

ML2 Regression Analysis

October 26, 2023

```
[1]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.metrics import r2_score, accuracy_score
import warnings
warnings.filterwarnings("ignore")

# Load the diabetes dataset
data = pd.read_csv("C:/Users/hp/Downloads/Practical_Data/diabetes.csv")
data.describe()
```

C:\Users\hp\anaconda3\lib\site-packages\scipy__init__.py:146: UserWarning: A NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (detected version 1.25.2

warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}")

```
[1]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin \
count	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479
std	3.369578	31.972618	19.355807	15.952218	115.244002
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000
75%	6.000000	140.250000	80.000000	32.000000	127.250000
max	17.000000	199.000000	122.000000	99.000000	846.000000

	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000
mean	31.992578	0.471876	33.240885	0.348958
std	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.078000	21.000000	0.000000
25%	27.300000	0.243750	24.000000	0.000000
50%	32.000000	0.372500	29.000000	0.000000
75%	36.600000	0.626250	41.000000	1.000000
max	67.100000	2.420000	81.000000	1.000000

```
[2]: data.skew()
```

```
[2]: Pregnancies      0.901674
      Glucose          0.173754
      BloodPressure    -1.843608
      SkinThickness    0.109372
      Insulin          2.272251
      BMI              -0.428982
      DiabetesPedigreeFunction  1.919911
      Age              1.129597
      Outcome          0.635017
      dtype: float64
```

```
[3]: data.kurt()
```

```
[3]: Pregnancies      0.159220
      Glucose          0.640780
      BloodPressure    5.180157
      SkinThickness    -0.520072
      Insulin          7.214260
      BMI              3.290443
      DiabetesPedigreeFunction  5.594954
      Age              0.643159
      Outcome          -1.600930
      dtype: float64
```

```
[4]: data.mode().iloc[0]
```

```
[4]: Pregnancies      1.000
      Glucose          99.000
      BloodPressure    70.000
      SkinThickness    0.000
      Insulin          0.000
      BMI              32.000
      DiabetesPedigreeFunction  0.254
      Age              22.000
      Outcome          0.000
      Name: 0, dtype: float64
```

```
[5]: X = data.drop('Outcome', axis=1)
      y = data['Outcome']
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
      ↪random_state=42)
```

```
[6]: linear_reg = LinearRegression()
      linear_reg.fit(X_train, y_train)
      y_pred_linear = linear_reg.predict(X_test)
      r2_linear = r2_score(y_test, y_pred_linear)
```

```
print(f"Linear Regression R-squared: {r2_linear}")

# Bivariate analysis - Logistic regression
logistic_reg = LogisticRegression()
logistic_reg.fit(X_train, y_train)
y_pred_logistic = logistic_reg.predict(X_test)
accuracy = accuracy_score(y_test, y_pred_logistic)
print(f"Logistic Regression Accuracy: {accuracy}")
```

Linear Regression R-squared: 0.255002811767418
Logistic Regression Accuracy: 0.7467532467532467

[]: