

Actividad 10: Programando Regresión Lineal Múltiple en Python

Gerardo Enrique Torres Flores 2064063

23 de marzo de 2025

1. Introducción

Regresión con Múltiples Variables. Es un método estadístico utilizado para modelar la relación entre una variable dependiente (o objetivo) y dos o más variables independientes (o predictoras). Es una **extensión de la regresión lineal simple**, que solo incluye una variable independiente. Vamos a extender el ejercicio utilizando más de una variable de entrada para el modelo. Esto le da mayor poder al algoritmo de Machine Learning, pues de esta manera podremos obtener predicciones más complejas. Nuestra ecuación de la “recta” ahora pasa a ser:

$$Y = b + m_1X_1 + m_2X_2 + \cdots + m_nX_n$$

y, por lo tanto, deja de ser una recta (en el caso de más de una variable), para ahora ser **un plano en 3D** o un hiperplano en dimensiones mayores.

2. Metodología

Para realizar esta actividad se siguieron los siguientes pasos:

1. Importación de librerías y configuración

Para realizar este ejercicio, se importan las librerías necesarias:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sb
4 import matplotlib.pyplot as plt
5 from mpl_toolkits.mplot3d import Axes3D
6 from matplotlib import cm
7 plt.style.use('ggplot')
8 from sklearn import linear_model
9 from sklearn.metrics import mean_squared_error, r2_score
```

2. Lectura y exploración de datos

Leemos el archivo CSV, examinamos sus dimensiones y mostramos los primeros registros:

```
1 data = pd.read_csv("./articulos_ml.csv")
2 data.shape
3 data.head()
4 data.describe()
```

3. Filtrado de datos y creación de nueva variable

Para simplificar, filtramos el conjunto de datos a ciertos rangos, y creamos una nueva variable que agrupa enlaces, comentarios e imágenes:

```
1 colores=['orange','blue']
2 tamanios=[30,60]
3
4 filtered_data = data[(data['Word count'] <= 3500) & (data['#
   Shares'] <= 80000)]
5
6 # Creamos una variable 'suma' que combina los enlaces,
   comentarios e imágenes
7 suma = (filtered_data["# of Links"]
8         + filtered_data['# of comments'].fillna(0)
9         + filtered_data['# Images video'])
```

4. Preparación de datos de entrenamiento

Generamos un DataFrame para las variables de entrada (en este caso, Word count y suma) y un vector con la variable objetivo:

```
1 dataX2 = pd.DataFrame()
2 dataX2["Word count"] = filtered_data["Word count"]
3 dataX2["suma"] = suma
4 XY_train = np.array(dataX2)
5 z_train = filtered_data['# Shares'].values
```

5. Entrenamiento del modelo de Regresión Lineal Múltiple

Creamos el objeto de regresión lineal y entrenamos el modelo con las variables de entrada y la variable objetivo:

```
1 regr2 = linear_model.LinearRegression()
2 regr2.fit(XY_train, z_train)
3 z_pred = regr2.predict(XY_train)
```

6. Evaluación del modelo y visualización

Se muestran los coeficientes, el error cuadrático medio y el puntaje de varianza. Posteriormente, se grafica el plano resultante junto con los datos reales y las predicciones:

```
1 print('Coefficients: \n', regr2.coef_)
2 print("Mean squared error: %.2f" % mean_squared_error(z_train
3     , z_pred))
4 print('Variance score: %.2f' % r2_score(z_train, z_pred))
5
6 fig = plt.figure()
7 ax = fig.add_subplot(111, projection='3d')
8
9 xx, yy = np.meshgrid(np.linspace(0, 3500, num=10),
10     np.linspace(0, 60, num=10))
11
12 nuevoX = (regr2.coef_[0] * xx)
13 nuevoY = (regr2.coef_[1] * yy)
14 z = (nuevoX + nuevoY + regr2.intercept_)
15
16 ax.plot_surface(xx, yy, z, alpha=0.2, cmap='hot')
17 ax.scatter(XY_train[:, 0], XY_train[:, 1], z_train, c='blue',
18     s=30, label="Datos reales")
19 ax.scatter(XY_train[:, 0], XY_train[:, 1], z_pred, c='red', s
20     =40, label="Predicciones")
```

```
18 ax.view_init(elev=30., azimuth=65)
19
20 ax.set_xlabel('Cantidad de Palabras')
21 ax.set_ylabel('Cantidad de Enlaces, Comentarios e Imágenes')
22 ax.set_zlabel('Compartido en Redes')
23 ax.set_title('Regresión Lineal con Múltiples Variables')
24 plt.tight_layout()
25 plt.show()
```

Luego, para predecir cuántos # Shares se obtienen con 2000 palabras y un total de 20 elementos (10 enlaces, 4 comentarios, 6 imágenes), hacemos:

```
1 z_Dosmil = regr2.predict([[2000, 10+4+6]])
2 print(int(z_Dosmil))
```

6. Comparación con el modelo Simple

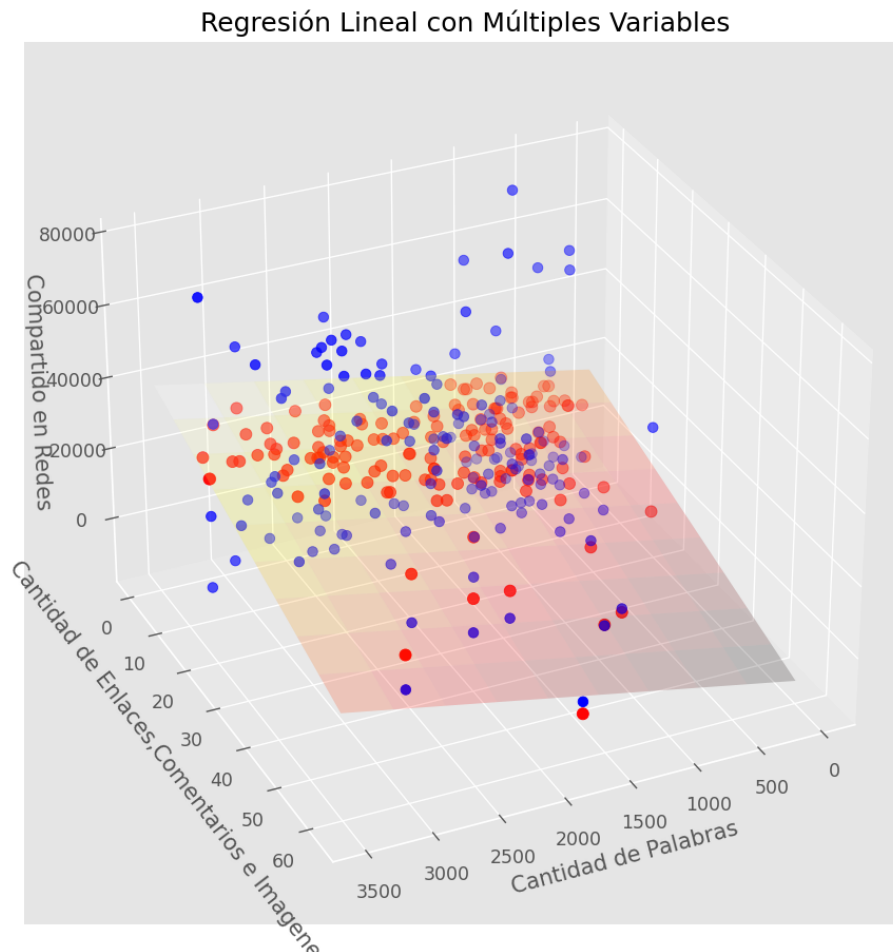
Finalmente, podemos comparar los resultados del modelo de Regresión Lineal Simple con la Regresión Lineal Múltiple:

```
1 # Restamos los errores calculados antes:
2 # Obviamente, "menos error" es mejor
3 mejoraEnError = mean_squared_error(y_train, y_pred) -
4                 mean_squared_error(z_train, z_pred)
5 print(mejoraEnError)
6
7 # También calculamos la mejora en la varianza:
8 mejoraEnVarianza = r2_score(z_train, z_pred) - r2_score(
9                 y_train, y_pred)
10 print(mejoraEnVarianza)
11 # Aunque no parezca mucho, recordemos que el valor más alto
12 # que se puede obtener es 1.0
13
14 # Finalmente, mejoramos en nuestra predicción de un
15 # artículo de 2.000 palabras,
16 # pues aunque disminuyen los "Shares" que obtendremos en el 2
17 # do modelo,
18 # seguramente ser un valor más cercano a la realidad
19 diferenciaComparar = z_Dosmil - y_Dosmil
20 print(int(diferenciaComparar))
```

3. Resultados

- **Coeficientes (Pendientes):** Se obtuvieron dos coeficientes correspondientes a las dos variables de entrada (**Word count** y **suma**).
- **Error Cuadrático Medio (MSE):** Indica la dispersión promedio de los errores (la diferencia entre los valores reales y los valores predichos).
- **Puntaje de Varianza (R^2):** Cuanto más se acerque a 1, mejor es el ajuste de nuestro modelo a los datos.

En la visualización 3D muestra cómo la combinación de las dos variables (**Word count** y **suma**) afecta la variable de salida (**# Shares**).



4. Conclusión

Gracias a que extendimos la regresión lineal simple a una **Regresión Lineal Múltiple**, incorporando más de una variable de entrada para el modelo predictivo, pudimos capturar relaciones más complejas entre los datos y mejorar la capacidad de predicción. El modelo entrenado no solo toma en cuenta la cantidad de palabras de un artículo, sino también el número total de enlaces, comentarios e imágenes, obteniendo un ajuste representado por un plano 3D. Los resultados, tanto numéricos (coeficientes, error y R^2) como visuales, demuestran cómo el modelo logra capturar la tendencia de los datos y, a su vez, se muestra capaz de predecir valores futuros como el número de compartidos para un artículo con características específicas.

```
1 Coefficients: [ 6.63216324 -483.40753769]
2 Mean squared error: 352122816.48
3 Variance score: 0.11
4 Prediction_2000_shares: 20518
5 Mejor_Error: 20765911.860715985
6 Mejora_Varianza: 0.052615337462582956
7 Diferencia_Predicci n: -2077
```