

Código de Regresión Lineal

Manuel Montufar Galnares

Junio 2024

El presente trabajo es para mostrar el código utilizado para hacer una regresión lineal sin usar un framework como sklearn o scipy, sino simples arreglos y operaciones que nos permitan predecir un valor.

La base de datos utilizada es un csv con información de varios vehículos y contiene un selling price que será nuestra variable a predecir. Cabe aclarar que esta base ya había sido en su mayoría limpiada para otro trabajo del semestre pasado, así que podré ahorrarme todo el trabajo del análisis exploratorio y limpieza de la base.

Dicho documento con la limpieza de la base también será subido al repositorio de github en caso de que el presente profesor o lector de este documento desee verlo, sin más que añadir comencemos.

```
"""
Importacion de librerias
"""

import polars as pl # Lectura y Transformacion de Datos
import matplotlib.pyplot as plt # Visualizacion de Datos
import seaborn as sns
import numpy as np # Manejo de arreglos y operaciones de vectores
```

Comenzamos importando las librerías que utilizaremos para el trabajo.

```
"""
Lectura de la base de datos
"""

# Establecemos la ruta de la base
ruta = r'C:\Users\Manuel\Montufar\Documents\ProyectosManu\Escuela\
Concentracion\Clase\Uresti\Implementacion\Vehiculos.csv'
# Definimos el data frame y leemos la base
df = pl.read_csv(f'{ruta}')
# Imprimimos el data frame y el tamaño de la base
print(df)
```

Hacemos la lectura de la base de datos con Polars y lo imprimimos para obtener el siguiente data frame:

shape: (558_837, 10)

year	make	model	body	...	color	interior	mmr	sellingprice
---	---	---	---	---	---	---	---	---
i64	str	str	str		str	str	i64	i64
2015	Kia	Sorento	SUV	...	White	black	20500	21500
2015	Kia	Sorento	SUV	...	White	beige	20800	21500
2014	BMW	3 Series	Sedan	...	Gray	black	31900	30000
2015	Volvo	S60	Sedan	...	White	black	27500	27750
2014	BMW	6 Series Gran Coupe	Sedan	...	Gray	black	66000	67000
...
2015	Kia	K900	Sedan	...	Silver	black	35300	33000
2012	Ram	2500	Crew Cab	...	White	black	30200	30800
2012	BMW	X5	SUV	...	Black	black	29800	34000
2015	Nissan	Altima	sedan	...	White	black	15100	11100
2014	Ford	F-150	SuperCrew	...	Gray	gray	29600	26700

Figura 1: Data Frame Crudo

```

"""
Análisis de Transformación de Datos
"""
# Esta base ya había sido previamente limpiada para un proyecto anterior,
# adjuntare el EDA de dicho trabajo

# Conteo de valores nulos por columna
nulos = df.null_count()
# Imprime el resultado
print(f"El conteo de valores nulos es: {nulos}")

# Encuentra la moda de la columna 'interior'
moda_interior = df.select(pl.col("interior").mode()).item()
print(f"La moda de la columna interior es: {moda_interior}")

# Rellena los valores nulos en la columna 'interior' con la moda
df = df.with_columns(
    pl.col("interior").fill_null(moda_interior)
)

# Verifica que ya no haya valores nulos
nulos_post = df.null_count()
print(f"Los valores nulos después de aplicar la moda son: {nulos_post}")

# Obtiene los valores únicos por columna
valores_unicos = {col: df[col].unique().to_list() for col in df.columns}

# Imprime los valores únicos
for col, valores in valores_unicos.items():
    print(f"Columna '{col}' tiene {len(valores)} valores únicos:
    {valores[:10]}{'...' if len(valores) > 10 else ''}")

# Vemos estadísticas descriptivas
print(df.describe())

```

Hacemos transformaciones de datos, checamos valores nulos y los rellenamos con la moda, contamos los valores unicos y vemos estadísticas descriptivas de las variables

Obtenemos el siguiente output:

```
El conteo de valores nulos es: shape: (1, 10)
```

year	make	model	body	...	color	interior	mmr	sellingprice
...
u32	u32	u32	u32	...	u32	u32	u32	u32

```

0 0 0 0 0 -- 0 17077 0 0

```

La moda de la columna interior es: black
 los valores nulos despues de aplicar la moda son: shape: (1, 10)

year	make	model	body	...	color	interior	mmr	sellingprice
...
u32	u32	u32	u32	...	u32	u32	u32	u32

```

0 0 0 0 0 -- 0 0 0 0

```

Columna 'year' tiene 34 valores únicos: [1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991] ...
 Columna 'make' tiene 61 valores únicos: ['Dodge', 'Isuzu', 'Infiniti', 'Plymouth', 'Gmc', 'Hummer', 'Saturu', 'BMW', 'Pontiac', 'Kia'] ...
 Columna 'model' tiene 973 valores únicos: ['pilot', 'XG300', 'Insight', 'taurus', 'focus', '89 tribeca', 'optima', 'RC 350', 'Prizm'] ...
 Columna 'body' tiene 27 valores únicos: ['Minivan', 'quad cab', 'c37 coupe', 'Wagon', 'cvt wagon', 'Crew Cab', 'T32 Sport Wagon', 'G Convertible', 'creeamax cab', 'Coupe'] ...
 Columna 'transmission' tiene 2 valores únicos: ['manual', 'automatic'] ...
 Columna 'odometer' tiene 172278 valores únicos: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10] ...
 Columna 'color' tiene 19 valores únicos: ['Purple', 'Pink', 'Grey', 'Gold', 'Green', 'Charcoal', 'brown', 'Turquoise', 'Lime', 'Brown', 'Red'] ...
 Columna 'interior' tiene 16 valores únicos: ['beange', 'burgundy', 'silver', 'gold', 'black', 'brown', 'yellow', 'tan', 'white', 'blue'] ...
 Columna 'mmr' tiene 1181 valores únicos: [25, 50, 75, 100, 125, 150, 175, 200, 225, 250] ...
 Columna 'sellingprice' tiene 1887 valores únicos: [1, 100, 125, 150, 175, 200, 225, 250, 275, 300] ...
 shape: (9, 11)

statistic	year	make	model	...	color	interior	mmr	sellingprice
...
str	f64	str	str	...	str	str	f64	f64

count	558837.0	558837	558837	...	558837	558837	558837.0	558837.0
null_count	0.0	0	0	...	0	0	0.0	0.0
mean	2018.83927	null	null	...	null	null	13769.27418	35811.326356
std	3.966864	null	null	...	null	null	9679.646165	9749.399466
min	1982.0	Acura	09-mar	...	Beige	beige	25.0	1.0
25%	2007.0	null	null	...	null	null	7180.0	6900.0
50%	2012.0	null	null	...	null	null	12250.0	12100.0
75%	2013.0	null	null	...	null	null	18300.0	18200.0
max	2015.0	Vw	yellow	...	Yellow	yellow	182000.0	200000.0

Figura 2: Transformación de Datos

```

"""
Identificación de Outliers en las variables a elegir
"""
# Selecciona las columnas relevantes
df = df.select(["mmr", "sellingprice"])

# Asignamos el mismo tipo de Dato
mmr = df['mmr'].to_numpy().astype(float)
sellingprice = df['sellingprice'].to_numpy().astype(float)

# Funcion para calcular el IQR y detectar outliers
def detectar_outliers(data):
    q1 = np.percentile(data, 25)
    q3 = np.percentile(data, 75)
    iqr = q3 - q1
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr
    return (data < lower_bound) | (data > upper_bound), lower_bound, upper_bound

# Identificar outliers
outliers_mmr, lower_bound_mmr, upper_bound_mmr = detectar_outliers(mmr)
outliers_sellingprice, lower_bound_sp, upper_bound_sp = detectar_outliers(sellingprice)

# Imprimir resultados
print(f"Numero de outliers en 'mmr': {np.sum(outliers_mmr)}")
print(f"Limites para outliers en 'mmr': {lower_bound_mmr} (inferior) y {upper_bound_mmr} (superior)")

print(f"Numero de outliers en 'sellingprice': {np.sum(outliers_sellingprice)}")
print(f"Limites para outliers en 'sellingprice': {lower_bound_sp} (inferior) y {upper_bound_sp} (superior)")

# Filtrar valores negativos en 'mmr'
valores_negativos_mmr = df.filter(pl.col("mmr") < 0)
print(f"Valores negativos en 'mmr': {valores_negativos_mmr.shape[0]}")
print(valores_negativos_mmr.head())

# Filtrar valores negativos en 'sellingprice'
valores_negativos_sellingprice = df.filter(pl.col("sellingprice") < 0)
print(f"Valores negativos en 'sellingprice': {valores_negativos_sellingprice.shape[0]}")
print(valores_negativos_sellingprice.head())

# Eliminar outliers basados en los limites
df_limpio = df.filter(
    (pl.col("mmr") >= 0) & (pl.col("mmr") <= 35100) &
    (pl.col("sellingprice") >= 0) & (pl.col("sellingprice") <= 35150)
)

print(f"Numero de filas despues de eliminar outliers: {df_limpio.shape[0]}")

```

Aquí verificamos outliers para nuestras variables asignadas ya que en la estadística descriptiva, había un cambio muy significativo entre el 75% y el máximo cuartil, nos encargamos de transformar estos outliers en valores que puedan ser utilizados en la regresión y estandarizamos en base a los límites y nos regresa el número de filas totales después de quitar y transformar outliers

Obtenemos el siguiente output

```

Número de outliers en 'mmr': 16315
Límites para outliers en 'mmr': -9700.0 (inferior) y 35100.0 (superior)
Número de outliers en 'sellingprice': 16354
Límites para outliers en 'sellingprice': -10050.0 (inferior) y 35150.0 (superior)
Valores negativos en 'mmr': 0
shape: (0, 2)

```

mmr	sellingprice
---	---
i64	i64

```


```

mmr	sellingprice
---	---
i64	i64

```


```

mmr	sellingprice
---	---
i64	i64

mmr	sellingprice
---	---
i64	i64

---	---
i64	i64

i64	i64

```

Valores negativos en 'sellingprice': 0
Valores negativos en 'sellingprice': 0
Valores negativos en 'sellingprice': 0
shape: (0, 2)

```

mmr	sellingprice
---	---
mmr	sellingprice
---	---
---	---
i64	i64
i64	i64

```

Número de filas después de eliminar outliers: 540900

```

Figura 3: Transformación de Outliers

```
"""
Construccion de un modelo de regresion lineal sin framework
"""

# Convertir a numpy arrays
X = df["mmr"].to_numpy().reshape(-1, 1) # Caracteristica
y = df["sellingprice"].to_numpy()      # Variable objetivo

# Calcular medias
x_mean = X.mean()
y_mean = y.mean()

# Restar las medias
X_centered = X - x_mean
y_centered = y - y_mean

# Calcular la pendiente (beta1) y la interseccion (beta0)
numerador = (X_centered.flatten() * y_centered).sum()
denominador = (X_centered.flatten() ** 2).sum()
beta1 = numerador / denominador
beta0 = y_mean - beta1 * x_mean

# Realizar predicciones
y_pred = beta0 + beta1 * X

# Evaluar el modelo
mse = ((y - y_pred.flatten()) ** 2).mean()
r2 = 1 - (ss_residual / ss_total)

# Imprimir resultados
print(f"Error cuadrático medio (MSE): {mse}")
print(f"Coeficiente de determinación (R^2): {r2}")
```

Hacemos el modelo de regresión utilizando lo siguiente: Convertimos las variables que elegimos a arreglos de numpy, calculamos las medias y las restamos al arreglo calculamos la pendiente y la intersección el numerador lo obtenemos con el la suma del producto de X despues de aplica el metodo flatten que regresa un vector de 1 dimensión por y, despues para el denominador, sera la suma del vector X al cuadrado, despues obtenemos las pendientes con las formulas de pendiente. Realizamos las predicciones con la ecuacion. Obetenmos el error cuadratico medio y la \hat{r}^2 para evaluar el modelo

El output:

```
Error cuadrático medio (MSE): 3085713.241441639
Coeficiente de determinación (R^2): 0.9675361060680938
```

Figura 4: Evaluación de Regresión

```
"""
Visualizacion de los datos
"""

# Configuracion de graficos
plt.figure(figsize=(10,6))

# Graficar los puntos de datos reales
plt.scatter(X,y,color='blue',label='Datos reales')

# Graficar la linea de regresion
plt.plot(X,y_pred,color='red',linewidth=2,label='Linea de regresion')

# Etiquetas y titulo
plt.xlabel('Precio de Mercado (mmr)')
plt.ylabel('Precio de Venta (sellingprice)')
plt.title('Regresion Lineal: Precio de Mercado vs. Precio de Venta')
plt.legend()

# Mostrar el grafico
plt.show()
```

Creamos el grafico con los valores reales y la linea de regresión para visualizar los datos resultantes
Grafico:

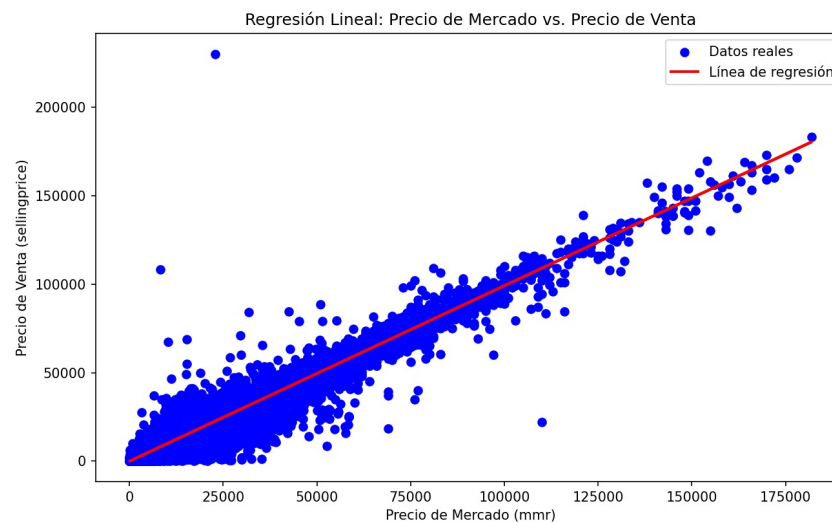


Figura 5: Grafico de Regresión

```
"""
Prediccion
"""

#####
# Definir un nuevo valor de mmr para la prediccion
nuevo_mmr = float(input('Ingresa el valor de mmr que consideres: '))
#####
# Inputa este valor con el valor que deseas predecir

# Calcular el precio de venta predicho usando la formula de la regresion
precio_predicho = beta0 + beta1 * nuevo_mmr

# Imprimir el resultado
print(f"El precio de venta predicho para un mmr de {nuevo_mmr} es {precio_predicho:.2f}")
```

Finalmente hacemos una predicción de valor de venta para un mmr específico con entrada de usuario Output:

```
Ingresa el valor de mmr que consideres: 17500
El precio de venta predicho para un mmr de 17500.0 es 17307.44
```

Figura 6: Predicción

Esto fue todo el código de regresión lineal, la evaluación del modelo pudo ser mejor, quizás utilizando sklearn podríamos obtener un mejor resultado, ya que cuando hice la práctica con Mathematica los MSE eran más bajos, pero quizás puede ser un tema de estandarización. Fue una práctica divertida para entender a fondo como funcionan algunos algoritmos clásicos en su forma más cruda y con las operaciones que hay por detrás.