



Alumno:	García González Aarón Antonio
Grupo:	3CV9
Unidad de Aprendizaje:	Procesamiento de lenguaje natural
Profesora:	Kolesnikova Olga
Tarea 3.	Algoritmo de aprendizaje automático por modelo de regresión lineal y logística con descenso gradiente mediante scikit-learn.
Fecha:	Martes 21 de abril de 2020

## Algoritmo de aprendizaje automático por modelo de regresión lineal

Valores de predicción de algunos elementos del conjunto de pruebas

Real	Prediccion	Error
687500	349042	49.2302
671500	556826	17.0773
330000	558064	69.1102
500000	700749	40.1497
635000	1.26425e+06	99.095
380000	201910	46.8658
310000	618175	99.4114
735000	810668	10.2949
331000	354575	7.12238
627500	681364	8.58385
320000	211705	33.8422
198000	334445	68.9117
476500	416287	12.6364
375000	686336	83.0228
429950	419839	2.35175
725000	741268	2.2438
390000	475573	21.9419
185000	93120.7	49.6645
300000	198498	33.8339
286300	141224	50.6727
324950	300760	7.4443
699900	569113	18.6865
330000	460679	39.5996
215000	193002	10.2318
557500	506231	9.19622

760000	890264	17.14
199000	229200	15.1759
208000	125376	39.723
280005	220651	21.1974
1.7e+06	1.1655e+06	31.4412
395000	437815	10.8391
325000	316705	2.55219
750000	678921	9.47722
230000	217593	5.39432
845000	1.13272e+06	34.0492
765000	832689	8.84817
455000	359128	21.0708
825000	775144	6.04312
460000	465206	1.13172
170000	53903.6	68.292
674250	832676	23.4966
637800	642132	0.67919
402000	651926	62.1706
1e+06	472238	52.7762
436000	572976	31.4166
553000	360633	34.7862
825000	1.30798e+06	58.5427
700000	718845	2.69214
538000	547551	1.77532
489000	415202	15.0916
207000	277811	34.2084

Valor de error general del programa

```
Exactitud: 0.7127615362214474
(base) Aarons-MacBook-Pro:regresion_lineal aarongarcia$
```

Código

```
1. from datetime import datetime
2. import pandas as pd
3. import numpy as np
4. import random
5. import math
6. from tabulate import tabulate
7. from sklearn import datasets, linear_model, metrics
8. from sklearn import preprocessing
9. from sklearn.model_selection import train_test_split
10.
11.
12. # funcion que obtiene los datos a partir de un documento .csv
13. # el campo fecha no lo desprecia
14. def getData(name_file):
15.     global θs
16.
17.     datos = pd.read_csv(name_file)
18.
19.     matriz = np.array(datos)
20.     matriz = matriz[:,1:] # quitamos la columna id
21.
22.     for f in matriz:
23.         f[0] = datetime.strptime(f[0], '%M/%d/%Y').date() # convertimos a fecha
24.
25.     minimo = min(matriz[:, 0]) # fecha mas antigua
26.
27.     for f in matriz:
28.         f[0] = int(abs(minimo - f[0]).days) # favorecemos en numero a las fechas mas cercanas
29.
30.     random.shuffle(matriz) # revolvemos los datos
31.
32.     Y = matriz[:, 1]
33.     X = matriz[:, [0]+[f for f in range(2,len(matriz[0]))]] # quitamos la columna de Y
34.
35.     θs = np.array( [0]*(len(X[0])+1) )
36.
37.     return X,Y
38.
39.
40. # funcion que a partir de la prediccion y el dato objetivo es decir el precio de las casas,
41. # imprime una tabla comparativa
42. def printPrediction(Y_predicciones, Y_test):
43.     tabla = []
44.     for i in range(0,len(Y_test)):
45.         error_particular = abs(100-(100/(Y_test[i])*Y_predicciones[i]))
46.         tabla.append([Y_test[i],Y_predicciones[i],error_particular])
47.
48.     print(tabulate(tabla,headers=['Real', 'Prediccion', 'Error' ],tablefmt="grid", numalign="center"))
49.
50.
51. def main():
52.     name_file = "data.csv"
53.     X, Y = getData(name_file)
```

```
54. X = preprocessing.scale(X) ## feature scaling
55.
56. X_train, X_test, Y_train, Y_test = train_test_split(X,Y,test_size = 0.3)
57.
58. reg = linear_model.LinearRegression() # creamos el objeto de tipo regresion linear
59. reg.fit(X_train, Y_train) # entrenamos el modelo
60.
61. Y_predicciones = reg.predict(X_test)
62.
63. printPrediction(Y_predicciones, Y_test)
64. print('Exactitud: {}'.format(reg.score(X_test, Y_test)))
65.
66.
67. if __name__ == "__main__":
68.     main()
```

## Algoritmo de aprendizaje automático por modelo de regresión logística

### Valores de predicción de algunos elementos del conjunto de pruebas

[illegible]

0	0
0	0
1	1
0	0
0	0
0	0
0	0
1	0
0	0
0	0
0	0
0	0
0	0
0	0
1	0
0	0
0	0
0	0
0	0
0	0
0	0
1	0
0	0
0	0
0	0
0	0
0	0
1	0
0	0
0	0
0	0
0	0

	0	0
	0	0
	1	1
	0	0
	0	0
	1	1
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	1	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	0	0
	1	1
	0	0

[illegible]

Valor de error general del programa

```
0.9855072463768116
(base) Aarons-MacBook-Pro:regresion_logistica aarongarcia$
```

Código

```
1. import nltk
2. import numpy as np
3. import pandas as pd
4. from sklearn.feature_extraction.text import CountVectorizer
5. from sklearn.naive_bayes import MultinomialNB
6. from sklearn.feature_extraction.text import TfidfVectorizer
7. from sklearn.model_selection import train_test_split
8. from sklearn import datasets
9. from sklearn.linear_model import LogisticRegression
10. from sklearn.metrics import precision_score
11.
12.
13. if __name__ == "__main__":
14.
15.     df=pd.read_csv('SMSSpamCollection.txt', sep='\t', names=[ 'Status', 'Message' ])
16.
17.     #print("Spam: ", len(df[df.Status=='spam']))
18.     #print("Ham: ", len(df[df.Status=='ham']))
19.
20.     df.loc[df["Status"]=="ham", "Status",]=0
21.     df.loc[df["Status"]=="spam", "Status",]=1
22.
23.     df_x=df["Message"]
24.     df_y=df["Status"]
25.
26.     cv = TfidfVectorizer(min_df=1, stop_words='english') # escalar los datos
27.
28.     df_x = cv.fit_transform(df_x).toarray()
29.     df_y = np.array(df_y.astype('int'))
30.
31.     x_train, x_test, y_train, y_test = train_test_split(df_x, df_y, test_size=0.3)
32.
33.     logreg = LogisticRegression()
34.     logreg.fit(x_train, y_train)
35.
36.     y_predicciones = logreg.predict(x_test)
37.
38.     precision = precision_score(y_test, y_predicciones)
39.     print(precision)
```