

Funciones auxiliares para el preprocesamiento del DataSet

```
[17]: def create_prep_dataset(index_path, n_elements):
    X = []
    y = []
    indexes = parse_index(index_path, n_elements)
    for i in range(n_elements):
        print("\rParsing email: {}".format(i+1), end = '')
        mail, label = parse_email(indexes[i])
        X.append(" ".join(['Subject']) + " ".join(mail['body']))
        y.append(label)
    return X, y

[18]: # Leer únicamente un subconjunto de 1000 correos electrónicos
X_train, y_train = create_prep_dataset("datasets/datasets/trec07p/full/index", 1000)
X_train
```

Parsing email: 1000

```
[18]: ['Subject\n\n\n\n\n\nndo feel dresur perform rise occasion\n\n\n\n\n\nntri viagra\nyour anxiety thing past will\nab back old self\n\n\n\n',
      'Subject\ni ve updat gulu i check mirrors\nit seem littl tyeo debianreadm file\n\n\nexample\nhttpgulususherbrookecadebianreadme\nftptpfdrdeb  
norgdebianreadme\nntest lenni access releas disttest the\ncurr test develop snapshot name etch packag which\nhav test unstab pass autom te  
st propog to\nnthi release\n\n\netch replac lenni like readmehtml\n\n\n\n\n\nyan morin\nconsult en logiciel libre\nnyanmorinsavoirfairelinuxcom\n5  
149941556\n\n\n\n nto unsubscrib email debianmirrorsrequestlistsdebianorg\nwith subject unsubscribe troubl contact listmasterlistsdebianorg\n\n',
      'Subjectmega authenticic i a g r a discount pricec i a l i s discount pricedo miss it click here\nhttpwwmwoujskhchumcom\n\n\n authent viagra\n\nmega authenticic i a g r a discount pricec i a l i s discount pricedo miss it click here\n\n',
      'Subject\nhey billi \n\n\nit realli fun go night \nvand talk said felt\ninsecur manhood i notic toilets\ny quit small area \nworri websit i tel  
l nmci secret weapon evera 3 inch trust girls\nlover bigger one ive 5 time mani chick \nsinc i use pill year ago the packag i used\nwa 6 month  
suppli one worth extra \ncent websit httpctmaycom\nring weekend go drink \nagain let know secret \nlater dude brad\n\n',
      'Subject\nsystem home it capabl link far i \nknew i within part with respect the\n affect technolog societi scienc ad agenc cashin\ncommer  
ci and\nphotograph paint electron canvas still seem like silenc \n white black light it didnt happen althrough far p\nfect espec preclud va  
st explan i can\nnt understand peopl reli so\nnavantgard art world added writer lawyer yet re\nach full potenti i imagin futur comput screen ho  
\nw would percept artifici imag busi apart there t\nhe computer sign shop the\nb filter unnecessary labor technolog comp\nani monopoli servic  
as servic root tre\nnd western cultur sinc dark age possibl unimag\nnin human think \n\n\nfriendlly not industri isnt welcom new peopl thought c
```

Aplicar vectorización a los Datos

```
[19]: vectorizer = CountVectorizer()
      X_train = vectorizer.fit_transform(X_train)
```

Aplicar vectorización a los Datos

```
[19]: vectorizer = CountVectorizer()
      X_train = vectorizer.fit_transform(X_train)
```

```
[20]: print(X_train.toarray())
print("\nFeatures", len(vectorizer.get_feature_names_out()))

[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
\nFeatures 21970
```

```
[21]: import pandas as pd
```

```
pd.DataFrame(X_train.toarray(), columns=[vectorizer.get_feature_names_out()])
```

[illegible]

1000 rows x 21970 columns

```
[22]: ['spam',  
       'ham',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'ham',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'spam',  
       'spam']
```

```
[23]: from sklearn.linear_model import LogisticRegression

      clf = LogisticRegression()
      clf.fit(X_train, y_train)
```

```
[23]: LogisticRegression
LogisticRegression()
```

```
[24]: # Lectura de un DataSet de correos nuevos.

# Leer 1500 correos de nuestro DataSet y quedarnos únicamente con los 500 ultimos correos electrónicos, los cuales no se han utilizado para el
X, y = create_prep_dataset("datasets/datasets/trec07p/full/index", 150)
X_test = X[100:]
y_test = y[100:]

Parsing email: 150
```

```
[25]: X_test = vectorizer.transform(X_test)
```

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

```
[26]: array(['spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam',  
         'spam', 'spam', 'ham', 'spam', 'spam', 'spam', 'spam', 'spam',  
         'ham', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'ham',  
         'spam', 'spam', 'ham', 'spam', 'spam', 'ham', 'spam', 'spam',  
         'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam',  
         'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam', 'spam',  
         'spam', 'spam'], dtype='<U4')
```

```
[27]: print("Predicción\n", y_pred)
      print("\nEtiquetas reales", y_test)
```

```
Predicción
['spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam'
'ham' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'ham' 'spam' 'spam' 'spam'
'spam' 'spam' 'spam' 'ham' 'spam' 'spam' 'ham' 'spam' 'spam' 'ham' 'spam'
'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam'
'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam' 'spam']
```

[illegible]

Evaluación de resultados

```
[28]: from sklearn.metrics import accuracy_score
      print("Accuracy: {:.3f}".format(accuracy_score(y_test, y_pred)))

Accuracy: 1.000
```

5.- Aumentando el DataSet

```
[29]: # Leer 20,000 correos electrónicos
      X, y = create_prep_dataset("datasets/datasets/trec07p/full/index", 20000)

      Parsing email: 20000
```

```
[30]: # Utilizamos 15,000 correos para entrenar el algoritmo y 5,000 para realizar pruebas
      X_train, y_train = X[:15000], y[:15000]
      X_test, y_test = X[15000:], y[15000:]
```

```
[31]: vectorizer = CountVectorizer()
      X_train = vectorizer.fit_transform(X_train)
```

```
[32]: clf = LogisticRegression()
      clf.fit(X_train, y_train)
```

```
[32]: ▼ LogisticRegression ⓘ ⓘ
      LogisticRegression()
```

```
[33]: X_test = vectorizer.transform(X_test)
      y_pred = clf.predict(X_test)
```

```
[34]: print("Accuracy: {:.3f}".format(accuracy_score(y_test, y_pred)))

Accuracy: 0.989
```

