

```
import pandas as pd
import numpy as np
```

+ Code

+ Text

```
data = pd.read_csv("diabetes.csv")
data.head()
```



	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

```
data.isnull().any()
```



```
Pregnancies      False
Glucose           False
BloodPressure     False
SkinThickness     False
Insulin           False
BMI               False
DiabetesPedigreeFunction False
Age               False
Outcome           False
dtype: bool
```

```
data.describe().T
```



	count	mean	std	min	25%	50%	75%	max
<b>Pregnancies</b>	768.0	3.845052	3.369578	0.000	1.00000	3.0000	6.00000	17.00
<b>Glucose</b>	768.0	120.894531	31.972618	0.000	99.00000	117.0000	140.25000	199.00
<b>BloodPressure</b>	768.0	69.105469	19.355807	0.000	62.00000	72.0000	80.00000	122.00
<b>SkinThickness</b>	768.0	20.536458	15.952218	0.000	0.00000	23.0000	32.00000	99.00
<b>Insulin</b>	768.0	79.799479	115.244002	0.000	0.00000	30.5000	127.25000	846.00
<b>BMI</b>	768.0	31.992578	7.884160	0.000	27.30000	32.0000	36.60000	67.10
<b>DiabetesPedigreeFunction</b>	768.0	0.471876	0.331329	0.078	0.24375	0.3725	0.62625	2.42
<b>Age</b>	768.0	33.240885	11.760232	21.000	24.00000	29.0000	41.00000	81.00
<b>Outcome</b>	768.0	0.348958	0.476951	0.000	0.00000	0.0000	1.00000	1.00

Glucose, BloodPressure, SkinThickness, Insulin, BMI

columns have values 0 which does not make sense, hence are missing values

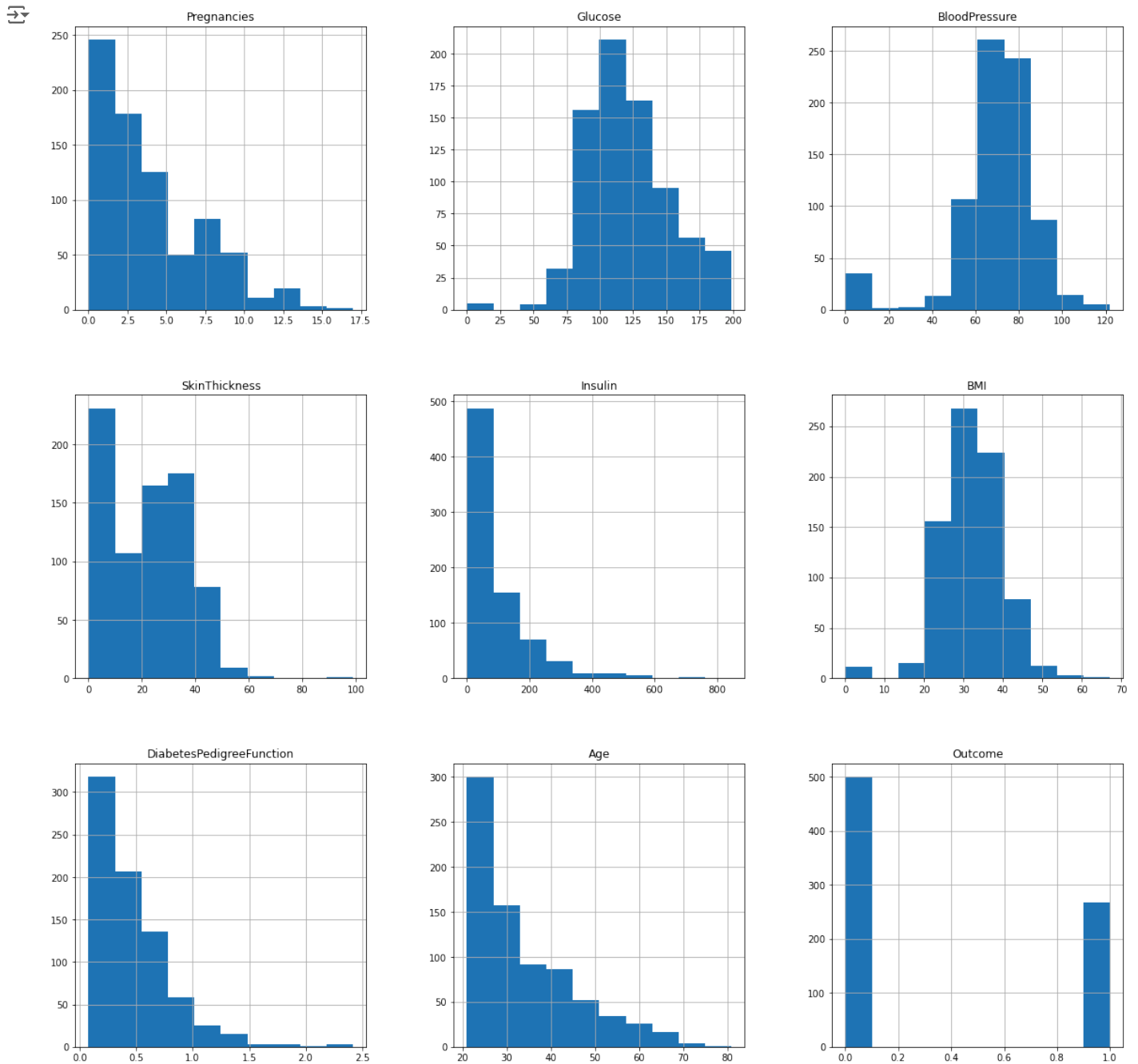
```
data_copy = data.copy(deep = True)
data_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = data_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']]
data_copy.isnull().sum()
```



```
Pregnancies      0
Glucose           5
BloodPressure     35
SkinThickness     227
Insulin           374
BMI               11
DiabetesPedigreeFunction 0
Age               0
Outcome           0
dtype: int64
```

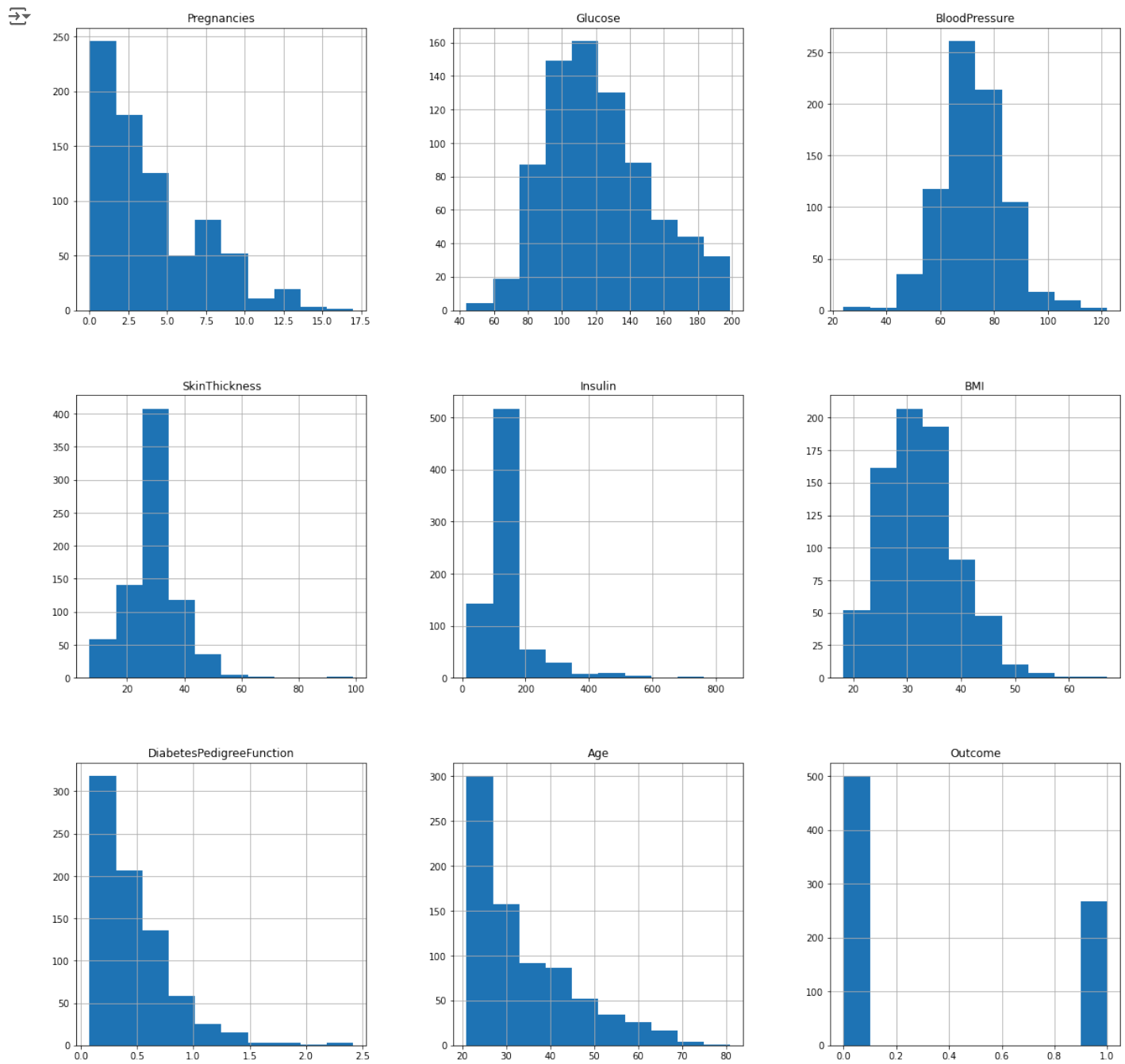
To fill these Nan values the data distribution needs to be understood

```
p = data.hist(figsize = (20,20))
```

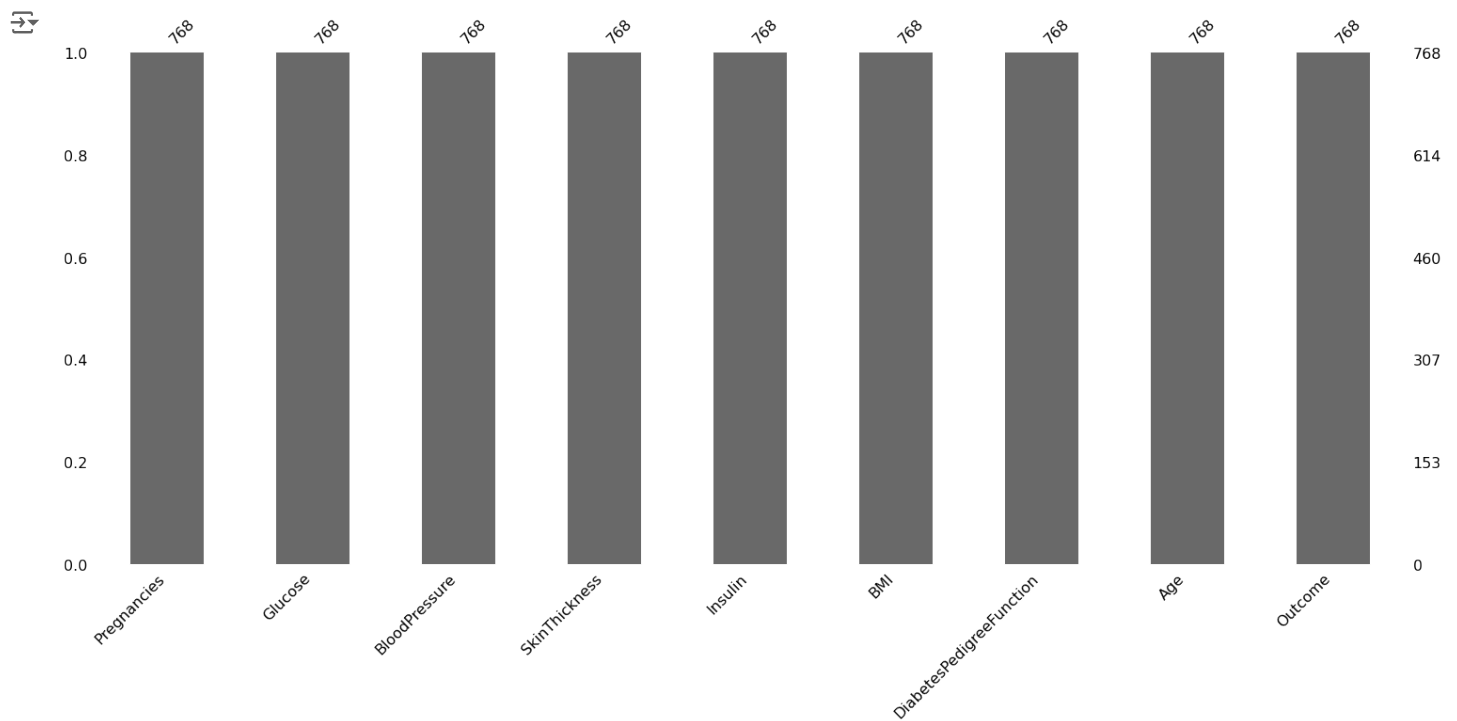


```
data_copy['Glucose'].fillna(data_copy['Glucose'].mean(), inplace = True)
data_copy['BloodPressure'].fillna(data_copy['BloodPressure'].mean(), inplace = True)
data_copy['SkinThickness'].fillna(data_copy['SkinThickness'].median(), inplace = True)
data_copy['Insulin'].fillna(data_copy['Insulin'].median(), inplace = True)
data_copy['BMI'].fillna(data_copy['BMI'].median(), inplace = True)
```

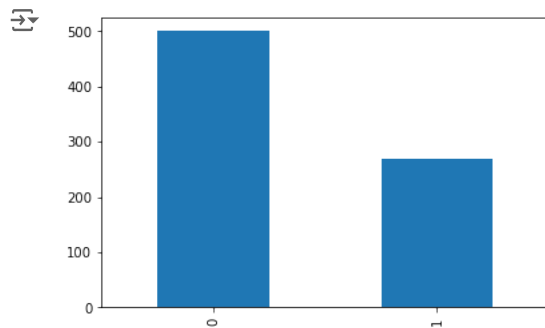
```
p = data_copy.hist(figsize = (20,20))
```



```
import missingno as msno
p = msno.bar(data)
```

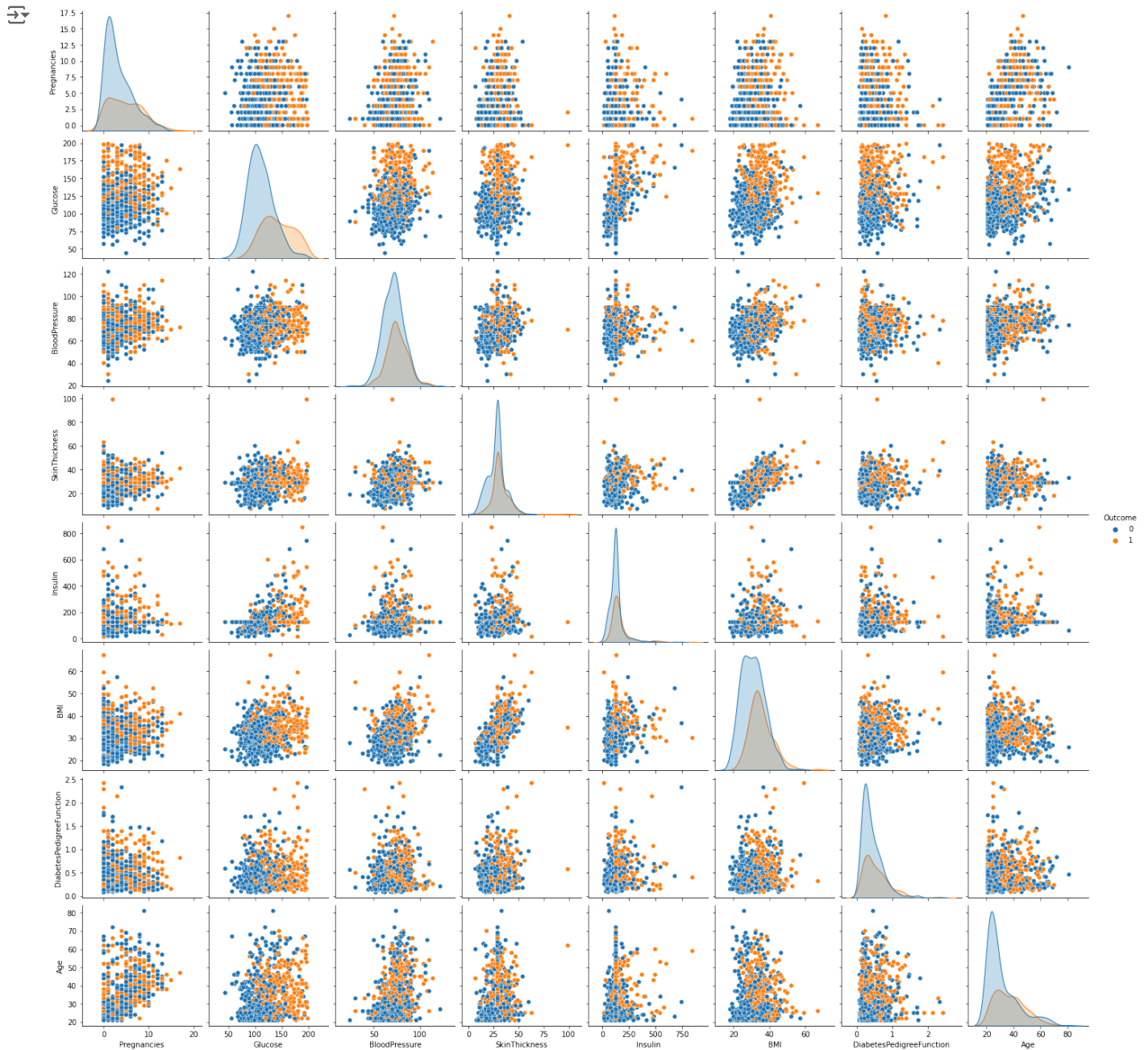


```
p=data.Outcome.value_counts().plot(kind="bar")
```

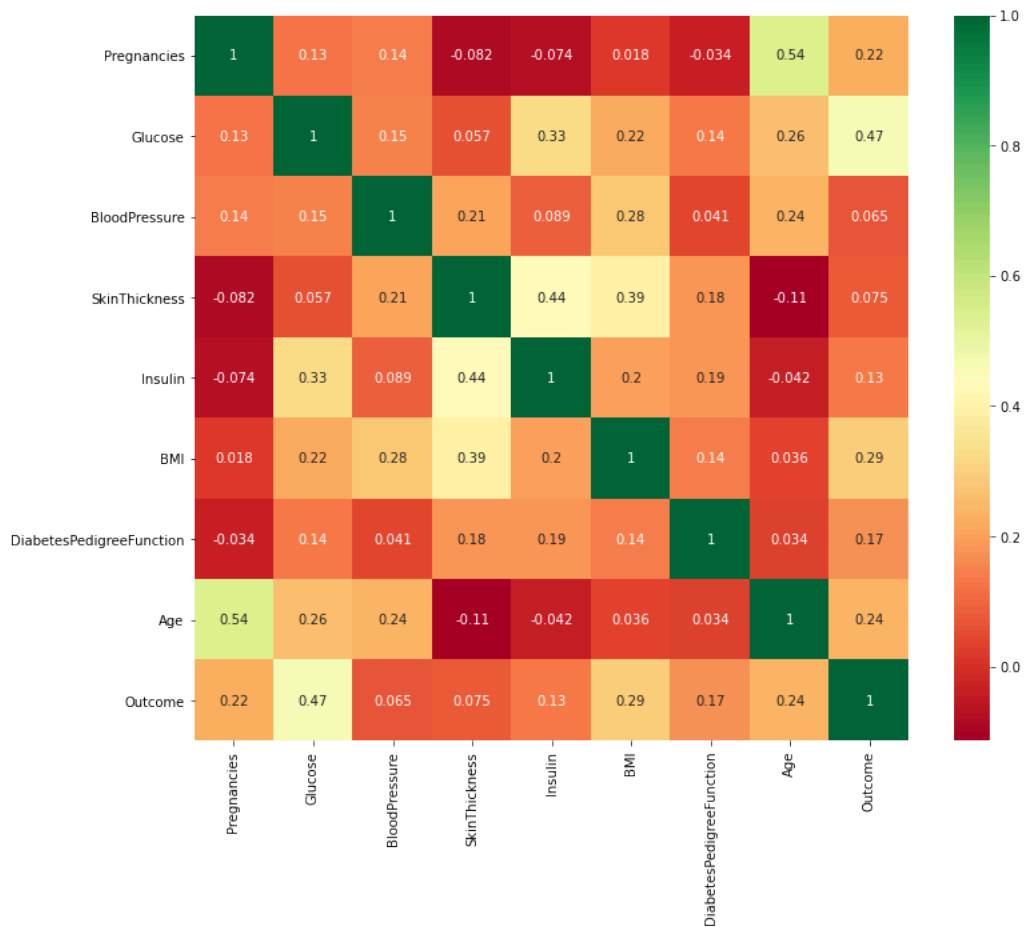


The above graph shows that the data is biased towards datapoints having outcome value as 0 where it means that diabetes was not present actually. The number of non-diabetics is almost twice the number of diabetic patients

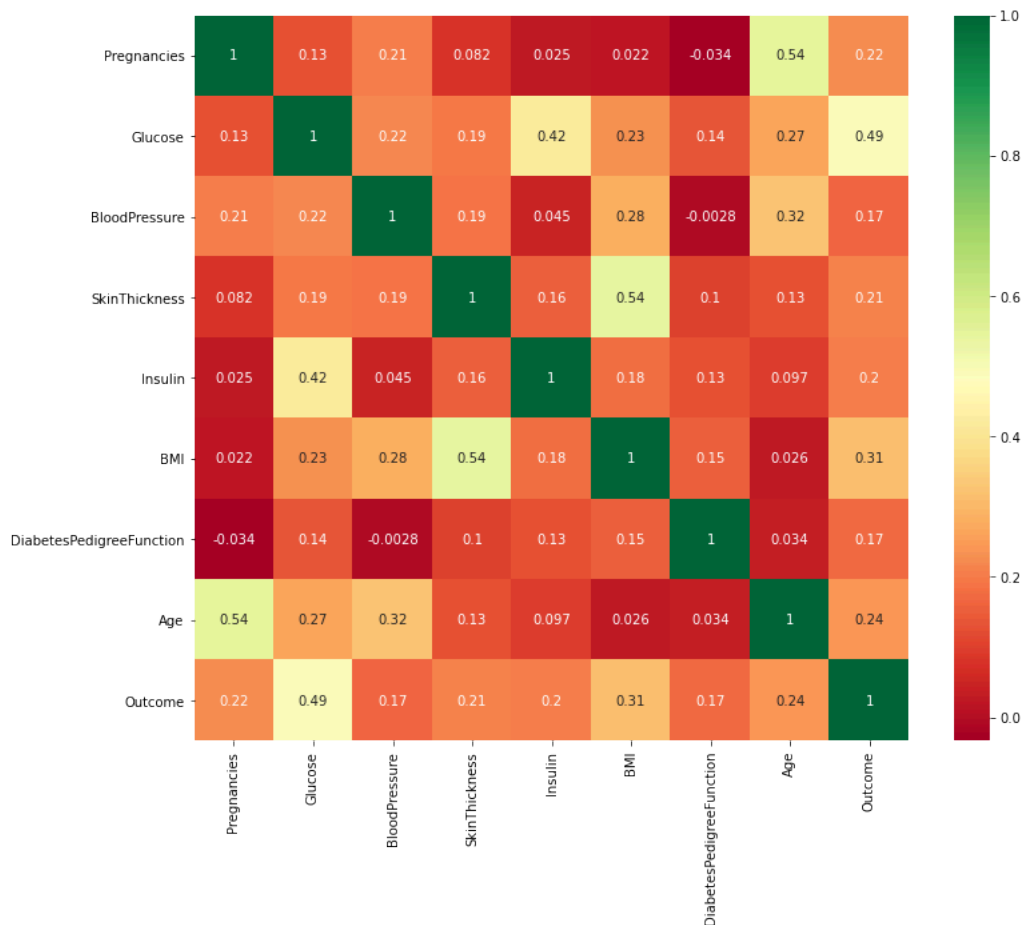
```
import seaborn as sns
p=sns.pairplot(data_copy, hue = 'Outcome')
```



```
import matplotlib.pyplot as plt
plt.figure(figsize=(12,10)) # on this line I just set the size of figure to 12 by 10.
p=sns.heatmap(data.corr(), annot=True,cmap='RdYlGn') # seaborn has very simple solution for heatmap
```



```
plt.figure(figsize=(12,10)) # on this line I just set the size of figure to 12 by 10.
p=sns.heatmap(data_copy.corr(), annot=True,cmap = 'RdYlGn') # seaborn has very simple solution for heatmap
```



```
from sklearn.preprocessing import StandardScaler
sc_X = StandardScaler()
X = pd.DataFrame(sc_X.fit_transform(data_copy.drop(["Outcome"], axis =1)),columns=['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThicknes
    'BMI', 'DiabetesPedigreeFunction', 'Age'])
```

```
X.head()
```



	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	0.639947	0.865108	-0.033518	0.670643	-0.181541	0.166619	0.468492	1.425995
1	-0.844885	-1.206162	-0.529859	-0.012301	-0.181541	-0.852200	-0.365061	-0.190672
2	1.233880	2.015813	-0.695306	-0.012301	-0.181541	-1.332500	0.604397	-0.105584
3	-0.844885	-1.074652	-0.529859	-0.695245	-0.540642	-0.633881	-0.920763	-1.041549
4	-1.141852	0.503458	-2.680669	0.670643	0.316566	1.549303	5.484909	-0.020496

```
y =data_copy.Outcome
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 1/3, random_state = 42, stratify=y)
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
train_scores = []
test_scores = []
```

```
for i in range(1,15):
    knn = KNeighborsClassifier(i)
    knn.fit(X_train, y_train)
    train_scores.append(knn.score(X_train, y_train))
    test_scores.append(knn.score(X_test, y_test))
```

```
max_test_score = max(test_scores)
```

```
test_score_index = [i for i, v in enumerate(test_scores) if v == max_test_score]
```

```
print('Max test score {} % and k = {}'.format(max_test_score*100, list(map(lambda x: x+1, test_score_index))))
```

```
↗ Max test score 76.5625 % and k = [11]
```

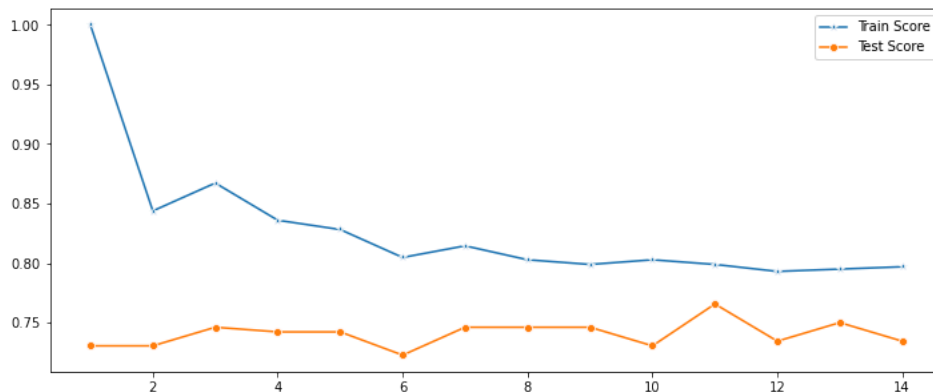
```
plt.figure(figsize=(12,5))
```

```
p = sns.lineplot(range(1,15),train_scores,marker='*',label='Train Score')
```

```
p = sns.lineplot(range(1,15),test_scores,marker='o',label='Test Score')
```

```
↗ C:\Users\Admin\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y.
warnings.warn(
```

```
C:\Users\Admin\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y.
warnings.warn(
```



```
# K=11
```

```
#Setup a knn classifier with k neighbors
```

```
knn = KNeighborsClassifier(11)
```

```
knn.fit(X_train,y_train)
```

```
knn.score(X_test,y_test)
```

```
↗ 0.765625
```

```
from mlxtend.plotting import plot_decision_regions
```

```
value = 20000
```

```
width = 20000
```

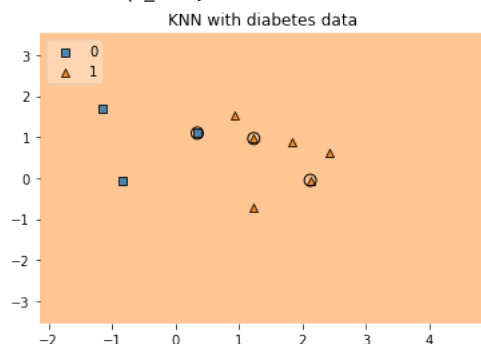
```
plot_decision_regions(X.values, y.values, clf = knn, legend = 2,filler_feature_values={2: value, 3: value, 4: value, 5: value, 6: value, 7: value, 8: value, 9: value, 10: value, 11: value, 12: value, 13: value, 14: value},
                    filler_feature_ranges={2: width, 3: width, 4: width, 5: width, 6: width, 7: width, 8: width, 9: width, 10: width, 11: width, 12: width, 13: width, 14: width},
                    X_highlight=X_test.values)
```

```
plt.title("KNN with diabetes data")
```

```
plt.show()
```

```
↗ C:\Users\Admin\anaconda3\lib\site-packages\mlxtend\plotting\decision_regions.py:244: UserWarning: No contour levels were found within the specified range.
ax.contour(xx, yy, Z, cset.levels,
```

```
C:\Users\Admin\anaconda3\lib\site-packages\mlxtend\plotting\decision_regions.py:316: MatplotlibDeprecationWarning: Using a string of size 1 to index arrays and matrices is deprecated and will raise an error in the future. Use an array, list, tuple, or string of size greater than one, instead.
ax.scatter(x_data,
```



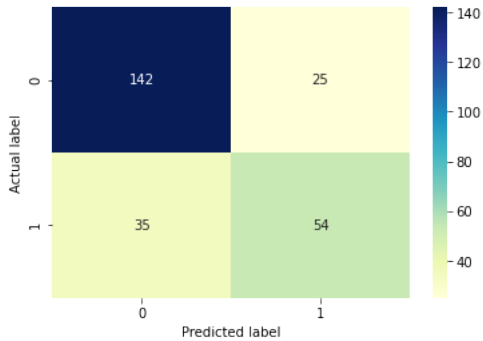


```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, fbeta_score
y_pred = knn.predict(X_test)
```

```
cnf_matrix = confusion_matrix(y_test, y_pred)
```

```
p = sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu", fmt='g')
plt.title('Confusion matrix', y=1.1)
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
```

```
Text(0.5, 15.0, 'Predicted label')
Confusion matrix
```



```
def model_evaluation(y_test, y_pred, model_name):
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    f2 = fbeta_score(y_test, y_pred, beta = 2.0)

    results = pd.DataFrame([[model_name, acc, prec, rec, f1, f2]],
                           columns = ["Model", "Accuracy", "Precision", "Recall",
                                     "F1 Score", "F2 Score"])
    results = results.sort_values(["Precision", "Recall", "F2 Score"], ascending = False)
    return results
```

```
model_evaluation(y_test, y_pred, "KNN")
```

```
Model Accuracy Precision Recall F1 Score F2 Score
0 KNN 0.765625 0.683544 0.606742 0.642857 0.62069
```

```
# Alternate way
from sklearn.metrics import classification_report
print(classification_report(y_test,y_pred))
```

```
precision recall f1-score support
0 0.80 0.85 0.83 167
1 0.68 0.61 0.64 89

accuracy 0.77 256
macro avg 0.74 0.73 0.73 256
weighted avg 0.76 0.77 0.76 256
```

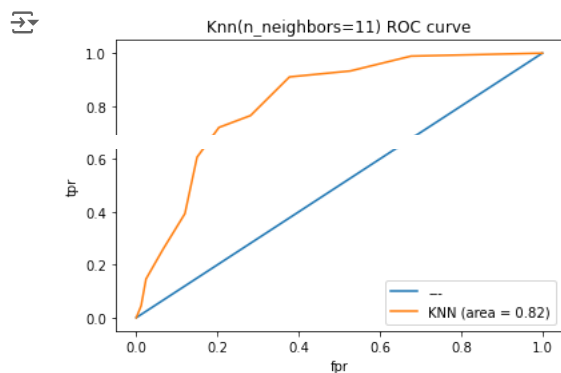
```
from sklearn.metrics import auc, roc_auc_score, roc_curve
```

```
y_pred_proba = knn.predict_proba(X_test)[:,-1]
fpr, tpr, threshold = roc_curve(y_test, y_pred_proba)
```

```
classifier_roc_auc = roc_auc_score(y_test, y_pred_proba)
plt.plot([0,1],[0,1], label = "---")
```

```
plt.plot(fpr, tpr, label = 'KNN (area = %0.2f)' % classifier_roc_auc)
plt.xlabel("fpr")
```

```
plt.ylabel("tpr")
plt.title('Knn(n_neighbors=11) ROC curve')
plt.legend(loc="lower right", fontsize = "medium")
plt.xticks(rotation=0, horizontalalignment="center")
plt.yticks(rotation=0, horizontalalignment="right")
```



#Hyper parameters tuning using GridSearchCV

```
from sklearn.model_selection import GridSearchCV
parameters_grid = {"n_neighbors": np.arange(0,50)}
knn= KNeighborsClassifier()
knn_GSV = GridSearchCV(knn, param_grid=parameters_grid, cv = 5)
knn_GSV.fit(X, y)
```

```
C:\Users\Admin\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py:548: FitFailedWarning: Estimator fit failed. The sc
Traceback (most recent call last):
  File "C:\Users\Admin\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py", line 531, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
  File "C:\Users\Admin\anaconda3\lib\site-packages\sklearn\neighbors\_base.py", line 1157, in fit
    return self._fit(X)
  File "C:\Users\Admin\anaconda3\lib\site-packages\sklearn\neighbors\_base.py", line 467, in _fit
    raise ValueError(
ValueError: Expected n_neighbors > 0. Got 0
```