



Universidad  
**Tecnológico**®





# Procesamiento de lenguaje natural y visión computacional

Modelos de lenguaje en  
n-gramas



Las convenciones utilizadas en la escritura a veces no reflejan la organización lógica de los símbolos textuales, por ejemplo, algunos signos de puntuación se escriben antes de cierta palabra y otros después, aunque no formen parte de ella y esto depende del lenguaje y del estilo de escritura. En ocasiones, en una cadena de caracteres o texto no están claros los límites lógicos entre las palabras o símbolos de puntuación y eso puede complicar el análisis de textos escritos.

En este tema conocerás lo siguiente:

- La importancia de la estandarización de textos.
- El concepto de token y tokenización.
- Qué son los unigramas, bigramas y n-gramas.
- La diferencia entre un q-grama y un n-grama.
- La utilidad de los modelos de lenguaje.
- La diferencia entre un conjunto de prueba y uno de entrenamiento.



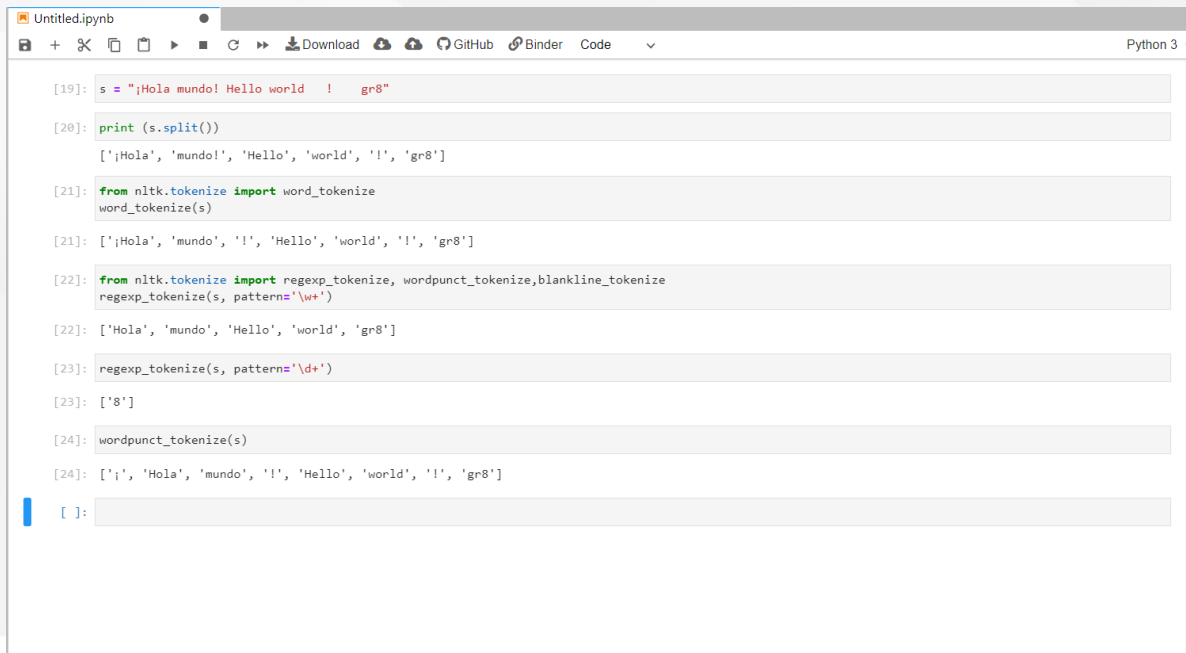


Es posible someter a los textos a un **proceso de estandarización**, el cual puede ayudar a mejorar la obtención de las entidades individuales (las palabras) y con ello, realizar el análisis del mensaje dentro del texto.

La **palabra** o **token** es la unidad mínima utilizada en procesamiento de textos y definirla depende de la tarea que se esté realizando. Por lo general, cada palabra se define como una unidad separada por espacios o signos de puntuación, sin embargo, los signos de puntuación pueden estimarse o no como palabras.

Los **tokens** son, por tanto, la ocurrencia de las palabras en un texto y la **tokenización** es el proceso que permite encontrar todas las palabras dentro de él. El conjunto de todas las palabras posibles dentro de un texto se conoce como **vocabulario**.

En la figura 1 se muestra la tokenización del texto: ¡Hola mundo! Hello world ! gr8", utilizando diferentes instrucciones.



```
Untitled.ipynb Python 3
[19]: s = "¡Hola mundo! Hello world ! gr8"
[20]: print (s.split())
      ['¡Hola', 'mundo!', 'Hello', 'world', '!', 'gr8']
[21]: from nltk.tokenize import word_tokenize
      word_tokenize(s)
[21]: ['¡Hola', 'mundo', '!', 'Hello', 'world', '!', 'gr8']
[22]: from nltk.tokenize import regexp_tokenize, wordpunct_tokenize, blankline_tokenize
      regexp_tokenize(s, pattern='\w+')
[22]: ['Hola', 'mundo', 'Hello', 'world', 'gr8']
[23]: regexp_tokenize(s, pattern='\d+')
[23]: ['8']
[24]: wordpunct_tokenize(s)
[24]: ['!', 'Hola', 'mundo', '!', 'Hello', 'world', '!', 'gr8']
[ ]:
```





La eliminación de las palabras frecuentes puede ayudar a construir un vocabulario que contenga solo aquellas palabras que le den significado o valor al texto. En la figura se observa un ejemplo de eliminación de palabras frecuentes en un texto, se procesa independientemente si está escrito en minúsculas o mayúsculas. Nota que las palabras “se, con, que, antes, de, tu, a, tus, esto, es, le” no forman parte del vocabulario final del ejemplo.

```
Untitled.ipynb Python 3 O
+ ✂ 📄 📄 ▶ ⏪ ⏩ ⏴ ⏵ Download 📁 📁 GitHub 🔗 Binder Code ▾

[15]: import re
      from nltk.corpus import stopwords

      texto = 'Se dice con frecuencia que antes de morir tu vida pasa frente a tus ojos. De hecho esto es cierto. Se le llama vivir.'

      #convertir la oración a minúsculas
      texto = texto.lower()

      #Elimina signos de puntuación y tokeniza
      palabras = re.findall(r'\w+', texto, flags = re.UNICODE)

      #Elimina palabras frecuentes
      palabras_importantes=[]
      for palabra in palabras:
          if palabra not in stopwords.words('spanish'):
              palabras_importantes.append(palabra)

      print (palabras_importantes)

['dice', 'frecuencia', 'morir', 'vida', 'pasa', 'frente', 'ojos', 'hecho', 'cierto', 'llama', 'vivir']

[ ]:
```





En ocasiones se requiere extraer las raíces de las palabras y construir el vocabulario con ellas, así es posible eliminar sufijos como “ando” o “iendo”. El **stemming** consiste en eliminar y reemplazar los sufijos de la raíz de la palabra. Este proceso tiene sus reglas dependiendo del idioma, en español se eliminan sufijos de acciones (ar, er, ir, ía, es, etc.), terminaciones plurales, géneros, entre otros. El proceso de extracción regresa una palabra sin la terminación o sufijos.

```
Untitled.ipynb Python 3
```

```
[22]: from nltk.tokenize import WordPunctTokenizer
texto = "No sabía como se ponía una lavadora hasta que conocí esta y es que es muy sencilla de utilizar! Todo un gustazo cuando estamos aprendiendo."
texto_tokenizado = WordPunctTokenizer().tokenize(texto)
print(texto_tokenizado)

['No', 'sabía', 'como', 'se', 'ponía', 'una', 'lavadora', 'hasta', 'que', 'conocí', 'esta', 'y', 'es', 'que', 'es', 'muy', 'sencilla', 'de', 'utiliza', 'r', '!', 'Todo', 'un', 'gustazo', 'cuando', 'estamos', 'aprendiendo', '.']
```

```
[23]: from nltk.stem import SnowballStemmer
stemmer = SnowballStemmer('spanish')
stemmed_text = [stemmer.stem(i) for i in texto_tokenizado]
print(stemmed_text)

['no', 'sab', 'com', 'se', 'pon', 'una', 'lavador', 'hast', 'que', 'conoc', 'esta', 'y', 'es', 'que', 'es', 'muy', 'sencill', 'de', 'utiliz', '!', 'to', 'd', 'un', 'gustaz', 'cuand', 'estam', 'aprend', '.']
```

```
[ ]:
```

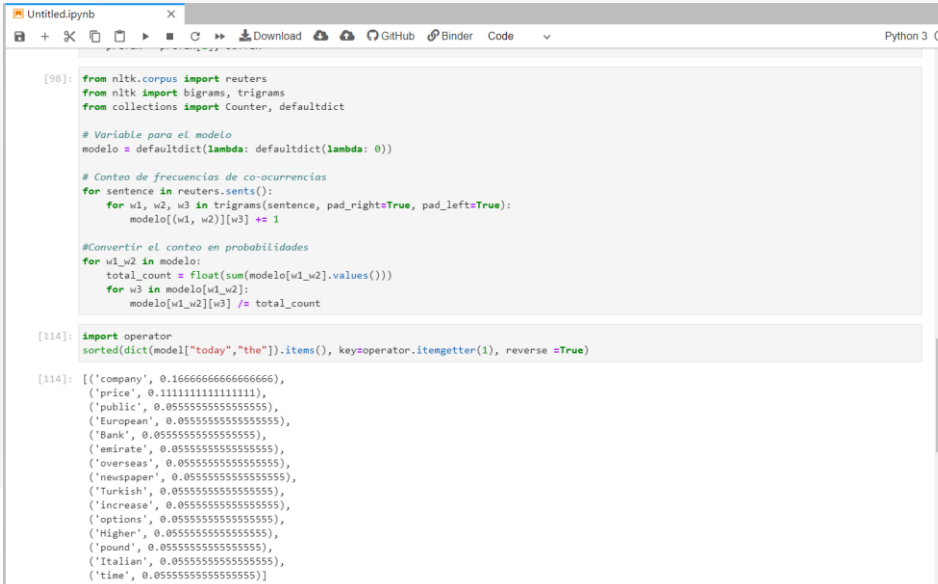




En general, el **n-grama** se puede definir como la secuencia de N elementos de una muestra determinada de texto o voz. Los elementos pueden ser letras o palabras. Algunos autores consideran que un **n-grama** es una secuencia de N palabras y un **q-grama** una secuencia de Q caracteres (incluidas las letras). Si el n-grama está formado solo por una palabra se llama **unigrama**, si está compuesto por dos, **bigrama**, y si está formado por tres elementos, **trigrama**.

Un **modelo de lenguaje** debe ser capaz de capturar la mayor cantidad de información posible para describir al lenguaje en términos de esa información. Es posible obtener un modelo estadístico de lenguaje basado en n-gramas. El modelo debe ser capaz de predecir la siguiente palabra en una oración, es decir,  $P(w|h)$ , la probabilidad de que una palabra  $w$  suceda dada una historia  $h$ .

En la figura se ve una implementación del modelo utilizando el corpus Reuters disponible en NLTK. El ejemplo utiliza trigramas para construir el modelo. Como resultado, se determina que la palabra más probable después del texto "today the" es "company" seguida de "price".



```
Untitled.ipynb
Python 3

[98]: from nltk.corpus import reuters
      from nltk import bigrams, trigrams
      from collections import Counter, defaultdict

      # Variable para el modelo
      modelo = defaultdict(lambda: defaultdict(lambda: 0))

      # Conteo de frecuencias de co-ocurrencias
      for sentence in reuters.sents():
          for w1, w2, w3 in trigrams(sentence, pad_right=True, pad_left=True):
              modelo[(w1, w2)][w3] += 1

      #Convertir el conteo en probabilidades
      for w1_w2 in modelo:
          total_count = float(sum(modelo[w1_w2].values()))
          for w3 in modelo[w1_w2]:
              modelo[w1_w2][w3] /= total_count

[114]: import operator
      sorted(dict(model["today", "the"]).items(), key=operator.itemgetter(1), reverse =True)

[114]: [('company', 0.16666666666666666),
      ('price', 0.11111111111111111),
      ('public', 0.05555555555555555),
      ('European', 0.05555555555555555),
      ('Bank', 0.05555555555555555),
      ('emirate', 0.05555555555555555),
      ('overseas', 0.05555555555555555),
      ('newspaper', 0.05555555555555555),
      ('Turkish', 0.05555555555555555),
      ('increase', 0.05555555555555555),
      ('options', 0.05555555555555555),
      ('Higher', 0.05555555555555555),
      ('pound', 0.05555555555555555),
      ('Italian', 0.05555555555555555),
      ('time', 0.05555555555555555)]
```



01

¿Cuál es la importancia del uso de n-gramas en el PLN?

02

¿El proceso de tokenización y n-gramas es un proceso que realizamos mentalmente los seres humanos durante una conversación casual?







En este tema se presentó el modelado del lenguaje y los n-gramas, que son dos de las herramientas más utilizadas en el procesamiento de lenguaje natural.

Los modelos de lenguaje asignan un valor de probabilidad a una oración o secuencia de palabras y con ello es posible predecir la palabra que sigue en un texto.

SHION T W NEWS EVER EXPERIENC  
YORK WHITE TOGETHER SURF WHAT  
GROWTH ABOUT JUST CONSULTANT  
T ALL SECOND NEXT LOVE QUICK I  
ERO FOCUS HOW FOUND THE LOCK  
BEHIND MORE V THEY TIME OPEN N  
N THAN MEGA GREAT WITH WELL  
AIL SEE MAKE FUN EVEN WAY SEE  
P AND NEW FULL BEAUTIFULL RANGE  
AT BETTER COLORS INVESTMENT N  
RN CATCH D JOK SHAPE BEACH A  
PE LITTLE UNDER OR KICK WISH A  
SMART HEN SAME POSSIBLE FACE  
SUN GIVE LIKE HAPPENS YOU BE  
SUN GIVE LIKE HAPPENS YOU BE





# Procesamiento de lenguaje natural y visión computacional

Clasificación de texto



La clasificación de textos es una de las tareas más importantes y típicas del aprendizaje automático supervisado. La asignación de categorías a documentos que pueden ser una página web, un libro de biblioteca, artículos multimedia, una galería, etc., tiene muchas aplicaciones, por ejemplo, filtrado de spam, enrutamiento de correo electrónico, análisis de opiniones o sentimientos, entre otros.

En este tema conocerás:

El proceso de clasificación de textos.

Los algoritmos más comunes utilizados para clasificar textos.

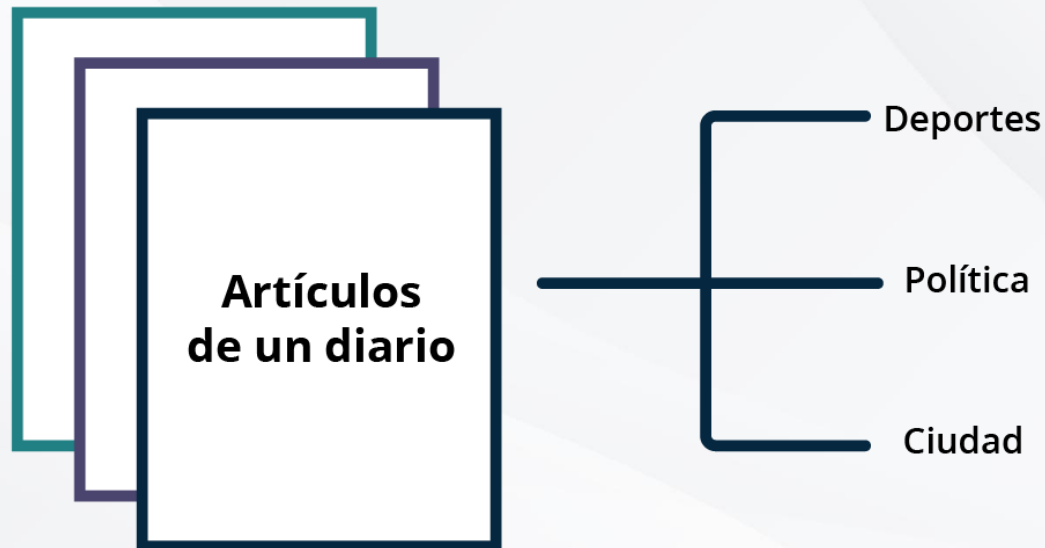
El clasificador bayesiano.

Los pasos necesarios para entrenar a un clasificador bayesiano.





Existen distintos tipos de algoritmos de aprendizaje automático: supervisado, no supervisado, por refuerzo y semisupervisado. La clasificación de textos es un ejemplo de tarea de aprendizaje automático supervisado, ya que un conjunto de datos etiquetado y formado por documentos de texto y sus etiquetas se utiliza para entrenar a un clasificador.





Un clasificador ingenuo bayesiano (*Naive Bayes*, por su nombre en inglés) o simplemente bayesiano está basado en el teorema de Bayes. Este clasificador asume que la presencia o ausencia de alguna característica particular no está relacionada con la presencia o ausencia de alguna otra. De tal manera que estas características contribuyen de forma independiente a la probabilidad de que tal objeto pertenezca a una clase independientemente de la presencia o ausencia de las otras características. Por ejemplo, una fruta puede clasificarse como un limón si es de color verde, redonda y aproximadamente de 4 cm de diámetro.

El teorema de Bayes permite calcular la probabilidad de que cierto objeto con ciertas características pertenezca a una clase dado cierto conocimiento *a priori*:

$$P(\text{clase} | \text{objeto}) = (P(\text{objeto} | \text{clase}) * P(\text{clase})) / P(\text{objeto})$$

Donde  $P(\text{clase} | \text{objeto})$  es la probabilidad de una clase dado el objeto con ciertas características.





El proceso de clasificación de textos incluye los siguientes pasos:

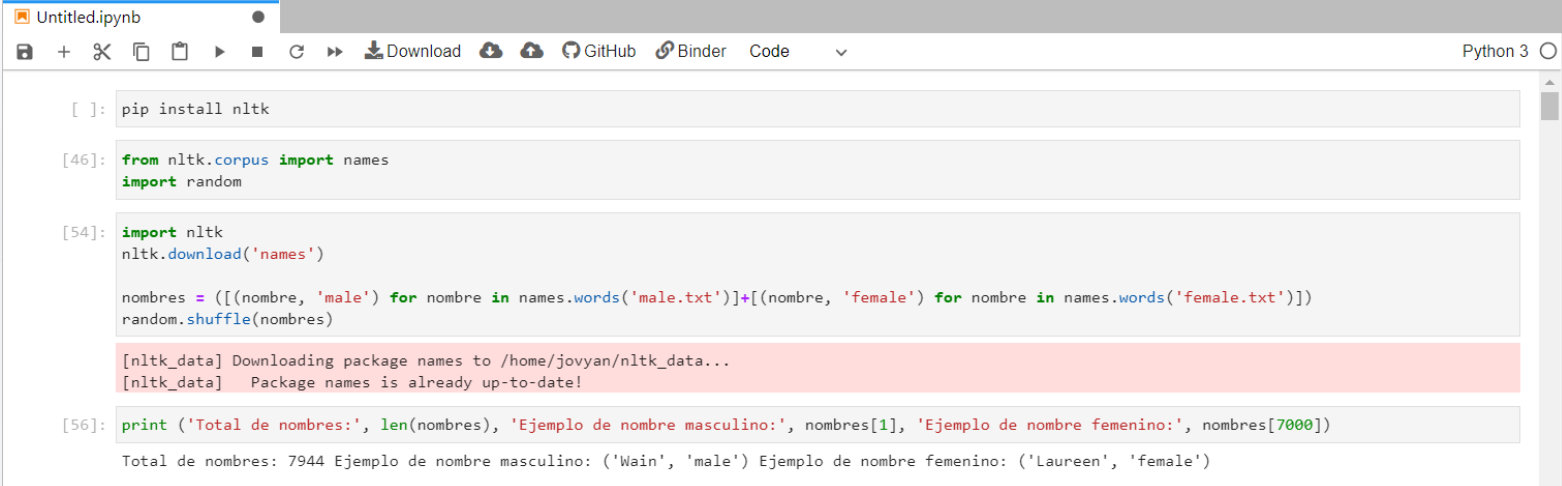
Preparación de los datos: en este paso se selecciona el conjunto de datos con el que se trabajará y se realiza un preprocesamiento básico al texto para después separarlo en el conjunto de entrenamiento y de prueba.

Extracción de características: el texto se transforma en elementos característicos que se usarán en el proceso de clasificación.

Entrenamiento del modelo: el clasificador se entrena con los datos etiquetados.

Evaluación del modelo: el clasificador se evalúa y, en caso de ser necesario, se realizan ajustes para mejorar el resultado.

El primer ejemplo trata de clasificar nombres de personas según el género. Se utilizará la librería NLTK (Bird, Klein y Loper, 2009) y el corpus *names*. En la figura se muestra la descarga de los datos y una breve exploración sobre ellos. En total hay 7,944 nombres clasificados como masculinos o femeninos.



```
Untitled.ipynb Python 3

[ ]: pip install nltk

[46]: from nltk.corpus import names
import random

[54]: import nltk
nltk.download('names')

nombres = [(nombre, 'male') for nombre in names.words('male.txt')] + [(nombre, 'female') for nombre in names.words('female.txt')]
random.shuffle(nombres)

[nltk_data] Downloading package names to /home/jovyan/nltk_data...
[nltk_data] Package names is already up-to-date!

[56]: print('Total de nombres:', len(nombres), 'Ejemplo de nombre masculino:', nombres[1], 'Ejemplo de nombre femenino:', nombres[7000])

Total de nombres: 7944 Ejemplo de nombre masculino: ('Wain', 'male') Ejemplo de nombre femenino: ('Laureen', 'female')
```





01

Descarga el corpus “*wine dataset*” dentro de scikit-learn y explóralo. El conjunto de datos contiene los resultados del análisis químico de vinos cultivados en la misma región en Italia, pero derivados de tres cultivares diferentes, consta de 13 atributos y 3 tipos de vino.

02

Separa las columnas de los datos en variables dependientes e independientes, es decir, atributos y etiquetas de clase.

Utiliza el código siguiente para separar los datos en el conjunto de entrenamiento y de prueba.

```
from sklearn.cross_validation import train_test_split
X_train, X_test, y_train, y_test = train_test_split(wine.data,
wine.target, test_size=0.3,random_state=109)
```

03

Genera el modelo utilizando el conjunto de datos de entrenamiento.

04

Con el modelo construido, realiza predicciones utilizando el conjunto de datos de prueba y almacena los resultados en una variable nombrada *y\_pred*.





Muchas de las tareas en procesamiento de lenguaje natural se pueden tratar como problemas de clasificación. La categorización de textos, donde el objetivo es asignar una clase a un texto particular, incluye tareas como análisis de sentimientos, detección de spam, identificación del lenguaje, identificación del sexo de un autor, entre otros. No olvides que los clasificadores supervisados utilizan corpus previamente etiquetados que les permiten construir modelos que predicen la categoría o etiqueta de una entrada en función de ciertas características. De igual forma, recuerda que para entrenar a un clasificador supervisado se debe separar el total de los datos en un conjunto de entrenamiento y otro de pruebas.







● Bird, S., Klein, E., y Loper, E. (2009). *Natural language processing with Python: analyzing text with the natural language toolkit*. Estados Unidos: O'Reilly Media.

