

Creado por:

Isabel Maniega

# Natural Language Processing (NLP)

<https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset>

El objetivo de este ejercicio:

- Los ordenadores trabajan con números, no con letras
- así que necesitamos NLP para tranasformar las palabras a números

```
In [1]: import warnings
warnings.filterwarnings("ignore")
```

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [4]: from sklearn.naive_bayes import MultinomialNB
```

## Cargar archivo .csv

```
In [5]: # https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset
```

```
In [6]: df = pd.read_csv("spam.csv",
                        sep=";", encoding='ISO-8859-1')
df.head(15)
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only ...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
5	spam	FreeMsg Hey there darling it's been 3 week's n...	NaN	NaN	NaN
6	ham	Even my brother is not like to speak with me. ...	NaN	NaN	NaN
7	ham	As per your request 'Melle Melle (Oru Minnamin...	NaN	NaN	NaN
8	spam	WINNER!! As a valued network customer you have...	NaN	NaN	NaN
9	spam	Had your mobile 11 months or more? U R entitle...	NaN	NaN	NaN
10	ham	I'm gonna be home soon and i don't want to tal...	NaN	NaN	NaN
11	spam	SIX chances to win CASH! From 100 to 20,000 po...	NaN	NaN	NaN
12	spam	URGENT!! You have won a 1 week FREE membership ...	NaN	NaN	NaN
13	ham	I've been searching for the right words to tha...	NaN	NaN	NaN
14	ham	I HAVE A DATE ON SUNDAY WITH WILL!!	NaN	NaN	NaN

```
In [7]: df = df.iloc[:, 0:2]
df.head()
```

	v1	v2
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

## Nombres para las columnas

```
In [8]: df.columns= ["Status", "Message"]
df.head()
```

	Status	Message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

```
In [10]: df.shape
```

```
Out[10]: (5572, 2)
```

```
In [11]: len(df)
```

```
Out[11]: 5572
```

## Vemos si nos faltan algunos datos

```
In [12]: df.Message.isnull().sum()
```

```
Out[12]: 0
```

```
In [13]: df.describe()
```

	Status	Message
count	5572	5572
unique	2	5169
top	ham	Sorry, I'll call later
freq	4825	30

## ¿Cuántos datos de "spam" en nuestros datos?

Forma 1

```
In [14]: df.head()
```

	Status	Message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

```
In [15]: df.Status.value_counts()
```

```
Out[15]: ham      4825
spam      747
Name: Status, dtype: int64
```

Forma 2

```
In [16]: df.iloc[:,0].value_counts()
```

```
Out[16]: ham      4825
spam      747
Name: Status, dtype: int64
```

Forma 3

```
In [17]: df_spam = df[df.Status == "spam"]
len(df_spam)
```

```
Out[17]: 747
```

Forma 4

```
In [18]: data = df[df.iloc[:,0] == "spam"]
len(data)
```

```
Out[18]: 747
```

## spam == 1 (True); ham == 0 (False)

Método 1

```
In [19]: df["Status"] = df["Status"].map({"ham": 0, "spam": 1})
df.head()
```

	Status	Message
0	0	Go until jurong point, crazy.. Available only ...
1	0	Ok lar... Joking wif u oni...
2	1	Free entry in 2 a wkly comp to win FA Cup fina...
3	0	U dun say so early hor... U c already then say...
4	0	Nah I don't think he goes to usf, he lives aro...

```
In [20]: df.shape
```

```
Out[20]: (5572, 2)
```

```
In [21]: X = df.Message
```

```
In [22]: y = df.Status
```

## Train, Test split

```
In [24]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0, test_size=0.2)
```

## Método 1: CountVectorizer

```
In [25]: from sklearn.feature_extraction.text import CountVectorizer
```

```
In [26]: cv = CountVectorizer()
```

```
In [28]: X_train = cv.fit_transform(X_train)
X_test = cv.transform(X_test)
```

```
In [29]: y_train = y_train.astype("int")
y_test = y_test.astype("int")
```

```
In [30]: y_train = np.array(y_train)
y_test = np.array(y_test)
```

```
In [31]: X_train
```

```
Out[31]: <4457x7612 sparse matrix of type '<class 'numpy.int64'>'
        with 58826 stored elements in Compressed Sparse Row format>
```

```
In [32]: X_test
```

```
Out[32]: <1115x7612 sparse matrix of type '<class 'numpy.int64'>'
        with 13975 stored elements in Compressed Sparse Row format>
```

```
In [33]: y_train
```

```
Out[33]: array([0, 0, 0, ..., 0, 0, 0])
```

```
In [34]: y_test
```

```
Out[34]: array([0, 0, 0, ..., 0, 0, 0])
```

## Un poco de Machine Learning

```
In [35]: clf = MultinomialNB()
```

```
In [36]: clf.fit(X_train, y_train)
```

```
Out[36]: ▼MultinomialNB
MultinomialNB()
```

```
In [38]: y_pred = clf.predict(X_test)
y_pred
```

```
Out[38]: array([0, 0, 0, ..., 0, 0, 0])
```

```
In [39]: from sklearn.metrics import accuracy_score
acc = accuracy_score(y_pred, y_test)
print(acc * 100)

98.7443946188341
```

```
In [40]: clf.score(X_test, y_test)
```

```
Out[40]: 0.9874439461883409
```

```
In [42]: aciertos = 0

for i in range(len(y_pred)):
    if y_pred[i] == y_test[i]:
        aciertos += 1

aciertos
```

```
Out[42]: 1101
```

```
In [43]: (aciertos/len(y_pred))*100
```

```
Out[43]: 98.7443946188341
```

## Calcular la matriz de confusión

```
In [44]: len(y_train)
```

```
Out[44]: 4457
```

Falsos Positivos

```
In [46]: FP = 0

for i in np.arange(len(y_test)):
    if y_test[i] == 0 and y_pred[i] == 1:
        FP += 1

FP
```

```
Out[46]: 2
```

Falsos Negativos

```
In [47]: FN = 0

for i in np.arange(len(y_test)):
    if y_test[i] == 1 and y_pred[i] == 0:
        FN += 1

FN
```

```
Out[47]: 12
```

True Positives

```
In [48]: TP = 0

for i in np.arange(len(y_test)):
    if y_test[i] == 1 and y_pred[i] == 1:
        TP += 1

TP
```

```
Out[48]: 154
```

True Negative

```
In [49]: TN = 0

for i in np.arange(len(y_test)):
    if y_test[i] == 0 and y_pred[i] == 0:
        TN += 1

TN
```

```
Out[49]: 947
```

```
In [51]: confusion_matrix = np.array([[TN, FP],
                                     [FN, TP]])
confusion_matrix
```

```
Out[51]: array([[947,  2],
               [ 12, 154]])
```

```
In [52]: ((TN + TP) / (TN+TP+FP+FN)) *100
```

```
Out[52]: 98.7443946188341
```

Forma con Sklearn

```
In [53]: from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred)
cm
```

```
Out[53]: array([[947,  2],
               [ 12, 154]])
```

## Ahora con: TfidfVectorizer

```
In [54]: X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0, test_size=0.2)
```

```
In [55]: X_train
```

```
Out[55]: 1114      No no:)this is kallis home ground.amla home to...
3589      I am in escape theatre now. . Going to watch K...
3095      We walked from my moms. Right on stagwood pass...
1012      I dunno they close ordi not... II v ma fan.
3320      Yo im right by yo work
```

```
4931      ...
3264      Match started.india <lt;#> for 2
1653      44 7732584351, Do you want a New Nokia 3510i c...
2607      I was at bugis juz now wat... But now i'm walk...
2732      :-) yeah! Lol. Luckily i didn't have a starrin...
Name: Message, Length: 4457, dtype: object
```

```
In [56]: X_test
```

```
Out[56]: 4456      Aight should I just plan to come up later toni...
690       Was the farm open?
944       I sent my scores to sophas and i had to do sec...
3768      Was gr8 to see that message. So when r u leavi...
1189      In that case I guess i'll see you at campus lodge
```

```
2906      ... ALRITE
1270      Sorry chikku, my cell got some problem thts y ...
3944      I will be gentle princess! We will make sweet ...
2124      Beautiful Truth against Gravity.. Read careful...
253       Ups which is 3days also, and the shipping comp...
Name: Message, Length: 1115, dtype: object
```

```
In [57]: y_train
```

```
Out[57]: 1114      0
690       0
944       0
3768      0
1189      0
```

```
4931      0
3264      1
1653      0
2607      0
2732      0
Name: Status, Length: 4457, dtype: int64
```

```
In [58]: y_test
```

```
Out[58]: 4456      0
690       0
944       0
3768      0
1189      0
```

```
2906      0
1270      0
3944      0
2124      0
253       0
Name: Status, Length: 1115, dtype: int64
```

```
In [59]: from sklearn.feature_extraction.text import TfidfVectorizer
```

```
In [60]: tv = TfidfVectorizer(stop_words = "english")
tv
```

```
Out[60]: ▼TfidfVectorizer
TfidfVectorizer(stop_words='english')
```

```
In [61]: X_train = tv.fit_transform(X_train)
X_test = tv.transform(X_test)
```

```
In [62]: y_train = y_train.astype("int")
y_test = y_test.astype("int")
```

```
In [63]: y_train = np.array(y_train)
y_test = np.array(y_test)
```

## Creamos el algoritmo

```
In [64]: clf = MultinomialNB()
```

```
In [65]: clf.fit(X_train, y_train)
```

```
Out[65]: ▼MultinomialNB
MultinomialNB()
```

```
In [66]: y_pred = clf.predict(X_test)
y_pred
```

```
Out[66]: array([0, 0, 0, ..., 0, 0, 0])
```

```
In [67]: from sklearn.metrics import accuracy_score
acc = accuracy_score(y_pred, y_test)
print(acc * 100)
```

```
96.59192825112108
```

```
In [68]: from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred)
cm
```

```
Out[68]: array([[949,  0],
               [ 38, 128]])
```

Creado por:

Isabel Maniega