# Analysis

July 3, 2023

```
[51]: import pandas as pd
      import seaborn as sns
      import matplotlib.pyplot as plt
      import statsmodels.api as sm
      import statsmodels.formula.api as smf
      import numpy as np
      from sklearn import metrics
```

```
[52]: df = pd.read_csv('diabetes-dataset.csv')
      df = df.rename(columns={'Outcome': 'DiabetesOutcome'})
```
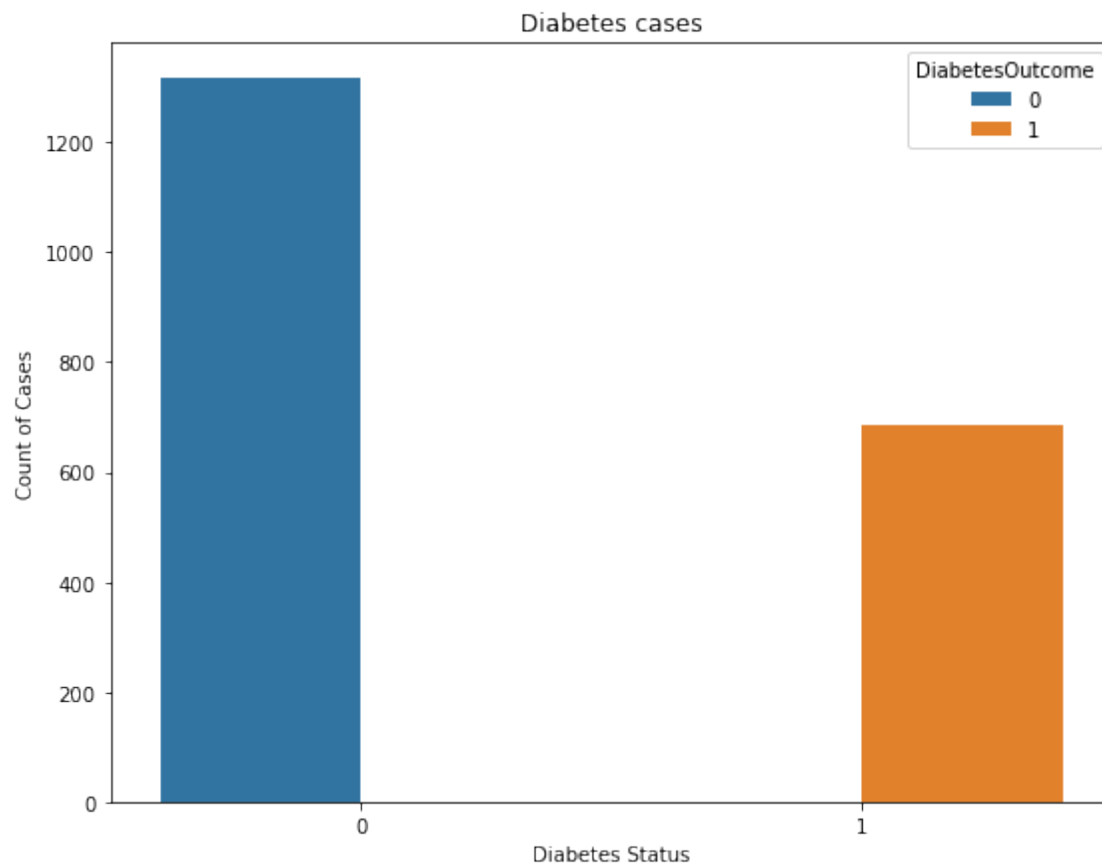
```
[53]: df.describe()
```

```
[53]:        Pregnancies      Glucose  BloodPressure  SkinThickness      Insulin
      count  2000.000000  2000.000000    2000.000000    2000.000000  2000.000000  \
      mean      3.703500   121.182500      69.145500      20.935000    80.254000
      std       3.306063    32.068636      19.188315      16.103243   111.180534
      min       0.000000     0.000000       0.000000       0.000000     0.000000
      25%       1.000000    99.000000      63.500000       0.000000     0.000000
      50%       3.000000   117.000000      72.000000      23.000000    40.000000
      75%       6.000000   141.000000      80.000000      32.000000   130.000000
      max      17.000000   199.000000     122.000000     110.000000   744.000000

                     BMI  DiabetesPedigreeFunction          Age  DiabetesOutcome
      count  2000.000000               2000.000000  2000.000000      2000.000000
      mean     32.193000                  0.470930    33.090500         0.342000
      std       8.149901                  0.323553    11.786423         0.474498
      min       0.000000                  0.078000    21.000000         0.000000
      25%      27.375000                  0.244000    24.000000         0.000000
      50%      32.300000                  0.376000    29.000000         0.000000
      75%      36.800000                  0.624000    40.000000         1.000000
      max      80.600000                  2.420000    81.000000         1.000000
```
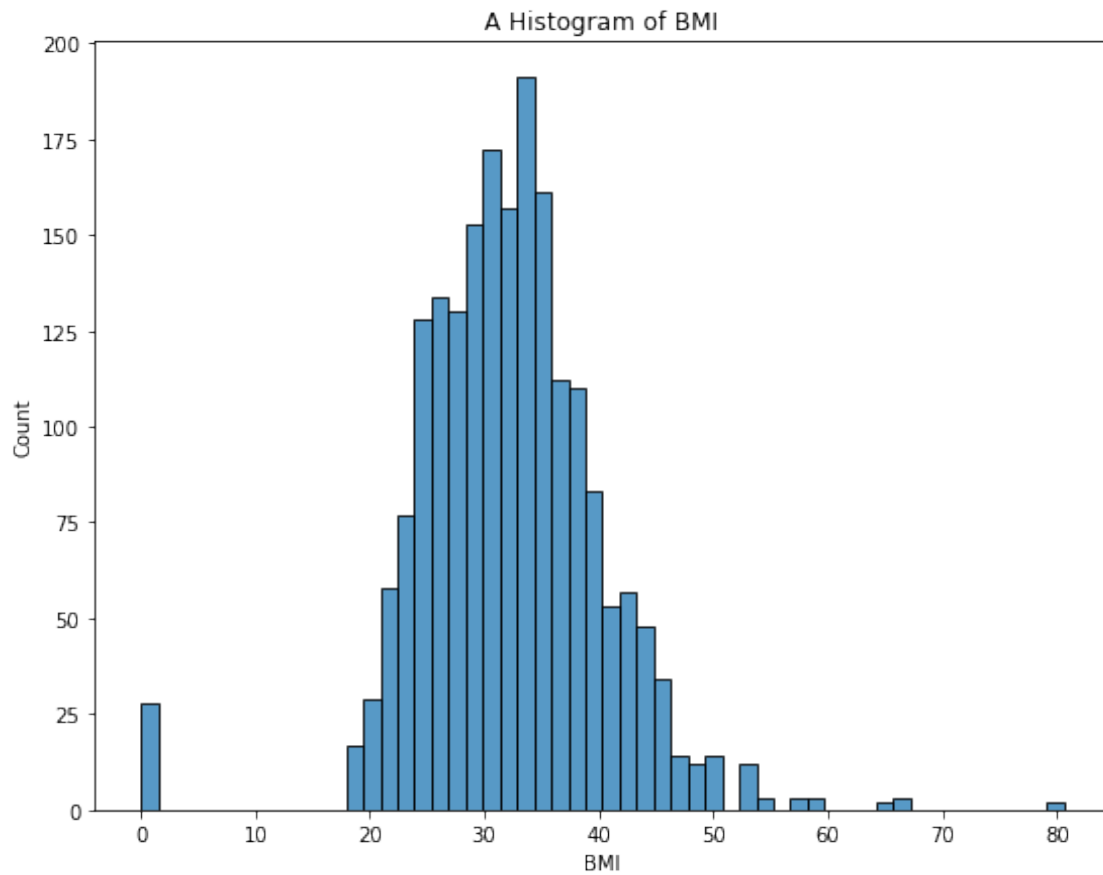
```
[54]: plt.figure(figsize=(9,7))
      sns.countplot(x='DiabetesOutcome', hue='DiabetesOutcome', data=df).
       ↪set(title='Diabetes cases', ylabel='Count of Cases', xlabel='Diabetes␣
       ↪Status')
```
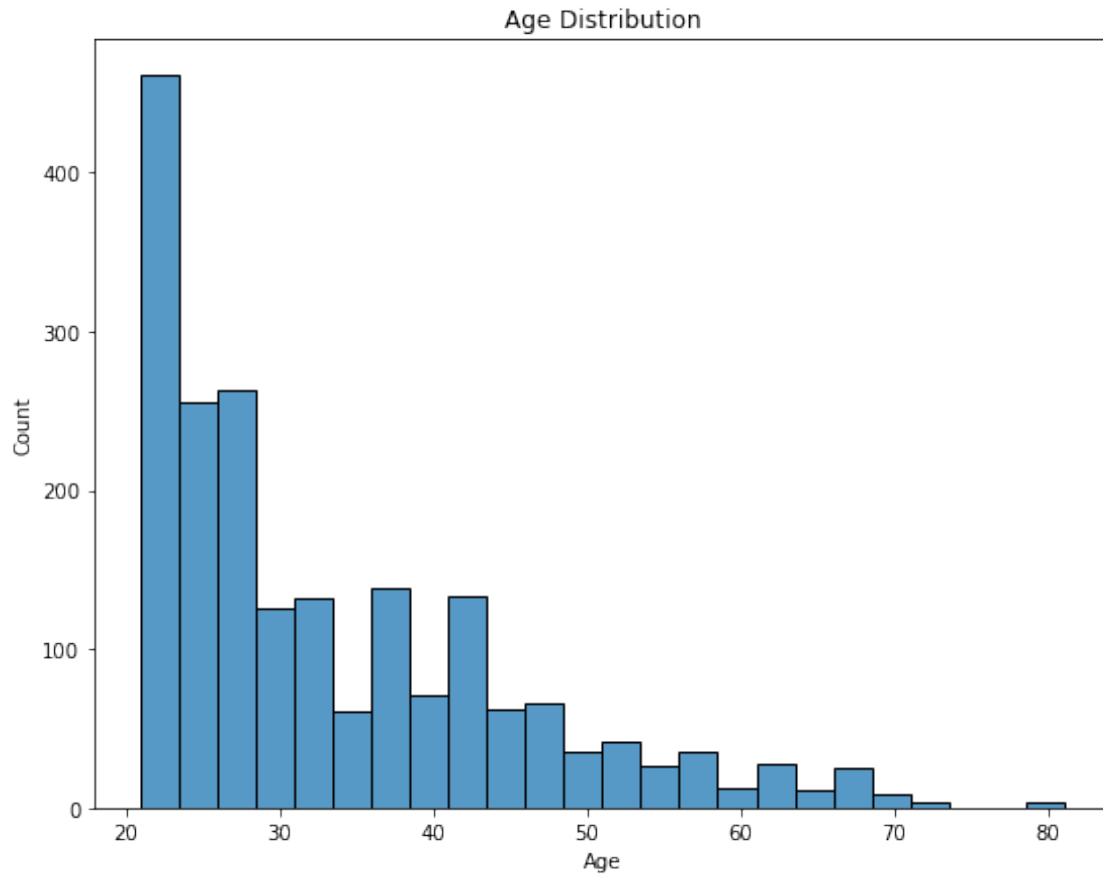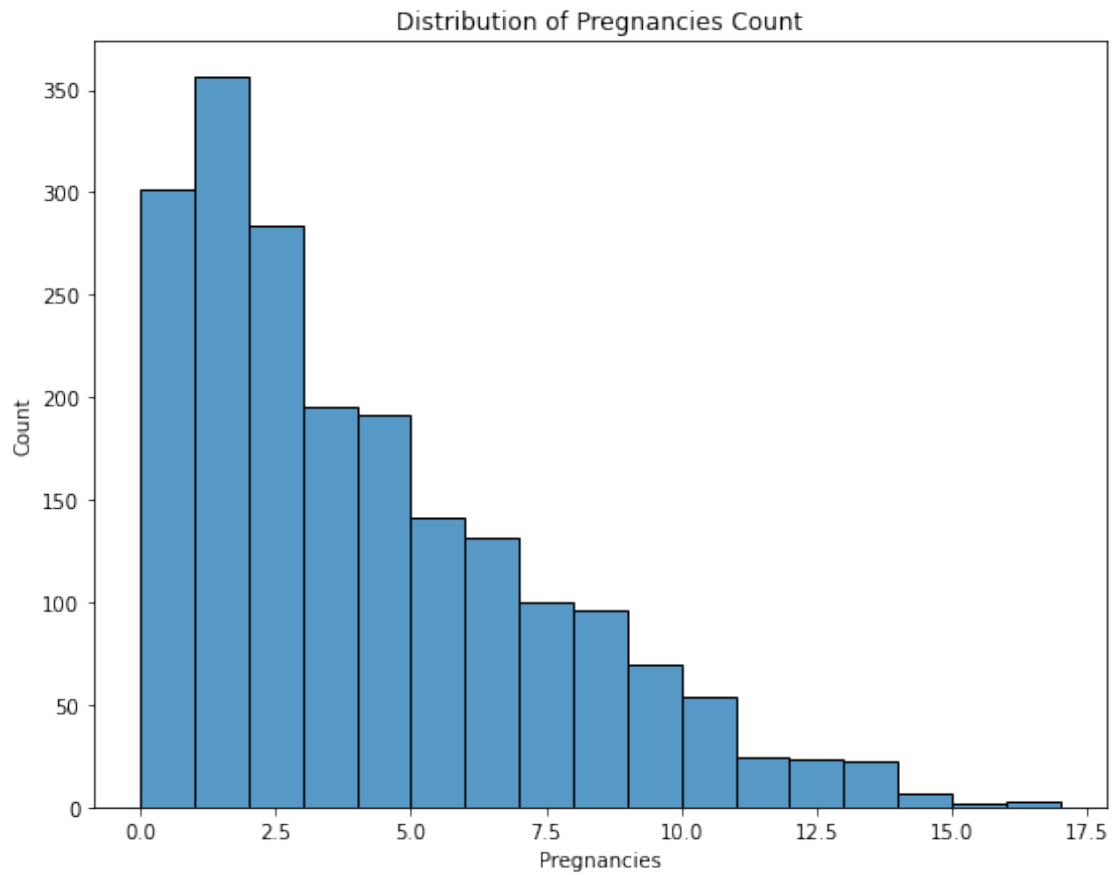
```
plt.show()
```



```
[55]: plt.figure(figsize=(9,7))
      sns.histplot(x='BMI', data=df).set(title='A Histogram of BMI', ylabel='Count',␣
       ↪xlabel='BMI')
      plt.show()
```

A Histogram of BMI

```
[56]: plt.figure(figsize=(9,7))
      sns.histplot(x='Age', data=df).set(title='Age Distribution', ylabel='Count',␣
       ↪xlabel='Age')
      plt.show()
```
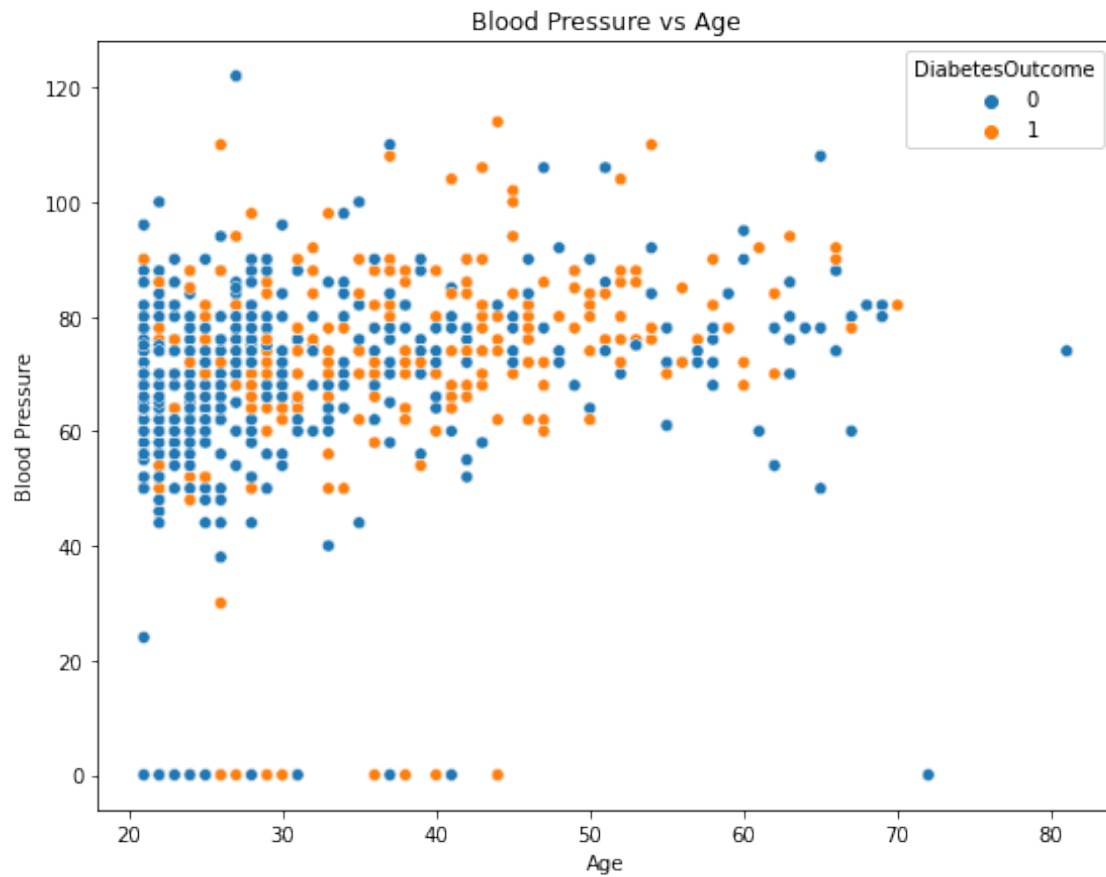
Age Distribution

```
[57]: plt.figure(figsize=(9,7))
      sns.histplot(x='Pregnancies', binwidth=1, data=df).set(title='Distribution of␣
      ↪Pregnancies Count', ylabel='Count', xlabel='Pregnancies')
      plt.show()
```
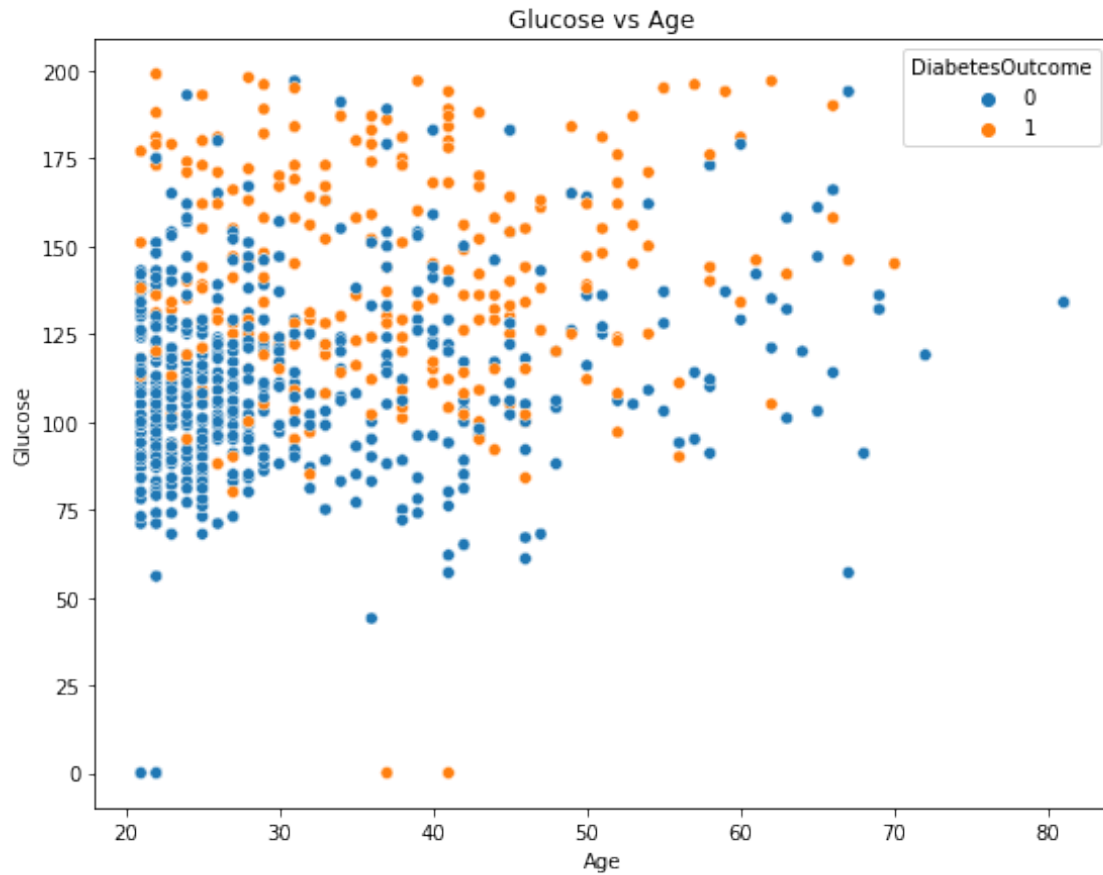
Distribution of Pregnancies Count

[58]:
```
corr = df.corr()
plt.figure(figsize=(9,7))
sns.heatmap(corr, xticklabels=corr.columns.values, yticklabels=corr.columns.
 ↪values, cmap="Greens", annot=True).set(title='Correlation Plot')
plt.show()
```

## Correlation Plot



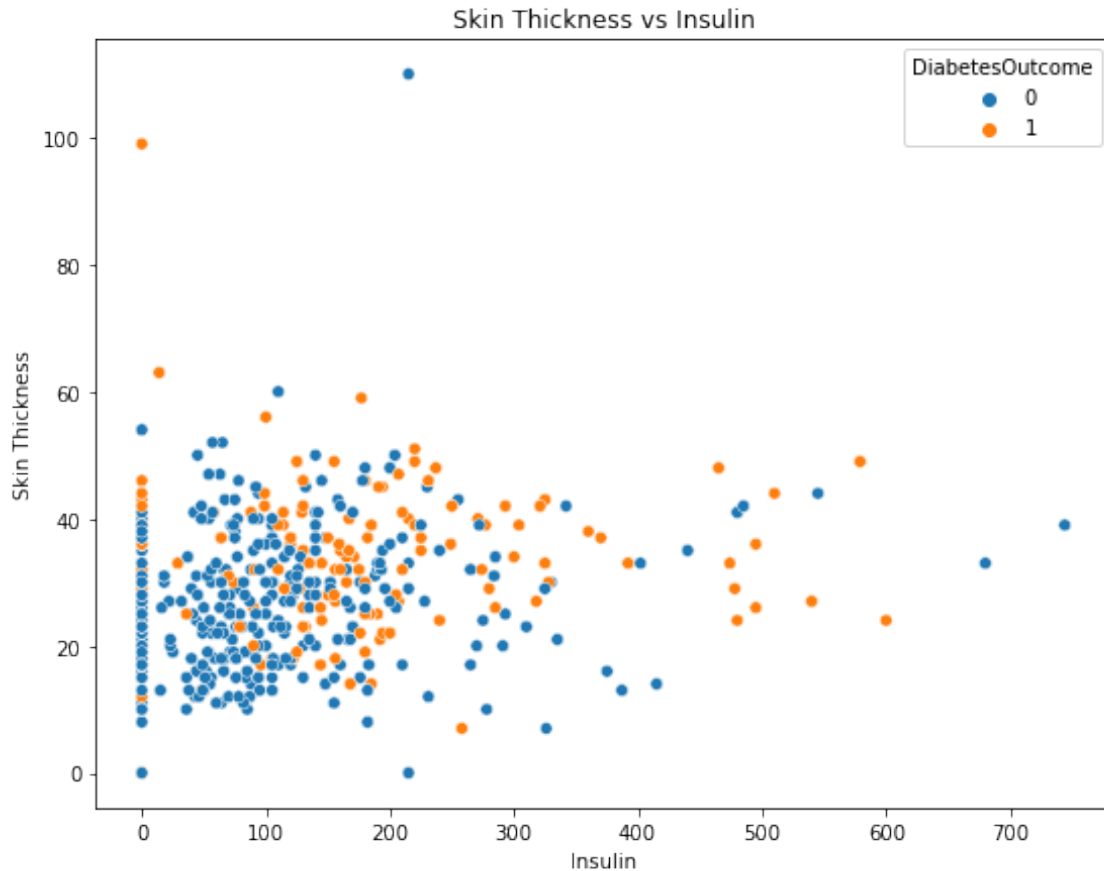|  | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | DiabetesOutcome |
|---|---|---|---|---|---|---|---|---|---|
| Pregnancies | 1 | 0.12 | 0.15 | -0.063 | -0.077 | 0.019 | -0.025 | 0.54 | 0.22 |
| Glucose | 0.12 | 1 | 0.14 | 0.062 | 0.32 | 0.23 | 0.12 | 0.25 | 0.46 |
| BloodPressure | 0.15 | 0.14 | 1 | 0.2 | 0.087 | 0.28 | 0.051 | 0.24 | 0.076 |
| SkinThickness | -0.063 | 0.062 | 0.2 | 1 | 0.45 | 0.39 | 0.18 | -0.11 | 0.076 |
| Insulin | -0.077 | 0.32 | 0.087 | 0.45 | 1 | 0.22 | 0.19 | -0.086 | 0.12 |
| BMI | 0.019 | 0.23 | 0.28 | 0.39 | 0.22 | 1 | 0.13 | 0.039 | 0.28 |
| DiabetesPedigreeFunction | -0.025 | 0.12 | 0.051 | 0.18 | 0.19 | 0.13 | 1 | 0.027 | 0.16 |
| Age | 0.54 | 0.25 | 0.24 | -0.11 | -0.086 | 0.039 | 0.027 | 1 | 0.24 |
| DiabetesOutcome | 0.22 | 0.46 | 0.076 | 0.076 | 0.12 | 0.28 | 0.16 | 0.24 | 1 |

```
[59]: plt.figure(figsize=(9,7))
      sns.scatterplot(df, y='BloodPressure', x='Age', hue='DiabetesOutcome').
        ↪set(title='Blood Pressure vs Age', ylabel='Blood Pressure', xlabel='Age')
      plt.show()
```

Blood Pressure vs Age

```
[60]: plt.figure(figsize=(9,7))
      sns.scatterplot(df, y='Glucose', x='Age', hue='DiabetesOutcome').
       ↪set(title='Glucose vs Age', ylabel='Glucose', xlabel='Age')
      plt.show()
```

Glucose vs Age

```
[61]: plt.figure(figsize=(9,7))
      sns.scatterplot(df, y='SkinThickness', x='Insulin', hue='DiabetesOutcome').
       ↪set(title='Skin Thickness vs Insulin', ylabel='Skin Thickness',␣
       ↪xlabel='Insulin')
      plt.show()
```

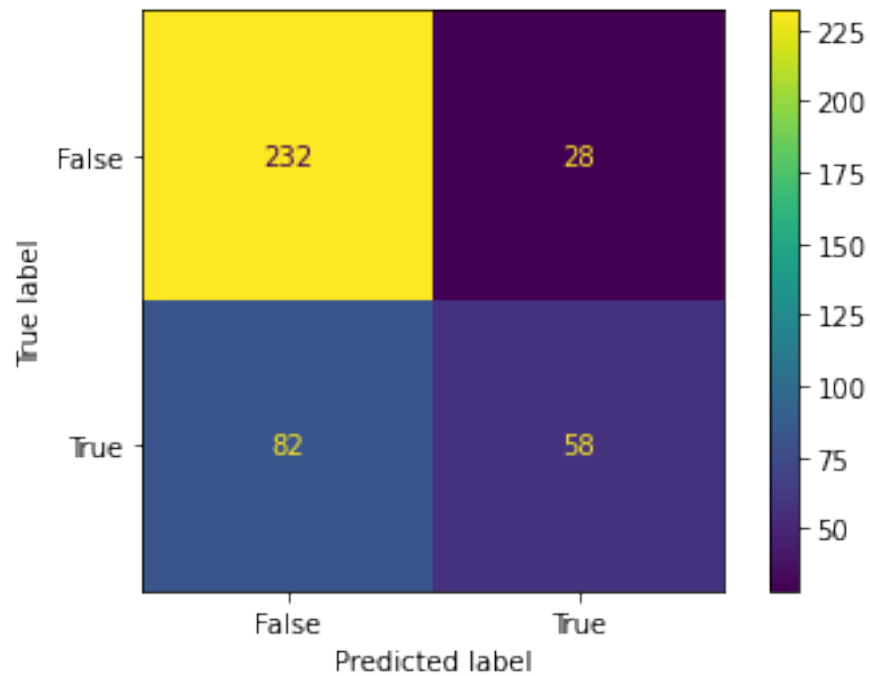Skin Thickness vs Insulin

```
[62]: df_train = df.sample(round(len(df)*0.8))
      df_test = df.drop(df_train.index)
```

```
[63]: formula = 'DiabetesOutcome ~␣
      ↪Pregnancies+Glucose+BloodPressure+SkinThickness+Insulin+BMI+DiabetesPedigreeFunction+Age'
```

```
[64]: model = smf.glm(formula=formula, data=df_train, family=sm.families.Poisson())
      result = model.fit()
```

```
[65]: result = result.predict(df_test)
      result[result > 0.5] = 1
      result[result <= 0.5] = 0
```

```
[66]: confusion_matrix = metrics.confusion_matrix(df_test.DiabetesOutcome, result)
      cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix =␣
      ↪confusion_matrix, display_labels = [False, True])
      cm_display.plot()
      plt.show()
```

```
[ ]:
```

```
[67]: accuracy = metrics.accuracy_score(df_test.DiabetesOutcome, result)
      print(accuracy)
```

```
0.725
```

```
[ ]:
```