- Existen dos modos de representar el problema que se traducen en dos arquitecturas de red distintas: se conocen con *Skip-Gram* o *Continuous Bag of Words* (CBOW) .
- En ambos casos el resultado es la asignación de una representación vectorial a cada palabra del vocabulario.

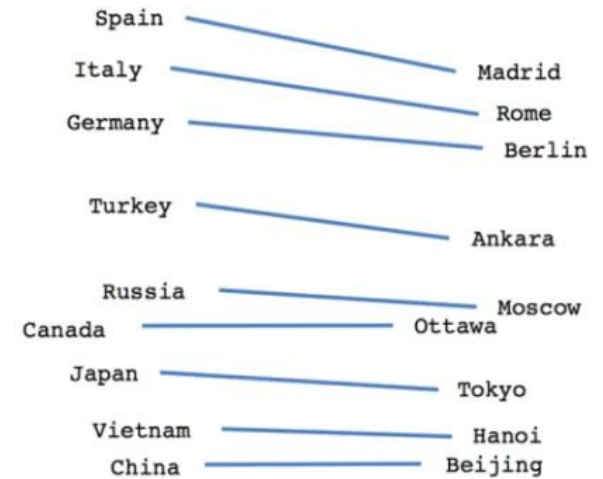
word2vec



Male-Female



Verb tense



Country-Capital

Otros modelos de word embeddings

- Stanford GloVe:
<https://nlp.stanford.edu/projects/glove>
- Facebook fastText:
<https://fasttext.cc>
- Google BERT:
<https://mccormickml.com/2019/05/14/BERT-word-embeddings-tutorial/#what-is-bert>

5. CASO DE USO: ANÁLISIS DE SENTIMIENTO



