```python
In [3]:  import pandas as pd
         import matplotlib.pyplot as plt
         import numpy as np

         %matplotlib inline
```

```python
In [5]:  data = pd.read_csv("desktop/data/diabetes.csv")
```

```python
In [6]:  data.shape
```

```
Out[6]:  (768, 9)
```

```python
In [7]:  data.head(5)
```

Out[7]:

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

```python
In [10]: # checking if any null value is present
         data.isnull().values.any()
```

```
Out[10]: False
```

```python
In [15]: import seaborn as sns
         import matplotlib.pyplot as plt
         corrmat = data.corr()
         top_corr_features = corrmat.index
         plt.figure(figsize=(20,20))
         # plotting heatmap
         g=sns.heatmap(data[top_corr_features].corr(),annot=True,cmap="RdYlGn")
```



```python
In [16]: data.corr()
```

Out[16]:

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| Pregnancies | 1.000000 | 0.129459 | 0.141282 | -0.081672 | -0.073535 | 0.017683 | -0.033523 | 0.544341 | 0.221898 |
| Glucose | 0.129459 | 1.000000 | 0.152590 | 0.057328 | 0.331357 | 0.221071 | 0.137337 | 0.263514 | 0.466581 |
| BloodPressure | 0.141282 | 0.152590 | 1.000000 | 0.207371 | 0.088933 | 0.281805 | 0.041265 | 0.239528 | 0.065068 |
| SkinThickness | -0.081672 | 0.057328 | 0.207371 | 1.000000 | 0.436783 | 0.392573 | 0.183928 | -0.113970 | 0.074752 |
| Insulin | -0.073535 | 0.331357 | 0.088933 | 0.436783 | 1.000000 | 0.197859 | 0.185071 | -0.042163 | 0.130548 |
| BMI | 0.017683 | 0.221071 | 0.281805 | 0.392573 | 0.197859 | 1.000000 | 0.140647 | 0.036242 | 0.292695 |
| DiabetesPedigreeFunction | -0.033523 | 0.137337 | 0.041265 | 0.183928 | 0.185071 | 0.140647 | 1.000000 | 0.033561 | 0.173844 |
| Age | 0.544341 | 0.263514 | 0.239528 | -0.113970 | -0.042163 | 0.036242 | 0.033561 | 1.000000 | 0.238356 |
| Outcome | 0.221898 | 0.466581 | 0.065068 | 0.074752 | 0.130548 | 0.292695 | 0.173844 | 0.238356 | 1.000000 |

```python
In [25]: diabetes_true_count = len(data.loc[data['Outcome'] == True])
         diabetes_false_count = len(data.loc[data['Outcome'] == False])
```

```python
In [26]: (diabetes_true_count,diabetes_false_count)
```

```
Out[26]: (268, 500)
```

```python
In [27]: #training test split
         from sklearn.model_selection import train_test_split
         features_columns = ['Pregnancies','Glucose','BloodPressure','SkinThickness','Insulin','BMI','DiabetesPedigreeFunctio
         n','Age']
         predicted_class = ['Outcome']
```

```python
In [47]: X = data[features_columns].values
         y = data[predicted_class].values

         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30, random_state=10)
```

```python
In [48]: print("total number of rows : {0}".format(len(data)))
         print("number of rows missing Glucose: {0}".format(len(data.loc[data['Glucose'] == 0])))
         print("number of rows missing BloodPressure: {0}".format(len(data.loc[data['BloodPressure'] == 0])))
         print("number of rows missing SkinThickness: {0}".format(len(data.loc[data['SkinThickness'] == 0])))
         print("number of rows missing insulin: {0}".format(len(data.loc[data['Insulin'] == 0])))
         print("number of rows missing BMI: {0}".format(len(data.loc[data['BMI'] == 0])))
         print("number of rows missing DiabetesPedigreeFunction: {0}".format(len(data.loc[data['DiabetesPedigreeFunction'] ==
         0])))
         print("number of rows missing Age: {0}".format(len(data.loc[data['Age'] == 0])))

         total number of rows : 768
         number of rows missing Glucose: 5
         number of rows missing BloodPressure: 35
         number of rows missing SkinThickness: 227
         number of rows missing insulin: 374
         number of rows missing BMI: 11
         number of rows missing DiabetesPedigreeFunction: 0
         number of rows missing Age: 0
```

```python
In [49]: from sklearn.impute import SimpleImputer
         fill_values = SimpleImputer(missing_values=0, strategy="mean")
         X_train = fill_values.fit_transform(X_train)
         X_test = fill_values.fit_transform(X_test)
```

```python
In [53]: ## Apply Algorithm
         from sklearn.ensemble import RandomForestClassifier
         random_forest_model = RandomForestClassifier(random_state=10)

         random_forest_model.fit(X_train, y_train.ravel())
```

```
Out[53]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                                criterion='gini', max_depth=None, max_features='auto',
                                max_leaf_nodes=None, max_samples=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=100,
                                n_jobs=None, oob_score=False, random_state=10, verbose=0,
                                warm_start=False)
```

```python
In [58]: predict_train_data = random_forest_model.predict(X_test)

         from sklearn import metrics

         print("Accuracy = {0:.3f}".format(metrics.accuracy_score(y_test,predict_train_data)))

         Accuracy = 0.766
```