



Tecnológico de Monterrey

Evidencia 1. Portafolio de análisis

Inteligencia artificial avanzada para la ciencia de datos I (Gpo 101)

Aylin Camacho Reyes - A01379272 (The cool other)

Ivan Mauricio Amaya Contreras

08 de Septiembre del 2022, Monterrey N.L.

Durante cinco semanas hemos visto como funcionan los modelos de predicción. En este caso se utilizó un conjunto de datos que contenía información sobre los componentes del vino, para el análisis de los datos se usaron frameworks de python que facilitaran el proceso.

Exploración de datos

Esta fase nos ayuda a conocer la información que nos están brindando, saber si tiene variables cualitativas o cuantitativas, como se encuentran distribuidos los datos por medio de la estadística descriptiva y si existen inconsistencias como valores nulos.

```
#exploring the data

#data types
print(df.info())

#statistical summary
print(df.describe())

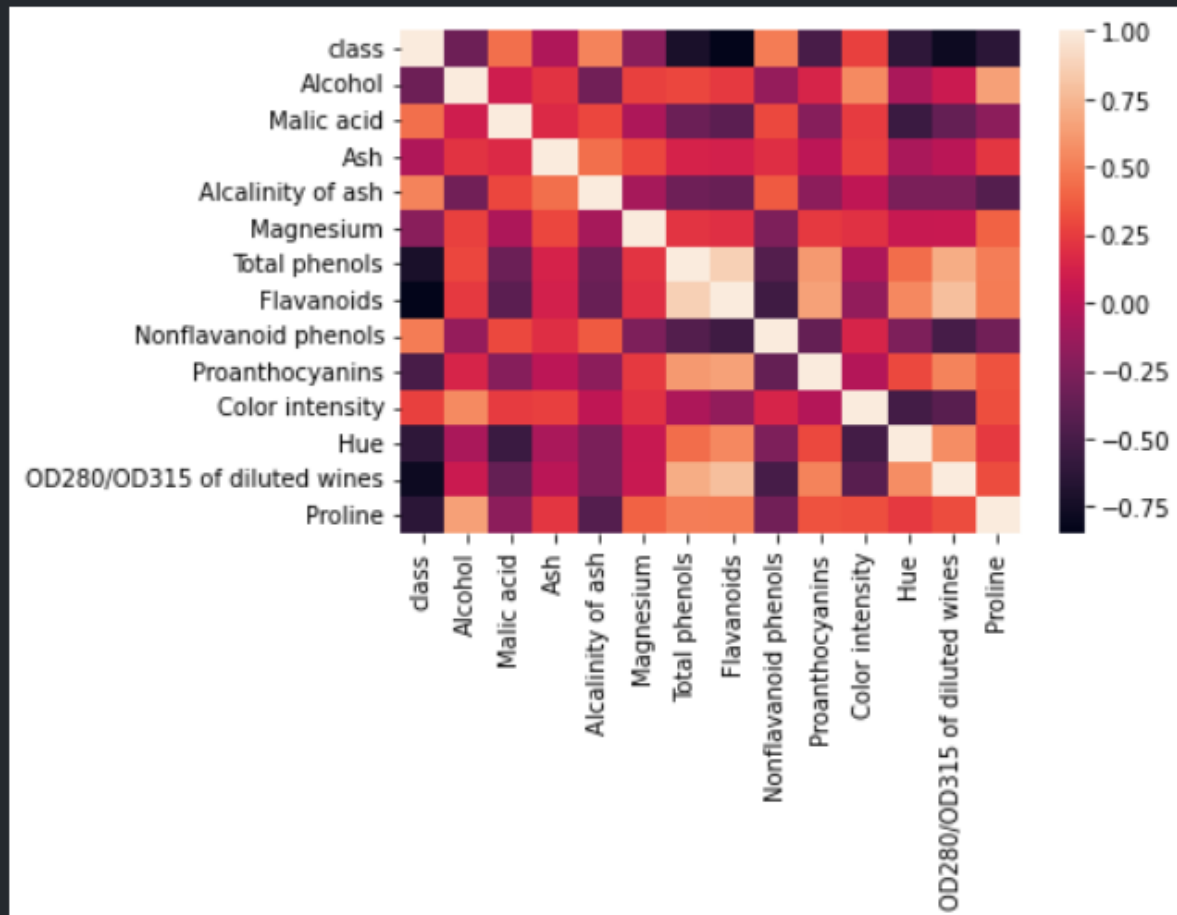
#inconsistencies
print(df.isnull().sum())
print(df.isna().sum())
```

También para explorar los datos se realizó un heat map para observar la relación entre las variables con el coeficiente de Pearson, en este caso no consideramos que las dos variables que mostraban mayor correlación fueran a ser las únicas que determinaran el valor de las variables dependientes.

```
#Correlacion de variables
corr = df.corr() # Coeficiente de Pearson
sb.heatmap(corr)
```

✓ 0.6s

<AxesSubplot:>



Separación del dataset (Train/Test/Validation)

Una vez entendidos los datos se pasan a la parte del modelado. Primero se separa el dataset entre variables dependientes e independientes, en este caso escogimos como variable dependiente el nivel de alcohol del vino y como variables independientes el resto de columnas a excepción de "class" ya que era una variable de tipo cualitativa.

```
#modeling

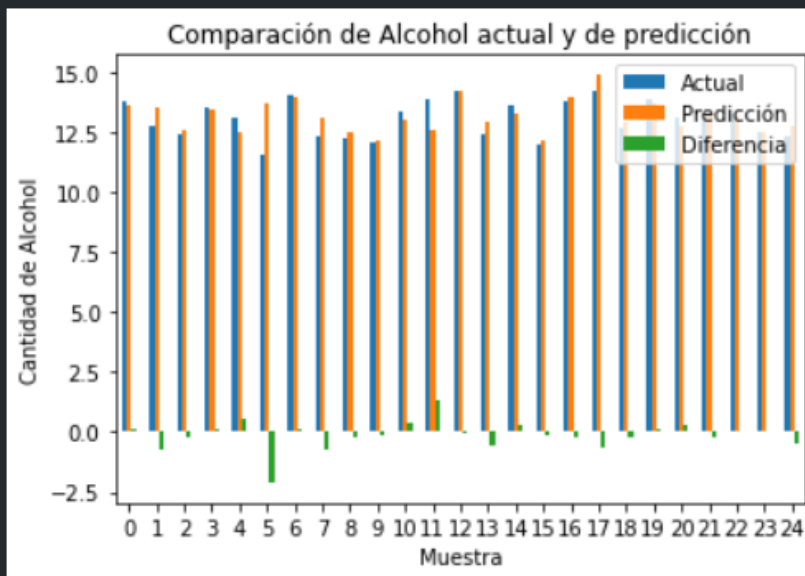
x = df[["Malic acid", "Ash", "Alcalinity of ash", "Magnesium", "Total phenols", "Flavanoids",
        "Nonflavanoid phenols", "Proanthocyanins", "Color intensity", "Hue", "OD280/OD315 of diluted wines", "Proline"]].values
y = df[["Alcohol"]].values

✓ 0.4s

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)

modelo_regresion = LinearRegression() # modelo de regresión
```

Para iniciar el modelo se declaran las variables independientes como “x” y las dependientes como “y”, para separar los datos en los que serán usados para prueba (20%) y los de validación (80%) con la ayuda de las librerías de sklearn.



Puntaje del r2 0.422535619638887

Error promedio 0.39482077365100327

Error de la raíz cuadrada del promedio of is 0.6283476534936716

Diagnóstico del grado de bias o sesgo

```
# calculate the bias and variance
bias = np.mean(y_test - y_pred)
print('Bias: %f' % bias)
```

✓ 0.7s

Bias: -0.057746

El sesgo negativo nos indica que gran parte de los datos están por debajo de la media, por lo que se encuentran por debajo del valor real.

Diagnóstico sobre el grado de varianza

```
variance = np.mean((y_test - y_pred)**2)
print('Variance: %f' % variance)
```

✓ 0.5s

Variance: 0.394821

Este valor nos indica que tanto cambiarían los valores de predicción si se cambian los datos de entrenamiento, este valor no es óptimo, pero se puede mejorar.

Diagnóstico sobre el nivel de ajuste del modelo

El modelo es overfitting ya que la varianza es alta, en este caso lo que se tendría que hacer para mejorar las predicciones del modelo es buscar disminuir la varianza para aumentar el bias o viceversa.