

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
from sklearn.metrics import confusion_matrix

df = pd.read_csv("/content/Diabetes.csv")

df.describe()

{"summary": "{\n  \"name\": \"df\",\n  \"rows\": 8,\n  \"fields\": [\n    {\n      \"column\": \"Number of times pregnant\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 269.85223453356366,\n        \"min\": 0.0,\n        \"max\": 768.0,\n        \"num_unique_values\": 8,\n        \"samples\": [\n          3.8450520833333335,\n          3.0,\n          768.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"Plasma glucose concentration a 2 hours in an oral glucose tolerance test\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 243.73802348295857,\n        \"min\": 0.0,\n        \"max\": 768.0,\n        \"num_unique_values\": 8,\n        \"samples\": [\n          120.89453125,\n          117.0,\n          768.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"Diastolic blood pressure (mm Hg)\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 252.8525053581062,\n        \"min\": 0.0,\n        \"max\": 768.0,\n        \"num_unique_values\": 8,\n        \"samples\": [\n          69.10546875,\n          72.0,\n          768.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"Triceps skin fold thickness (mm)\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 263.7684730531098,\n        \"min\": 0.0,\n        \"max\": 768.0,\n        \"num_unique_values\": 7,\n        \"samples\": [\n          768.0,\n          20.536458333333332,\n          32.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"2-Hour serum insulin (mu U/ml)\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 350.26059167945886,\n        \"min\": 0.0,\n        \"max\": 846.0,\n        \"num_unique_values\": 7,\n        \"samples\": [\n          79.79947916666667,\n          127.25,\n          768.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"Body mass index (weight in kg/(height in m)^2)\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 262.05117817552093,\n        \"min\": 0.0,\n        \"max\": 768.0,\n        \"num_unique_values\": 8,\n        \"samples\": [\n          31.992578124999998,\n          32.0,\n          768.0\n        ]\n      }\n    }\n  ]\n}"

```

```

768.0\n          ],\n          \"semantic_type\": \"\",\n          \"description\": \"\"\n        },\n        {\n          \"column\":\n          \"Diabetes pedigree function\",\n          \"properties\": {\n            \"dtype\": \"number\",\n            \"std\": 271.3005221658502,\n            \"min\": 0.078,\n            \"max\": 768.0,\n            \"num_unique_values\": 8,\n            \"samples\": [\n              0.47187630208333325,\n              0.3725,\n              768.0\n            ],\n            \"semantic_type\": \"\",\n            \"description\": \"\"\n          },\n          {\n            \"column\": \"Age (years)\",\n            \"properties\": {\n              \"dtype\": \"number\",\n              \"std\":\n              260.1941178528413,\n              \"min\": 11.760231540678685,\n              \"max\": 768.0,\n              \"num_unique_values\": 8,\n              \"samples\": [\n                33.240885416666664,\n                29.0,\n                768.0\n              ],\n              \"semantic_type\": \"\",\n              \"description\": \"\"\n            },\n            {\n              \"column\":\n              \"Class variable (0 or 1)\",\n              \"properties\": {\n                \"dtype\": \"number\",\n                \"std\": 271.3865920388932,\n                \"min\": 0.0,\n                \"max\": 768.0,\n                \"num_unique_values\":\n                5,\n                \"samples\": [\n                  0.3489583333333333,\n                  1.0,\n                  0.47695137724279896\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n              }\n            }\n          }\n        ],\n        \"type\": \"dataframe\"}

```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 768 entries, 0 to 767
```

```
Data columns (total 9 columns):
```

```
#    Column
```

```
Non-Null Count  Dtype
```

```
---  ---
```

```
-----
```

```
0    Number of times pregnant
```

```
768 non-null    int64
```

```
1    Plasma glucose concentration a 2 hours in an oral glucose  
tolerance test  768 non-null    int64
```

```
2    Diastolic blood pressure (mm Hg)
```

```
768 non-null    int64
```

```
3    Triceps skin fold thickness (mm)
```

```
768 non-null    int64
```

```
4    2-Hour serum insulin (mu U/ml)
```

```
768 non-null    int64
```

```
5    Body mass index (weight in kg/(height in m)^2)
```

```
768 non-null    float64
```

```
6    Diabetes pedigree function
```

```
768 non-null    float64
```

```
7    Age (years)
```

```
768 non-null    int64
```

```
8    Class variable (0 or 1)
```

```
768 non-null    int64
```

```
dtypes: float64(2), int64(7)
```

```
memory usage: 54.1 KB
```

```
df.head()
```

```
{
  "summary": {
    "name": "df",
    "rows": 768,
    "fields": [
      {
        "column": "Number of times pregnant",
        "properties": {
          "dtype": "number",
          "std": 3,
          "min": 0,
          "max": 17,
          "num_unique_values": 17,
          "samples": [6, 1, 3],
          "semantic_type": ""
        },
        "description": "Plasma glucose concentration a 2 hours in an oral glucose tolerance test",
        "properties": {
          "dtype": "number",
          "std": 31,
          "min": 0,
          "max": 199,
          "num_unique_values": 136,
          "samples": [151, 101, 112],
          "semantic_type": ""
        },
        "description": "Diastolic blood pressure (mm Hg)",
        "properties": {
          "dtype": "number",
          "std": 19,
          "min": 0,
          "max": 122,
          "num_unique_values": 47,
          "samples": [86, 46, 85],
          "semantic_type": ""
        },
        "description": "Triceps skin fold thickness (mm)",
        "properties": {
          "dtype": "number",
          "std": 15,
          "min": 0,
          "max": 99,
          "num_unique_values": 51,
          "samples": [7, 12, 48],
          "semantic_type": ""
        },
        "description": "2-Hour serum insulin (mu U/ml)",
        "properties": {
          "dtype": "number",
          "std": 115,
          "min": 0,
          "max": 846,
          "num_unique_values": 186,
          "samples": [52, 41, 183],
          "semantic_type": ""
        },
        "description": "Body mass index (weight in kg/(height in m)^2)",
        "properties": {
          "dtype": "number",
          "std": 7.884160320375446,
          "min": 0.0,
          "max": 67.1,
          "num_unique_values": 248,
          "samples": [19.9, 31.0, 38.1],
          "semantic_type": ""
        },
        "description": "Diabetes pedigree function",
        "properties": {
          "dtype": "number",
          "std": 0.3313285950127749,
          "min": 0.078,
          "max": 2.42,
          "num_unique_values": 517,
          "samples": [1.731, 0.426, 0.138],
          "semantic_type": ""
        },
        "description": "Age (years)",
        "properties": {
          "dtype": "number",
          "std": 11,
          "min": 21,
          "max": 81
        }
      }
    ]
  }
}
```



```

0.349,\n          \"num_unique_values\": 5,\n          \"samples\": [\n
0.34,\n          0.315,\n          0.245\n          ],\n
\"semantic_type\": \"\", \n          \"description\": \"\"\n          }\n
n      },\n      {\n          \"column\": \"Age (years)\",\n
\"properties\": {\n          \"dtype\": \"number\", \n          \"std\":\n
16,\n          \"min\": 23,\n          \"max\": 63,\n
\"num_unique_values\": 5,\n          \"samples\": [\n          27,\n
23,\n          30\n          ],\n          \"semantic_type\": \"\", \n
\"description\": \"\"\n          }\n      },\n      {\n          \"column\":\n
\"Class variable (0 or 1)\",\n          \"properties\": {\n
\"dtype\": \"number\", \n          \"std\": 0,\n          \"min\": 0,\n
\"max\": 1,\n          \"num_unique_values\": 2,\n          \"samples\":\n
[\n          1,\n          0\n          ],\n          \"semantic_type\":\n
\"\", \n          \"description\": \"\"\n          }\n      }\n      ]\n
n}], \"type\": \"dataframe\"}

```

```

num_instances, num_features = df.shape
print(f\"Number of instances: {num_instances}\")
print(f\"Number of features: {num_features}\")

```

```
print(df.dtypes)
```

```

Number of instances: 768
Number of features: 9
Number of times pregnant
int64
Plasma glucose concentration a 2 hours in an oral glucose tolerance
test      int64
Diastolic blood pressure (mm Hg)
int64
Triceps skin fold thickness (mm)
int64
2-Hour serum insulin (mu U/ml)
int64
Body mass index (weight in kg/(height in m)^2)
float64
Diabetes pedigree function
float64
Age (years)
int64
Class variable (0 or 1)
int64
dtype: object

df['Class variable (0 or 1)'].value_counts()

```

```
0    500
```

```
1    268
```

```
Name: Class variable (0 or 1), dtype: int64
```

```
X = df.drop('Class variable (0 or 1)', axis=1)
```

```
y = df['Class variable (0 or 1)']
```

```
0    1
```

```
1    0
```

```
2    1
```

```
3    0
```

```
4    1
```

```
..
```

```
763  0
```

```
764  0
```

```
765  0
```

```
766  1
```

```
767  0
```

```
Name: Class variable (0 or 1), Length: 768, dtype: int64
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
features = ['Number of times pregnant', 'Plasma glucose concentration  
a 2 hours in an oral glucose tolerance test', 'Diastolic blood  
pressure (mm Hg)', 'Triceps skin fold thickness (mm)', '2-Hour serum  
insulin (mu U/ml)', 'Body mass index (weight in kg/(height in m)^2)',  
'Diabetes pedigree function', 'Age (years)']  
label = 'Class variable (0 or 1)'
```

```
for col in features:
```

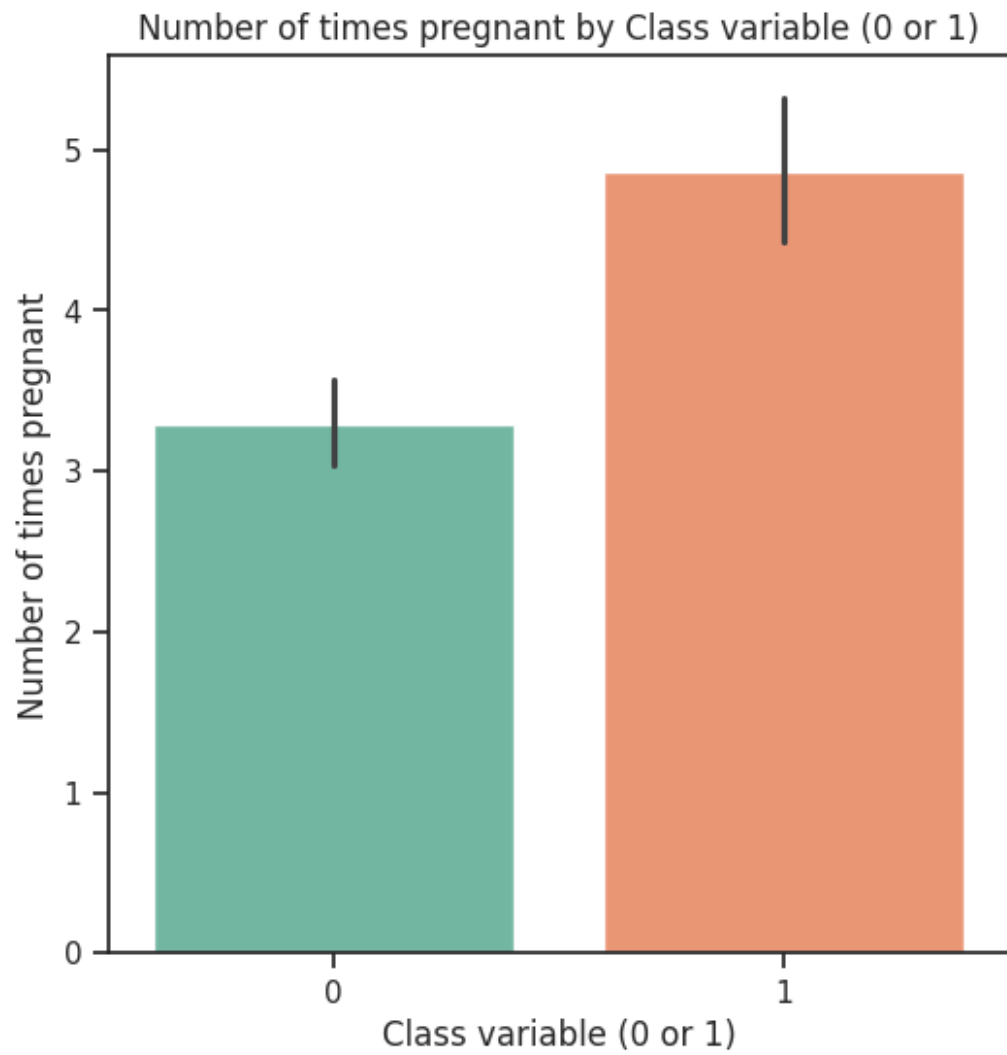
```
    plt.figure(figsize=(6, 6))
```

```
    sns.barplot(x=label, y=col, hue=label, data=df, palette='Set2',
```

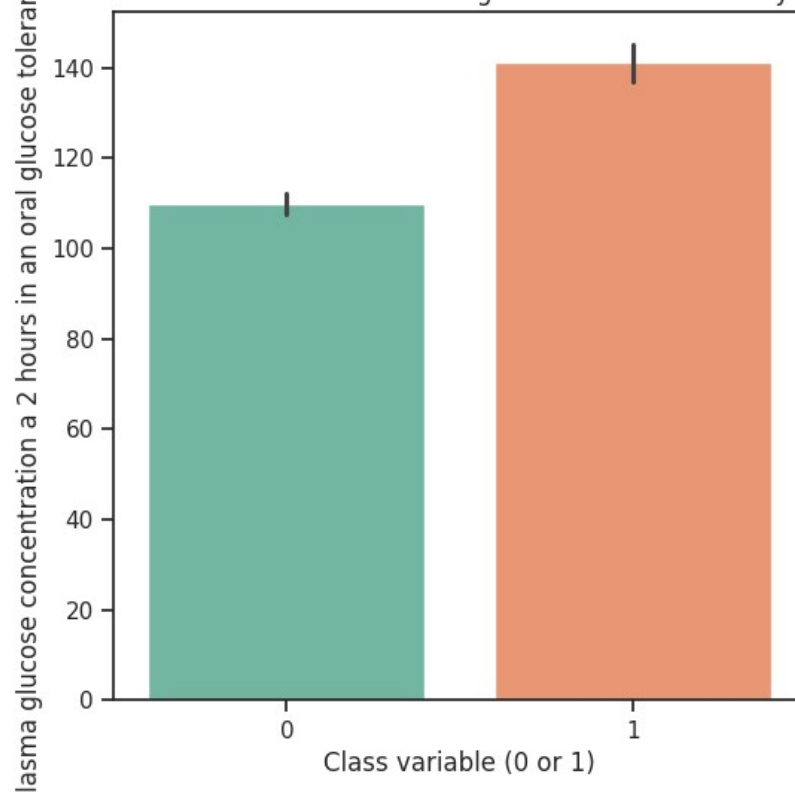
```
    legend=False)
```

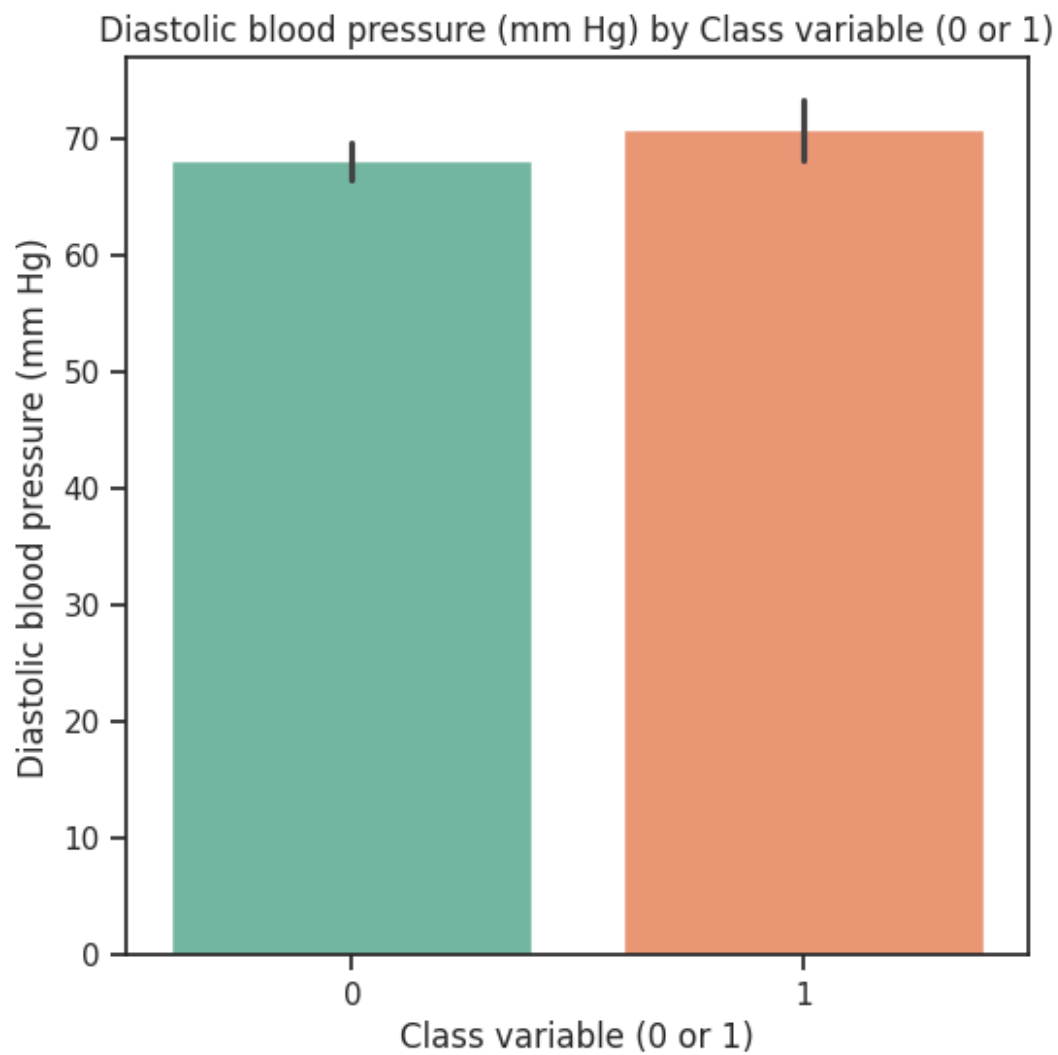
```
    plt.title(f'{col} by {label}')
```

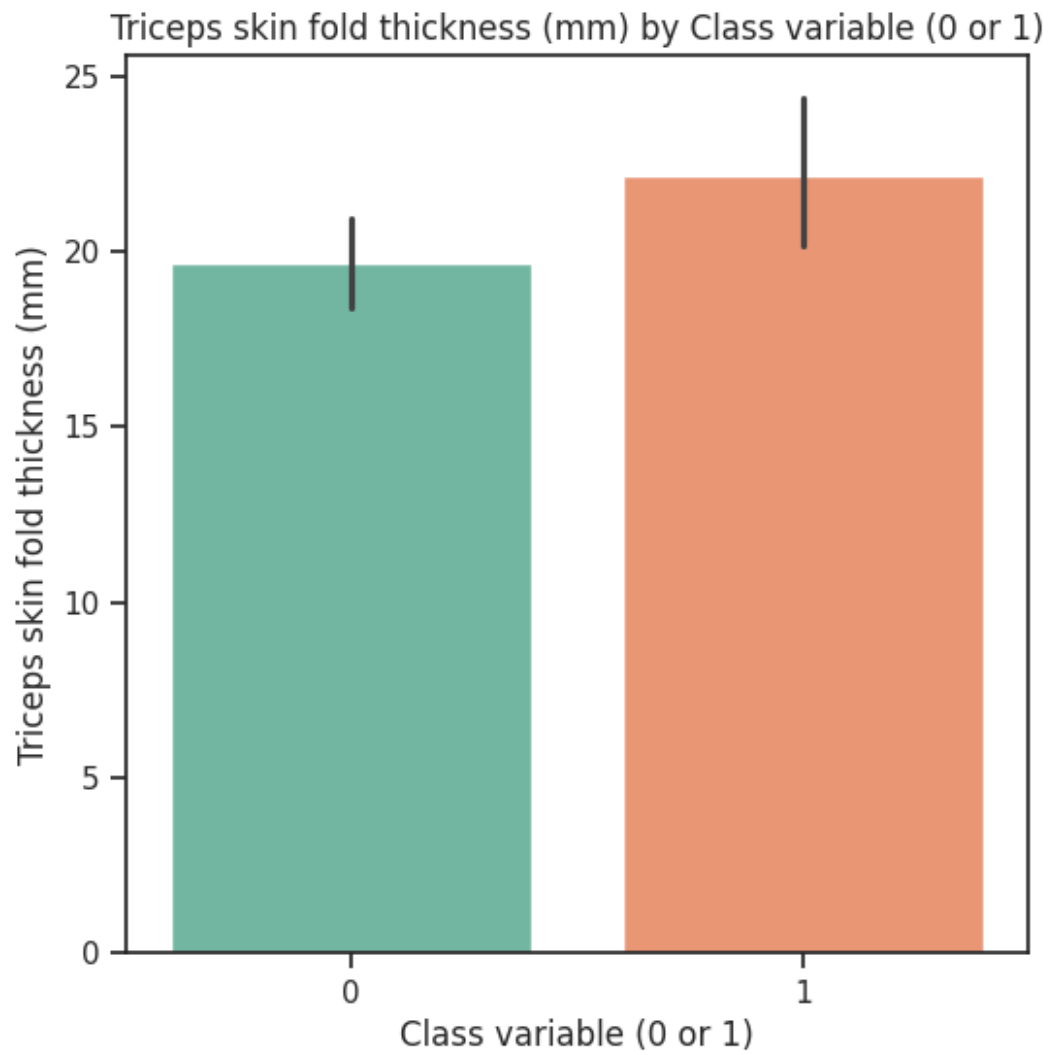
```
    plt.show()
```

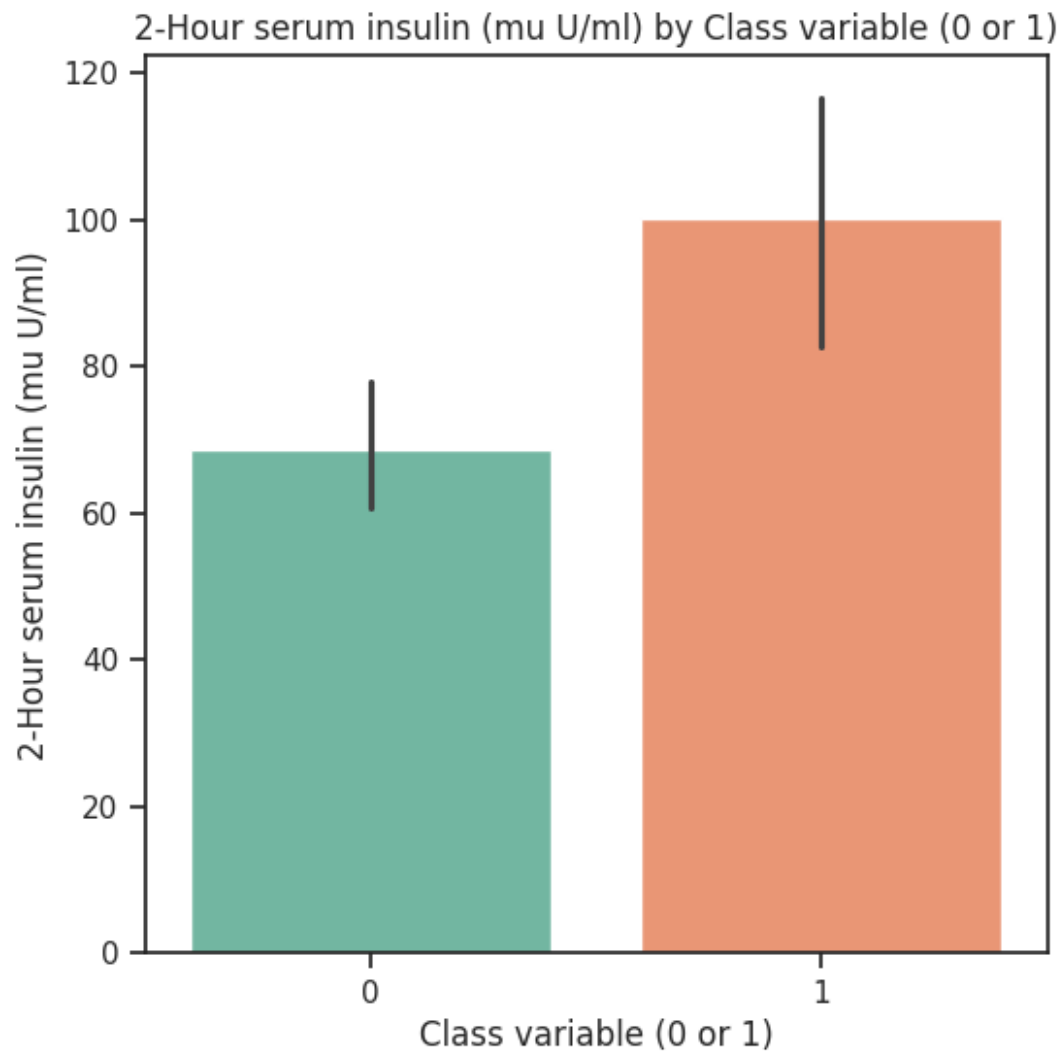


Plasma glucose concentration a 2 hours in an oral glucose tolerance test by Class variable (0 or 1)

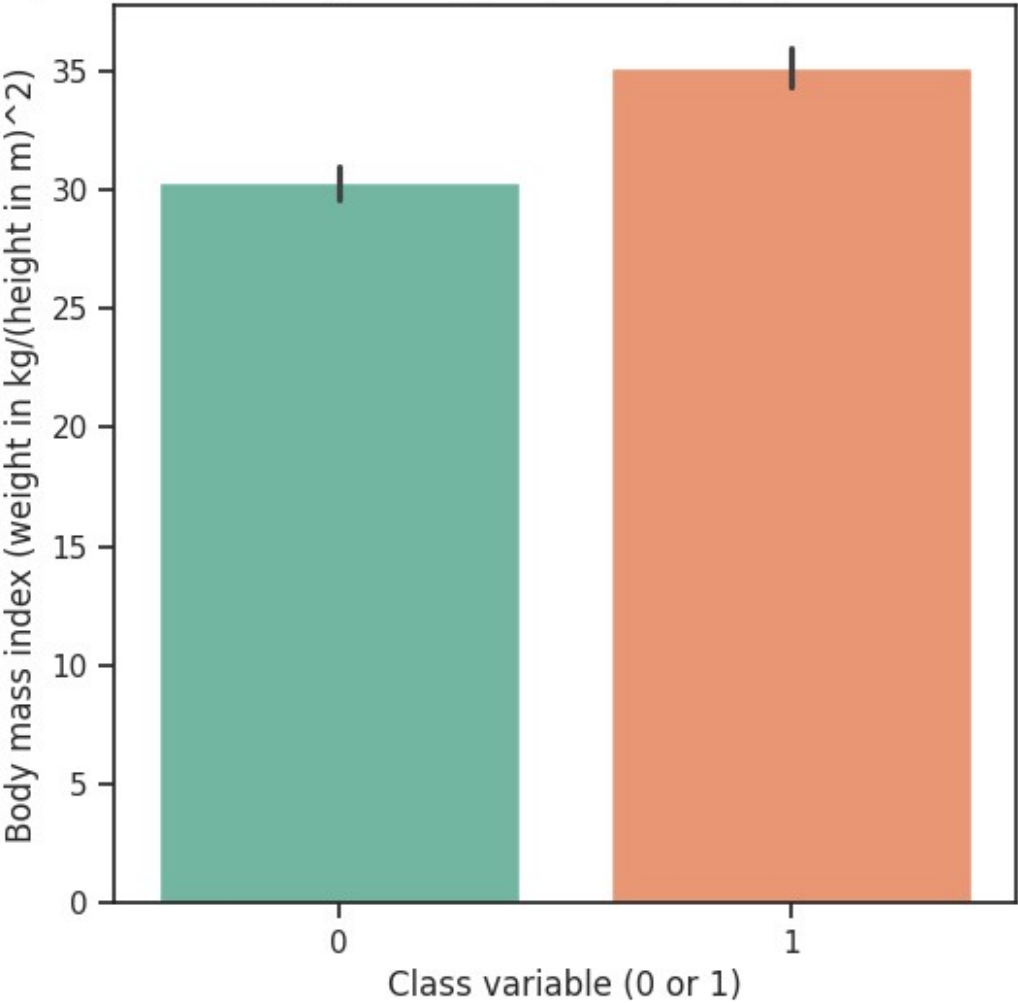


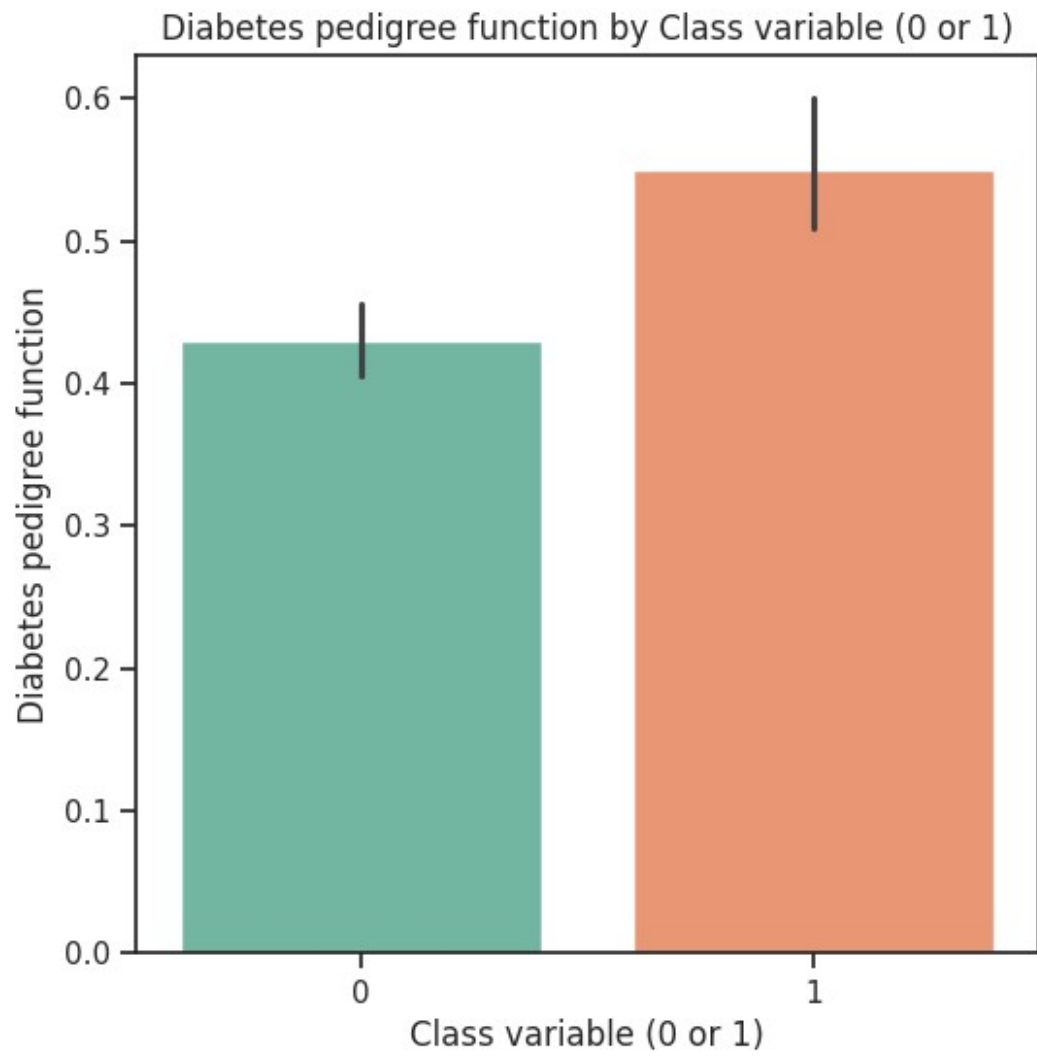


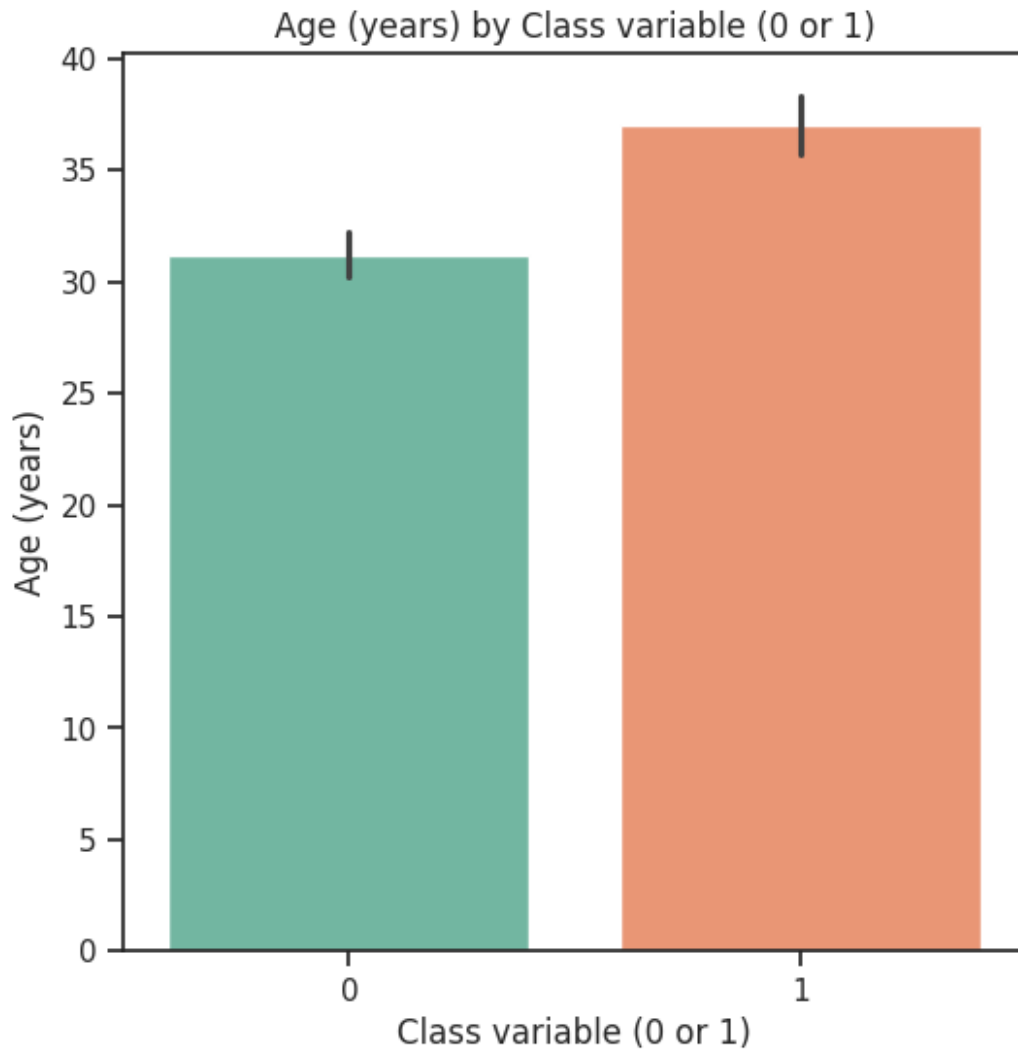




Body mass index (weight in kg/(height in m)^2) by Class variable (0 or 1)







```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

model = LogisticRegression(max_iter = 1000)
model.fit(X_train, y_train)
LogisticRegression(max_iter=1000)

y_pred = model.predict(X_test)

cm = confusion_matrix(y_test,y_pred)
print(cm)

accuracy = accuracy_score(y_test, y_pred)
print('Accuracy of the binary classifier = {:.3f}'.format(accuracy))
```

```
print(classification_report(y_test, y_pred))
```

```
[[78 21]  
 [18 37]]
```

Accuracy of the binary classifier = 0.747

	precision	recall	f1-score	support
0	0.81	0.79	0.80	99
1	0.64	0.67	0.65	55
accuracy			0.75	154
macro avg	0.73	0.73	0.73	154
weighted avg	0.75	0.75	0.75	154