

# Act9

Daniel Rojas

March 2025

## 1 Introducción

La regresión lineal simple es un algoritmo supervisado de aprendizaje automático que busca describir una serie de datos mediante una recta de una sola variable. La ecuación de la recta con una sola variable está descrita de la siguiente forma:

$$y = mX + b \quad (1)$$

Donde  $y$  es la variable dependiente, o los valores predichos,  $X$  la independiente, o representa los valores de entrada, y la  $b$  representa la intersección con el eje  $y$  de la recta.

Para encontrar la recta de mejor ajuste, se busca minimizar el valor de la función del error cuadrático medio, que básicamente lo que hace es obtener la media de la diferencia al cuadrado de cada predicción hecha con su respectivo valor real.

## 2 Metodología

Para completar la actividad, se usó de base el capítulo de regresión lineal del libro Aprende Machine Learning, adjuntado al final en las referencias.

De la siguiente forma se pueden cargar las librerías y dataset a utilizar:

```
1 # Imports necesarios
2 import numpy as np
3 import pandas as pd
4 import seaborn as sb
5 import matplotlib.pyplot as plt
6 plt.rcParams['figure.figsize'] = (16, 9)
7 plt.style.use('ggplot')
8 from sklearn import linear_model
9 from sklearn.metrics import mean_squared_error, r2_score
```

```
1 data = pd.read_csv("articulos_ml.csv")
```

Para empezar, hay que realizar un preprocesado de los datos. Primero hay que decidir que columna del dataset usaremos como la feature, o la variable  $X$ , y

cual como nuestra etiqueta, o la variable y. En este caso, se eligió como variable feature la cantidad de palabras de cada artículo, y como etiqueta la cantidad de veces que dicho artículo fue compartido.

Por otro lado, hay algunos datos anómalos que se alejan demasiado de los datos normales, como se puede observar en la siguiente gráfica:

```
1 #Visualizar que hay algunos datos demasiado alejados de la media de
  ↪ datos (1808 palabras), o casos an malos
2 colores=['orange','blue']
3 tamanios=[30,60]
4
5 f1 = data['Word count'].values
6 f2 = data['# Shares'].values
7
8 # Vamos a pintar en 2 colores los puntos por debajo de la media de
  ↪ Cantidad de Palabras
9 asignar=[]
10 for index, row in data.iterrows():
11     if (row['Word count']>1808):
12         asignar.append(colores[0])
13     else:
14         asignar.append(colores[1])
15
16 plt.scatter(f1, f2, c=asignar, s=tamanios[0])
17 plt.show()
```

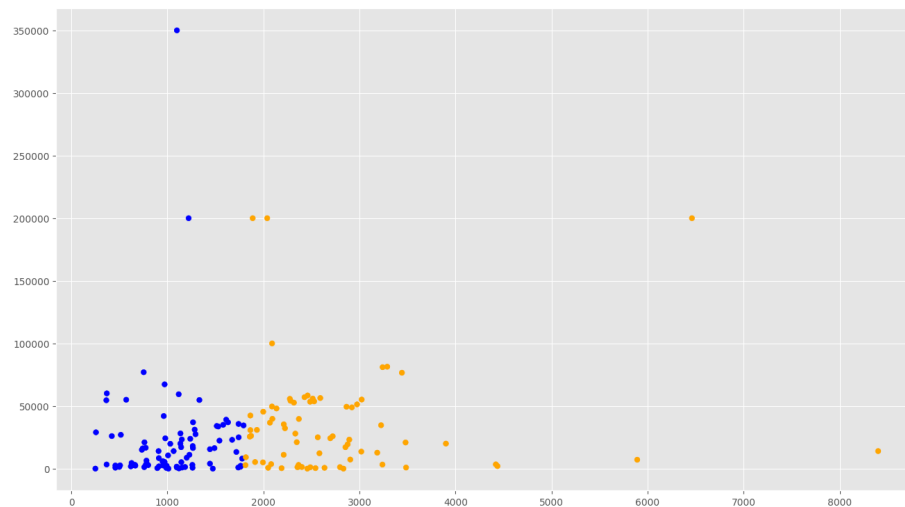


Figure 1: Gráfica de cantidad de palabras vs número de veces que se compartió un artículo

Se puede observar que mas o menos cuando las palabras son mayores a las 3500, se empiezan a observar cada vez menos datos. Lo mismo sucede cuando la cantidad de compartidos supera los 80,000. Para obtener una mejor estimación,

lo ideal sería enfocarnos en generar la regresión en donde están las mayorías de valores. Para filtrar los datos y visualizar el resultado, se puede ejecutar la siguiente línea de código:

```

1 #Datos recortados con el máximo de palabras restringido a 3500 y
  ↳ el máximo de compartidos limitado a 80000
2 filtered_data = data[(data['Word count'] <= 3500) & (data['# Shares
  ↳ ''] <= 80000)]
3
4 f1 = filtered_data['Word count'].values
5 f2 = filtered_data['# Shares'].values
6
7 # Vamos a pintar en colores los puntos por debajo y por encima de
  ↳ la media de Cantidad de Palabras
8 asignar=[]
9 for index, row in filtered_data.iterrows():
10     if(row['Word count']>1808):
11         asignar.append(colores[0])
12     else:
13         asignar.append(colores[1])
14
15 plt.scatter(f1, f2, c=asignar, s=tamamos[0])
16 plt.show()

```

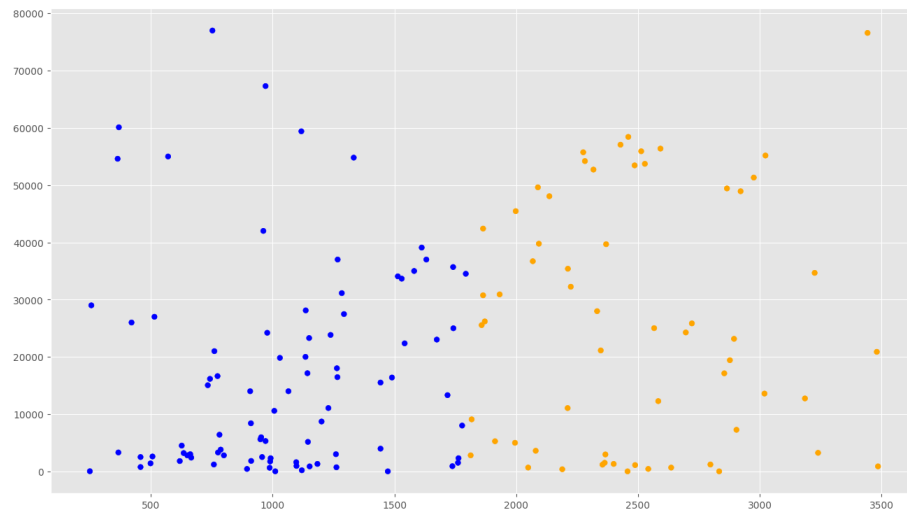


Figure 2: Gráfica de cantidad de palabras vs número de veces que se compartió un artículo filtrada

Ahora, solo queda separar los datos en el dataset de features y el de etiquetas:

```

1 #Obtener el dataset de entrenamiento
2 dataX =filtered_data[["Word count"]]
3 X_train = np.array(dataX)
4 y_train = filtered_data['# Shares'].values

```

A continuación, se entrenó un modelo de regresión lineal con el dataset filtrado, y se observó su error cuadrado y el puntaje de covarianza para conocer su rendimiento:

```
1 #Se crea el objeto de regresión lineal
2 regr = linear_model.LinearRegression()
3
4 #Se entrena el modelo
5 regr.fit(X_train, y_train)
6
7 #Se hacen predicciones para calcular el error cuadrado y el puntaje
  ↳ de varianza
8 y_pred = regr.predict(X_train)
9
10 #Se obtiene la pendiente de la recta
11 print('Coefficients: \n', regr.coef_)
12 #Se obtiene el valor de la intersección con el eje y
13 print('Independent term: \n', regr.intercept_)
14 # Error Cuadrado Medio
15 print("Mean squared error: %.2f" % mean_squared_error(y_train,
  ↳ y_pred))
16 # Puntaje de Varianza. El mejor puntaje es un 1.0
17 print('Variance score: %.2f' % r2_score(y_train, y_pred))
```

Lo anterior da como resultado una pendiente de 5.70, la intersección con el eje y en 11200.30, un error cuadrado de 372888728.34, y puntaje de covarianza de 0.06.

Se puede observar la recta de mejor ajuste junto con los puntos del dataset de la siguiente forma:

```
1 plt.scatter(X_train[:,0], y_train, c=asignar, s=tamamos[0]) #
  ↳ Crear la gráfica de dispersión de puntos
2 plt.plot(X_train[:,0], y_pred, color='red', linewidth=3) #gráfica
  ↳ la recta de mejor ajuste
3
4 plt.xlabel('Cantidad de Palabras')
5 plt.ylabel('Compartido en Redes')
6 plt.title('Regresión Lineal')
7
8 plt.show()
```

### 3 Resultados

Los valores del error cuadrado medio y de la varianza muestran que la recta de mejor ajuste no describe con mucha precisión los datos analizados. Esto tiene sentido, ya que si se observa la gráfica, los datos están demasiado dispersos como para hacer predicciones certeras con una recta.

Si se observa la imagen de la recta graficada, se puede observar que, por lo general, se puede observar que tiene una pendiente muy poco inclinada, y tiene sentido ya que visualmente se puede observar que conforme se aumenta el número

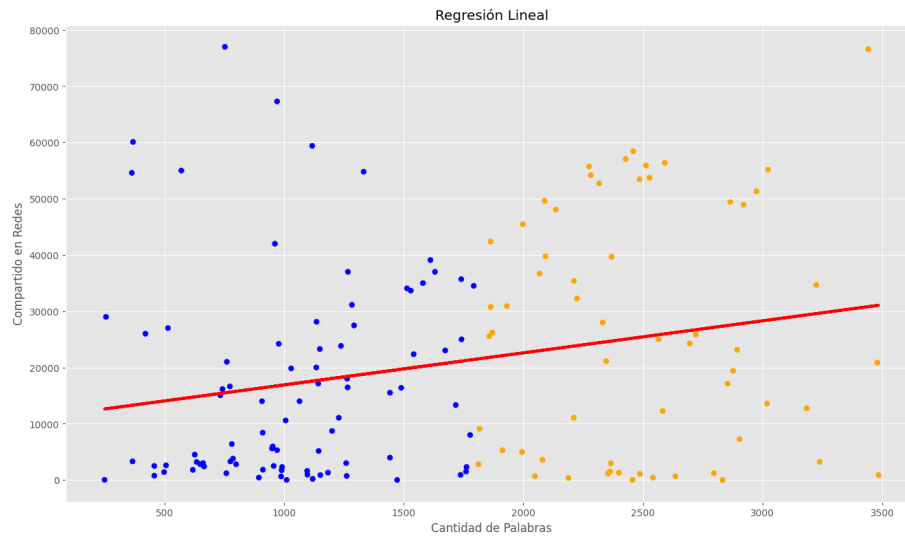


Figure 3: Gráfica de cantidad de palabras vs número de veces que se compartió un artículo filtrada con la recta de mejor ajuste

de palabras, el número de compartidos aumenta un poco; aunque los datos están demasiado dispersos como para obtener respuestas únicamente visualizando los datos.

## 4 Conclusión

La regresión lineal simple es un método que puede servir en casos donde sepamos que los datos tienden a organizarse de manera lineal al usar solamente una feature.

Por otro lado, puede ser útil para visualizar de forma general la relación lineal que poseen los datos, en vez de usarse para hacer predicciones.

Por último, la regresión lineal simple claramente es un algoritmo muy fácil de implementar, y no requiere de muchos recursos computacionales para realizarla. Sin embargo, sus resultados pueden no ser muy acordes a la realidad dependiendo del caso.

## 5 Referencias

Bagnato, J. I. (2020). Aprende Machine Learning.