

# Regresión Lineal: Predicción de Precios de Autos ¶

Una tarea bastante común en el contexto de Machine Learning es la predicción de una variable según algunos descriptores. En esta ocasión partiremos de un dataset que involucra ocho variables (algunas numéricas, otras categóricas) y el precio de un automóvil, el cual se desea estimar empleando un modelo de regresión lineal.

El objetivo de este cuaderno es que usted aprenda:

- Algunas maneras de lidiar con descriptores categóricos no-numéricos.
- Uso de la librería Pandas para manejar datos tabulares en formato .csv.
- Uso de la librería Scikit Learn para preprocesamiento y entrenamiento de un modelo de regresión lineal.
- Funcionamiento e implementación del algoritmo Descenso de Gradiente Estocástico para un modelo lineal básico.

Debe completar las celdas vacías y seguir las instrucciones anotadas en el cuaderno.

La fecha límite de entrega es el día **9 de septiembre** y se realizará a través de Bloque Neón.

```
In [ ]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
import sklearn
from sklearn.preprocessing import OrdinalEncoder, MinMaxScaler, RobustScaler
from sklearn.model_selection import train_test_split
```

```
In [ ]: df = pd.read_csv('carDataset/toyota.csv')
df.head()
```

```
In [ ]: df.info()
```

Podemos observar algunas categorías numéricas, y otras de tipo *object* (string).

```
In [ ]: categ = ['model', 'year', 'transmission', 'fuelType', 'engineSize']
str_categ = ['model', 'transmission', 'fuelType']
numer = ['price', 'mileage', 'tax', 'mpg']
```

# Parte 1

A continuación se realizó el ajuste de los datos para un modelo de regresión lineal definido mediante la librería SciKit-Learn. Se observan dos métricas de precisión, en este caso Error Cuadrático Medio (RMSE) y  $R^2$  (puede leer más sobre esta métrica en [ScikitLearn: r2\\_score \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2\\_score.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2_score.html)). Se le entregarán los dos primeros casos como ejemplo, usted debe implementar los demás a partir de estos primeros.

## Casos:

- Caso 1: Codificación One-Hot (MinMaxScaler)
- Caso 2: Codificación Ordinal (MinMaxScaler)
- Caso 3: Codificación One-Hot (RobustScaler)
- Caso 4: Codificación Ordinal (RobustScaler)
- Caso 5: Entrenar con un 30% (Codificación Ordinal, MinMaxScaler)
- Caso 6: Entrenar con un 55% (Codificación Ordinal, MinMaxScaler)
- Caso 7: Entrenar con un 80% (Codificación Ordinal, MinMaxScaler)

Analice los resultados.

## Caso 1:

### Codificación One-Hot

Una de las formas de codificar categorías no numéricas se conoce como *one-hot encoding*, en donde se crea una columna para cada valor distinto que exista en la característica que estamos codificando y, para cada registro, marcar con un 1 la columna a la que pertenezca dicho registro y dejar las demás con 0.

Igualmente, en este caso se realizará un escalamiento de los datos utilizando un `MinMaxScaler`, investigue más sobre esta función en: [ScikitLearn: MinMaxScaler \(https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html).

```
In [ ]: df_ohe = pd.get_dummies(df)
        scaler = MinMaxScaler()
        df_scl_ohe = scaler.fit_transform(df_ohe)
        df_scl_ohe = pd.DataFrame(df_scl_ohe, columns = df_ohe.columns)
        df_scl_ohe.head()
```

### Separación de Datos

```
In [ ]: X = df_scl_oe.drop(['price'], axis=1)
y = df_scl_oe['price']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
0.2, random_state=22)
```

```
In [ ]: from sklearn.metrics import mean_squared_error, r2_score
from sklearn.linear_model import LinearRegression

model_1 = LinearRegression()

model_1.fit(X_train, y_train)
preds = model_1.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, preds))
rs = model_1.score(X_test, y_test)
print(f'RMSE: {round(rmse,3)}')
print(f'R^2: {round(rs,3)}\n')
```

## Caso 2

### Codificación Ordinal

```
In [ ]: oe = OrdinalEncoder()
df_oe = df.copy()
df_oe[str_categ] = oe.fit_transform(df_oe[str_categ])

x = df_oe.values #returns a numpy array
min_max_scaler = MinMaxScaler()
x_scaled = min_max_scaler.fit_transform(x)
df_scl_oe = pd.DataFrame(x_scaled, columns=df_oe.columns)
df_scl_oe.head()
```

### Separación de Datos

```
In [ ]: X = df_scl_oe.drop(['price'], axis=1)
y = df_scl_oe['price']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
0.2, random_state=22)
```

```
In [ ]: model_2 = LinearRegression()

model_2.fit(X_train, y_train)
preds = model_2.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, preds))
rs = model_2.score(X_test, y_test)
print(f'RMSE: {round(rmse,3)}')
print(f'R^2: {round(rs,3)}\n')
```

## Caso 3

### Codificación One-Hot

In [ ]:

### Separación de Datos

In [ ]:

## Caso 4

### Codificación Ordinal

In [ ]:

### Separación de Datos

In [ ]:

## Caso 5

### Codificación Ordinal

In [ ]:

### Separación de Datos

In [ ]:

## Caso 6

### Codificación Ordinal

In [ ]:

## Separación de Datos

In [ ]:

## Caso 7

### Codificación Ordinal

In [ ]:

## Separación de Datos

In [ ]:

## Parte 2

Ahora usted debe desarrollar su propia implementación del método Descenso de Gradiente estocástico. Para esta ocasión debe utilizar una codificación One-Hot, un RobustScaler y realizar el entrenamiento a partir del 80% de los datos. Al final del entrenamiento, en la lista `errores` se deben tener los valores de la función de error para cada iteración y así poder observar el progreso gráficamente.

```
In [ ]: df_ohe = pd.get_dummies(df)
        scaler = RobustScaler()
        df_scl_ohe = scaler.fit_transform(df_ohe)
        df_scl_ohe = pd.DataFrame(df_scl_ohe, columns = df_ohe.columns)

        X = df_scl_ohe.drop(['price'], axis=1)
        y = df_scl_ohe['price']
        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
0.2, random_state=22)
```



```
In [ ]: plt.figure(figsize=(15,5))  
plt.ploterrores)  
plt.xlabel('# Iteraciones')  
plt.ylabel('Función de Error')  
plt.title('Descenso de Gradiente Estocástico')
```

```
In [ ]:
```