

Regresión Lineal IA ARLG

Abraham López

Marzo 2025

1 Introducción

1.1 ¿Qué es la regresión lineal?

La regresión lineal es un método estadístico que permite modelar la relación entre una variable dependiente Y y una o más variables independientes X . El objetivo es encontrar la línea recta (en el caso de una sola variable independiente) que mejor se ajuste a los datos observados. En términos más simples, busca modelar cómo cambia una variable (la dependiente) en función de otra variable (la independiente). La ecuación de la regresión lineal simple se puede expresar como:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

donde:

- Y es la variable dependiente.
- X es la variable independiente.
- β_0 es la intersección con el eje Y .
- β_1 es la pendiente de la recta.
- ϵ es el término de error.

1.2 ¿Qué es la regresión lineal en el machine learning?

En el machine learning, los programas de computación denominados algoritmos analizan grandes conjuntos de datos y trabajan hacia atrás a partir de esos datos para calcular la ecuación de regresión lineal. Los científicos de datos primero entrenan el algoritmo en conjuntos de datos conocidos o etiquetados y, a continuación, utilizan el algoritmo para predecir valores desconocidos. Los datos de la vida real son más complicados que el ejemplo anterior. Es por eso que el análisis de regresión lineal debe modificar o transformar matemáticamente los valores de los datos para cumplir con los siguientes cuatro supuestos:

1.2.1 Relación lineal

Debe existir una relación lineal entre las variables independientes y las dependientes. Para determinar esta relación, los científicos de datos crean una gráfica de dispersión (una colección aleatoria de valores x e y) para ver si caen a lo largo de una línea recta. De lo contrario, puede aplicar funciones no lineales, como la raíz cuadrada o el logaritmo, para crear matemáticamente la relación lineal entre las dos variables.

1.2.2 Independencia residual

Los científicos de datos utilizan residuos para medir la precisión de la predicción. Un residuo es la diferencia entre los datos observados y el valor previsto. Los residuos no deben tener un patrón identificable entre ellos. Por ejemplo, no querrá que los residuos crezcan con el tiempo. Puede utilizar diferentes pruebas matemáticas, como la prueba de Durbin-Watson, para determinar la independencia residual. Puede usar datos ficticios para reemplazar cualquier variación de datos, como los datos estacionales.

1.2.3 Normalidad

Las técnicas de representación gráfica, como las gráficas Q-Q, determinan si los residuos se distribuyen normalmente. Los residuos deben caer a lo largo de una línea diagonal en el centro de la gráfica. Si los residuos no están normalizados, puede probar los datos para detectar valores atípicos aleatorios o valores que no sean típicos. Eliminar los valores atípicos o realizar transformaciones no lineales puede solucionar el problema.

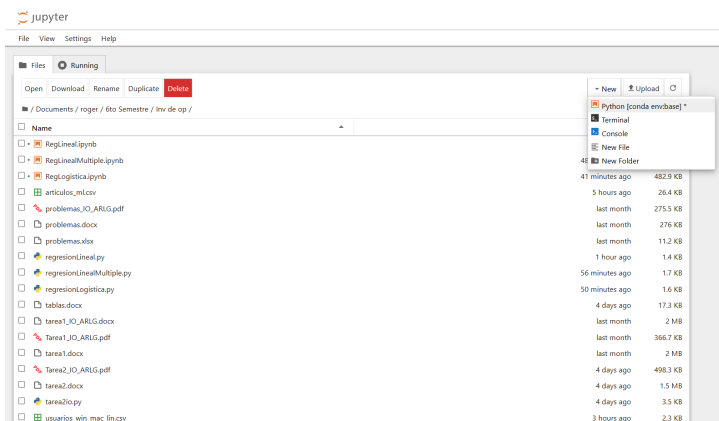
1.2.4 Homocedasticidad

La homocedasticidad supone que los residuos tienen una variación constante o desviación estándar de la media para cada valor de x . De lo contrario, es posible que los resultados del análisis no sean precisos. Si no se cumple esta suposición, es posible que tenga que cambiar la variable dependiente. Dado que la variación se produce de forma natural en grandes conjuntos de datos, tiene sentido cambiar la escala de la variable dependiente. Por ejemplo, en lugar de usar el tamaño de la población para predecir la cantidad de estaciones de bomberos en una ciudad, podría usar el tamaño de la población para predecir la cantidad de estaciones de bomberos por persona.

2 Metodología

En este ejemplo, se realizó un análisis de regresión lineal simple para explorar la relación entre el número de palabras (Word count) y el número de compartidos (Shares) en artículos de Inteligencia Artificial. El objetivo es determinar si el número de palabras en un artículo podía predecir su popularidad, medida en términos de compartidos.

Primeramente, creamos un archivo .ipynb en Jupyter para poder hacer el código y correrlo dentro de esta plataforma.



Descargamos el archivo .csv necesario para realizar las operaciones.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	Title	url	Word count	# of Links	# of common # Images v Elapsed da	# Shares																	
2	What is the	https://ble	1888	1	2	34	200000																
3	10 Companies Using		1742	9		9	5	25000															
4	How Artificial Intelligence		862	6	0	1	10	42000															
5	Obtain and the Block		1221	3		2	68	200000															
6	Nasa finds entire solar		2039	1	104	4	131	200000															
7	5 ways Data Science is		781	0		1	14	21000															
8	200 unives	https://az-	6462	600	26	2	170	200000															
9	How Much	https://cr	753	3	0	1	78	77000															
10	Tech cons	https://w	1118	2		1	62	58400															
11	Artificial Intelligence is		1581	4		2	60	35000															
12	Facebook robots shut		2090	1	95	1	267	100000															
13	A visual test	http://w	1098	4		15	1002	300000															
14	10 Breakthrough Tech		3800	30		10	62	20000															
15	Google CEO Sundar P		256	0	27	4	95	28000															
16	How machine learning		1267	2		4	124	37000															
17	New AI can	https://w	971	10		1	228	67000															
18	Researcher	https://h	268	5		1	222	60100															
19	What's	https://ec	635	3	1	2	12	3200															
20	Which One to Choose		1631	6	6	9	180	37000															
21	A computer was asked		571	4		1	277	55000															
22	AI is Inven	https://w	1333	8		4	284	54800															
23	An Artificial	https://w	364	5		2	313	54800															
24	Denmark's largest law		115	2		5	173	27000															
25	How Artificial Intelligence		1135	8	0	1	146	20000															
26	Machine Learning Zon		666	4	0	3	18	2400															

Importamos lo necesario y así poder comenzar el código.

```

Jupyter RegLineal Last Checkpoint: 1 hour ago
File Edit View Run Kernel Settings Help Trusted
JupyterLab Python [conda envbase]

[5]: import numpy as np
import pandas as pd
import seaborn as sb
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from matplotlib import cm
%matplotlib inline
plt.rcParams['figure.figsize'] = (16, 9)
plt.style.use('ggplot')
from sklearn import linear_model
from sklearn.metrics import mean_squared_error, r2_score

data = pd.read_csv("./articulos_ml.csv")
data.shape

[5]: (161, 8)

[7]: data.head()

[7]:

```

	Title	url	Word count	# of Links	# of comments	# Images video	Elapsed days	# Shares
0	What is Machine Learning and how do we use it ...	https://blog.signals.network/what-is-machine-...	1888	1	2.0	2	34	200000
1	10 Companies Using Machine Learning in Cool Ways	NaN	1742	9	NaN	9	5	25000
2	How Artificial Intelligence Is Revolutionizing...	NaN	962	6	0.0	1	10	42000
3	Dbrain and the Blockchain of Artificial Intell...	NaN	1221	3	NaN	2	68	200000
4	Nasa finds entire solar system filled with eig...	NaN	2039	1	104.0	4	131	200000

El primer paso consiste en cargar los datos y realizar una exploración inicial para entender su estructura y contenido.

```
[9]: data.describe()
```

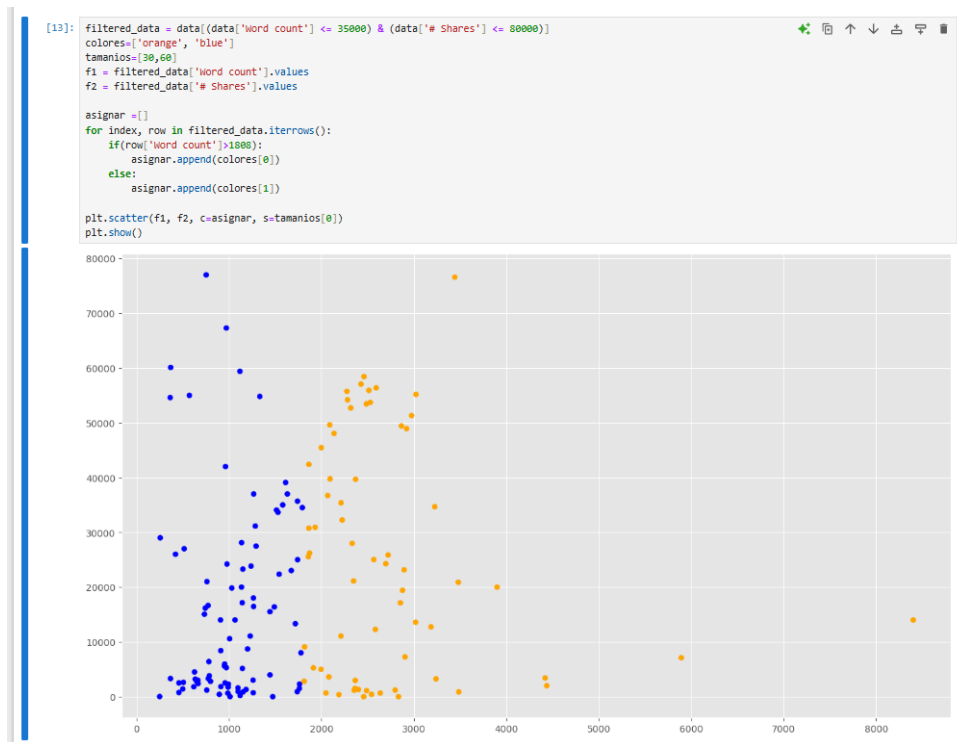
```
[9]:
```

	Word count	# of Links	# of comments	# Images video	Elapsed days	# Shares
count	161.000000	161.000000	129.000000	161.000000	161.000000	161.000000
mean	1808.260870	9.739130	8.782946	3.670807	98.124224	27948.347826
std	1141.919385	47.271625	13.142822	3.418290	114.337535	43408.006839
min	250.000000	0.000000	0.000000	1.000000	1.000000	0.000000
25%	990.000000	3.000000	2.000000	1.000000	31.000000	2800.000000
50%	1674.000000	5.000000	6.000000	3.000000	62.000000	16458.000000
75%	2369.000000	7.000000	12.000000	5.000000	124.000000	35691.000000
max	8401.000000	600.000000	104.000000	22.000000	1002.000000	350000.000000

Se realizan gráficos para visualizar la distribución de los datos y detectar patrones o valores atípicos.



Se preparan los datos para crear un gráfico de dispersión que muestre la relación entre las variables.



Se entrena un modelo de regresión lineal utilizando los datos filtrados.



Por último, predecimos el número de shares para un artículo con 2000 palabras.

Cantidad de palabras

```
[23]: y_dosmil = regr.predict([[2000]])  
      print(int(y_dosmil))  
      20661
```

3 Resultados

Tras realizar el análisis de regresión lineal, se obtuvieron los siguientes resultados:

- **Coefficiente (pendiente):** 2.04
Esto indica que, por cada palabra adicional en un artículo, se espera que el número de shares aumente en aproximadamente 2.04.
- **Término independiente (intercept):** 16575.98
Este valor representa el número de shares esperado cuando el número de palabras es 0. En este caso, si un artículo tiene 0 palabras, se esperarían aproximadamente 16,576 shares.
- **Mean Squared Error (MSE):** 382,325,815.45
El MSE es una medida del error promedio entre los valores reales y los predichos. Un valor alto como este sugiere que el modelo tiene un error significativo en sus predicciones, lo que indica que no se ajusta bien a los datos.
- **Variance Score (R^2):** 0.012777
El coeficiente de determinación (R^2) indica que solo el 1.28% de la variabilidad en el número de shares es explicada por el número de palabras. Esto sugiere que hay otros factores no incluidos en el modelo que influyen en el número de shares.

4 Conclusión

El análisis de regresión lineal que se llevó a cabo nos permitió investigar la relación entre el número de palabras en los artículos y la cantidad de shares que reciben, específicamente en el ámbito del machine learning y la inteligencia artificial. Los resultados que obtuvimos muestran que el número de palabras tiene un efecto bastante limitado en la predicción de los shares, como lo indica el bajo coeficiente de determinación ($R^2 = 0.012777$), que solo explica el 1.28% de la variabilidad en los datos. Además, el elevado error cuadrático medio ($MSE = 382,325,815.45$) sugiere que el modelo no es muy preciso en sus predicciones.

Estos hallazgos indican que el número de palabras por sí solo no es un buen predictor del número de shares. Es probable que otros factores, como la cantidad de enlaces, imágenes, comentarios o incluso el tema del artículo, tengan un impacto considerable en la popularidad de los mismos. Por lo tanto, sería recomendable ampliar el análisis para incluir estas variables adicionales en un modelo de regresión múltiple o explorar técnicas más avanzadas de machine learning que nos ayuden a captar mejor la complejidad de los datos.

En resumen, aunque el modelo de regresión lineal nos da una primera idea, es fundamental adoptar un enfoque más robusto y multidimensional para mejorar la precisión de las predicciones y entender mejor los factores que influyen en el éxito de los artículos en términos de shares.