

Capstone Project: Healthcare - Data Science PGP

****Student Name-Sourabh Singh Thakur****

****Problem Statement:**

- NIDDK (National Institute of Diabetes and Digestive and Kidney Diseases) research creates knowledge about and treatments for the most chronic, costly, and consequential diseases.
 - The dataset used in this project is originally from NIDDK. The objective is to predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset.
 - Build a model to accurately predict whether the patients in the dataset have diabetes or not.
- Dataset Description: The datasets consists of several medical predictor variables and one target variable (Outcome). Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and more.

Variables - Description

- Pregnancies - Number of times pregnant
- Glucose - Plasma glucose concentration in an oral glucose tolerance test
- BloodPressure - Diastolic blood pressure (mm Hg)
- SkinThickness - Triceps skinfold thickness (mm)
- Insulin - Two hour serum insulin
- BMI - Body Mass Index
- DiabetesPedigreeFunction - Diabetes pedigree function
- Age - Age in years
- Outcome - Class variable (either 0 or 1). 268 of 768 values are 1, and the others are 0

****Week 1:**

Data Exploration:

-Perform descriptive analysis. Understand the variables and their corresponding values. On the columns below, a value of zero does not make sense and thus indicates missing value:

- Glucose
- BloodPressure
- SkinThickness
- Insulin
- BMI

-Visually explore these variables using histograms. Treat the missing values accordingly.

-There are integer and float data type variables in this dataset. Create a count (frequency) plot describing the data types and the count of variables.

****Solution:**

Week 1:

Data Exploration:

(1) Read Data and Perform descriptive analysis:

In [1]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="white", color_codes=True)
sns.set(font_scale=1.2)
```

In [2]:

```
df = pd.read_csv("D:\\DATASET//healthcare.csv")
```

In [3]:

```
df.head()
```

Out[3]:

	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

According to problem statement, a value of zero in the following columns indicates missing value:

- Glucose
- BloodPressure
- SkinThickness
- Insulin
- BMI

We will replace zeros in these columns with null values.

```
In [4]: cols_with_null_as_zero = ['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']
df[cols_with_null_as_zero] = df[cols_with_null_as_zero].replace(0, np.NaN)
```

```
In [5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Pregnancies      768 non-null    int64  
 1   Glucose          763 non-null    float64 
 2   BloodPressure    733 non-null    float64 
 3   SkinThickness    541 non-null    float64 
 4   Insulin          394 non-null    float64 
 5   BMI              757 non-null    float64 
 6   DiabetesPedigreeFunction 768 non-null    float64 
 7   Age              768 non-null    int64  
 8   Outcome          768 non-null    int64  
dtypes: float64(6), int64(3)
memory usage: 54.1 KB
```

```
In [6]: df.isnull().sum()
```

```
Out[6]: Pregnancies      0
Glucose          5
BloodPressure    35
SkinThickness    227
Insulin          374
BMI              11
DiabetesPedigreeFunction 0
Age              0
Outcome          0
dtype: int64
```

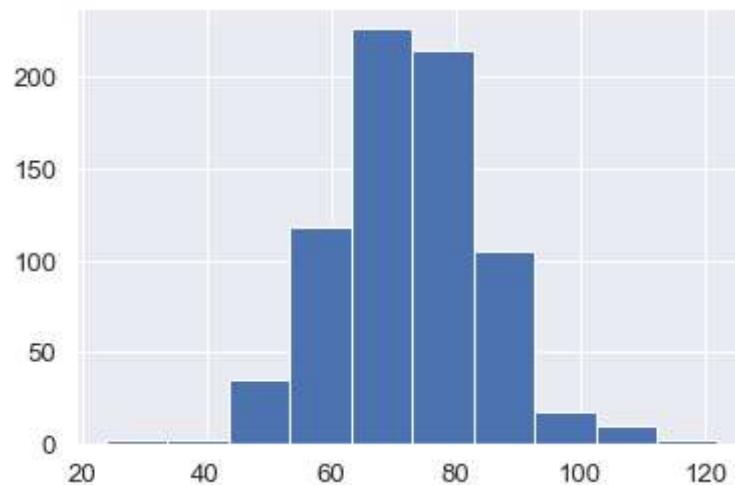
In [7]: `df.describe()`

Out[7]:

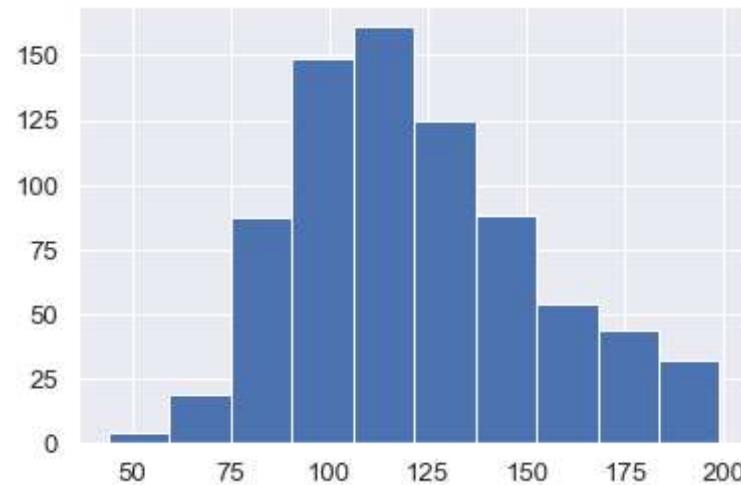
	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	763.000000	733.000000	541.000000	394.000000	757.000000	768.000000	768.000000	768.000000
mean	3.845052	121.686763	72.405184	29.153420	155.548223	32.457464	0.471876	33.240885	0.348958
std	3.369578	30.535641	12.382158	10.476982	118.775855	6.924988	0.331329	11.760232	0.476951
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	64.000000	22.000000	76.250000	27.500000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	29.000000	125.000000	32.300000	0.372500	29.000000	0.000000
75%	6.000000	141.000000	80.000000	36.000000	190.000000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

**(2) Visually explore these variables using histograms and treat the missing values accordingly:

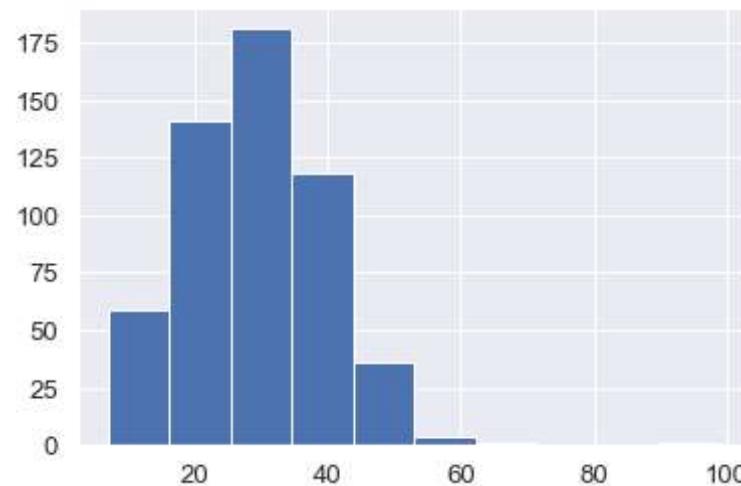
In [8]: `df['BloodPressure'].hist();`



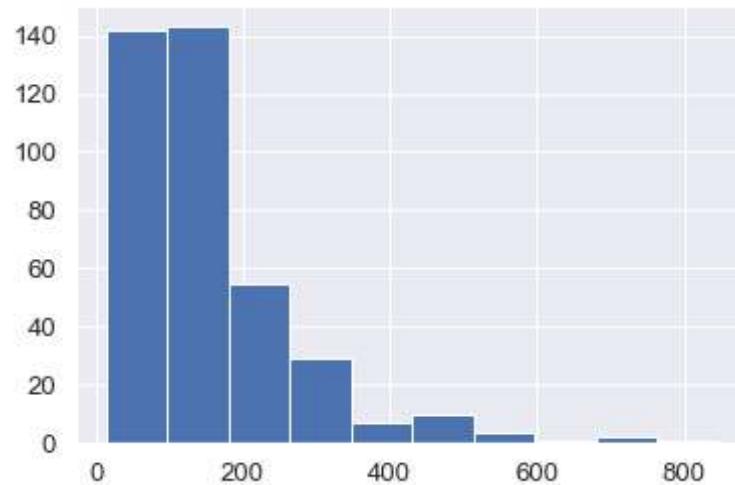
```
In [9]: df['Glucose'].hist();
```



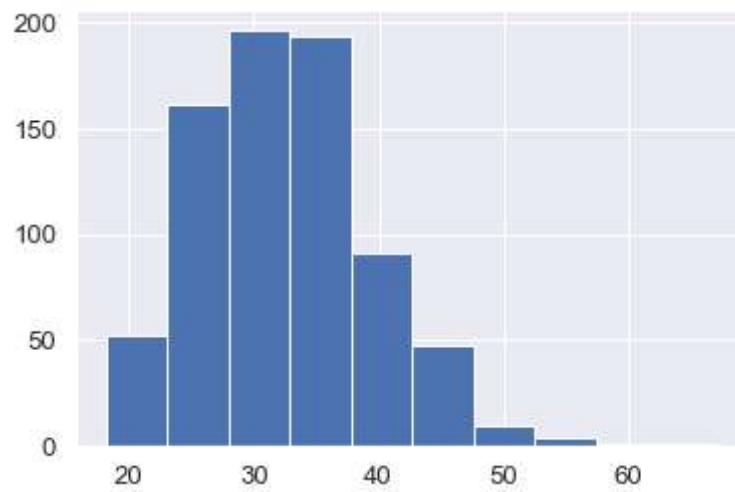
```
In [10]: df['SkinThickness'].hist();
```



```
In [11]: df['Insulin'].hist();
```



```
In [12]: df['BMI'].hist();
```



From above histograms, it is clear that ****Insulin**** has highly skewed data distribution and remaining 4 variables have relatively balanced data distribution therefore we will treat missing values in these 5 variables as below:-

- Glucose - replace missing values with mean of values.
- BloodPressure - replace missing values with mean of values.
- SkinThickness - replace missing values with mean of values.
- Insulin - replace missing values with median of values.
- BMI - replace missing values with mean of values.

```
In [13]: df['Insulin'] = df['Insulin'].fillna(df['Insulin'].median())
```

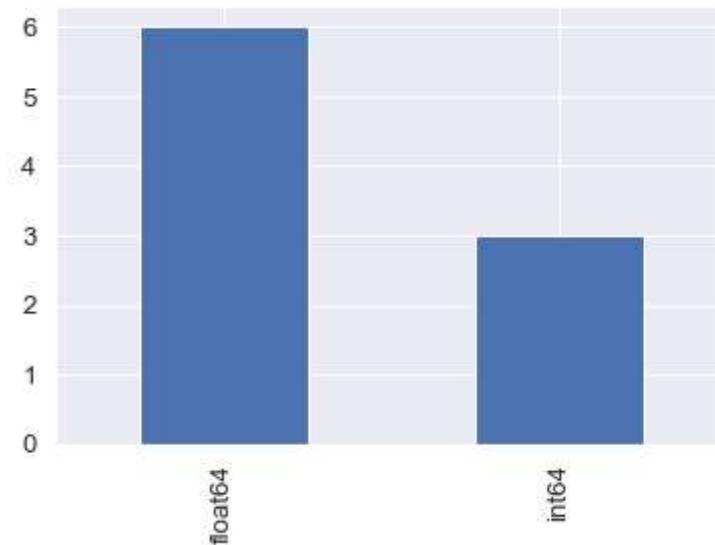
```
In [14]: cols_mean_for_null = ['Glucose', 'BloodPressure', 'SkinThickness', 'BMI']
df[cols_mean_for_null] = df[cols_mean_for_null].fillna(df[cols_mean_for_null].mean())
```

```
In [15]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Pregnancies      768 non-null    int64  
 1   Glucose          768 non-null    float64 
 2   BloodPressure    768 non-null    float64 
 3   SkinThickness    768 non-null    float64 
 4   Insulin          768 non-null    float64 
 5   BMI              768 non-null    float64 
 6   DiabetesPedigreeFunction 768 non-null    float64 
 7   Age              768 non-null    int64  
 8   Outcome          768 non-null    int64  
dtypes: float64(6), int64(3)
memory usage: 54.1 KB
```

****3) Create a count (frequency) plot describing the data types and the count of variables:**

```
In [16]: df.dtypes.value_counts().plot(kind = 'bar');
```



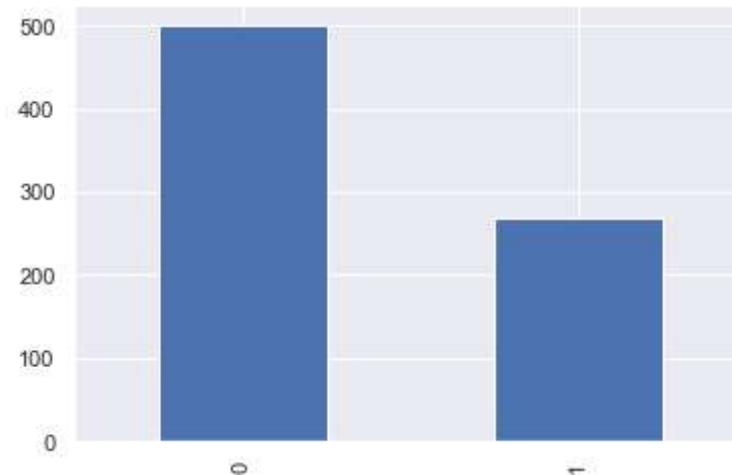
****Week 2:**

Data Exploration:

(1) Check the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan future course of action:

```
In [22]: df['Outcome'].value_counts().plot(kind='bar')
df['Outcome'].value_counts()
```

```
Out[22]: 0    500
1    268
Name: Outcome, dtype: int64
```



**Since classes in Outcome is little skewed so we will generate new samples using SMOTE (Synthetic Minority Oversampling Technique) for the class '1' which is under-represented in our data. We will use SMOTE out of many other techniques available since:

- It generates new samples by interpolation.
- It doesn't duplicate data.

```
In [23]: df_X = df.drop('Outcome', axis=1)
df_y = df['Outcome']
print(df_X.shape, df_y.shape)
```

```
(768, 8) (768,)
```

```
In [25]: from imblearn.over_sampling import SMOTE
```

```
In [26]: df_X_resampled, df_y_resampled = SMOTE(random_state=500).fit_resample(df_X, df_y)
print(df_X_resampled.shape, df_y_resampled.shape)
```

```
(1000, 8) (1000,)
```

```
In [27]: df_y_resampled.value_counts().plot(kind='bar')
df_y_resampled.value_counts()
```

```
Out[27]: 1    500
0    500
Name: Outcome, dtype: int64
```



**2) Create scatter charts between the pair of variables to understand the relationships. Describe your findings:

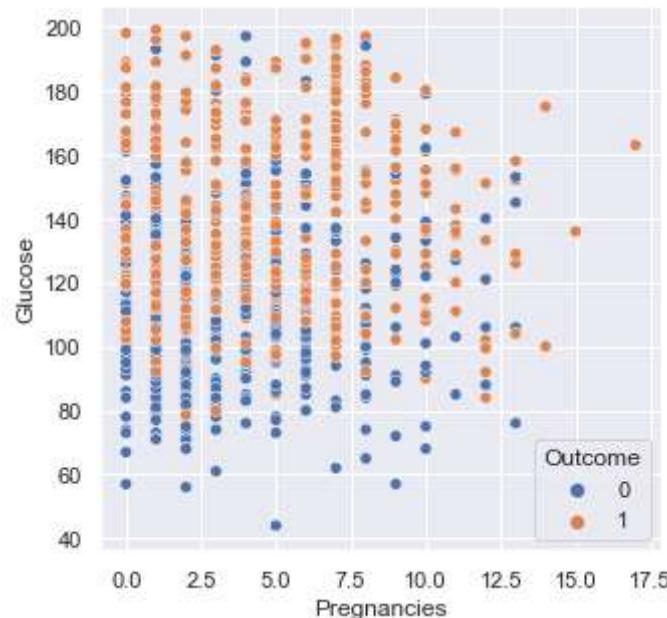
```
In [28]: df_resampled = pd.concat([df_X_resampled, df_y_resampled], axis=1)
df_resampled
```

Out[28]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome	
0	6	148.000000	72.000000	35.000000	125.000000	33.600000		0.627000	50	1
1	1	85.000000	66.000000	29.000000	125.000000	26.600000		0.351000	31	0
2	8	183.000000	64.000000	29.153420	125.000000	23.300000		0.672000	32	1
3	1	89.000000	66.000000	23.000000	94.000000	28.100000		0.167000	21	0
4	0	137.000000	40.000000	35.000000	168.000000	43.100000		2.288000	33	1
...	
995	4	132.146483	72.823785	29.153420	125.000000	25.859463		0.244872	58	1
996	2	106.554819	62.296787	16.561448	50.374299	24.487812		0.653070	22	1
997	2	145.658166	81.170917	43.512752	178.134564	43.633479		0.450657	43	1
998	8	145.042368	79.915264	45.643123	129.894080	37.931776		0.628950	40	1
999	3	165.024147	68.975853	36.585346	233.293088	33.297502		0.561254	29	1

1000 rows × 9 columns

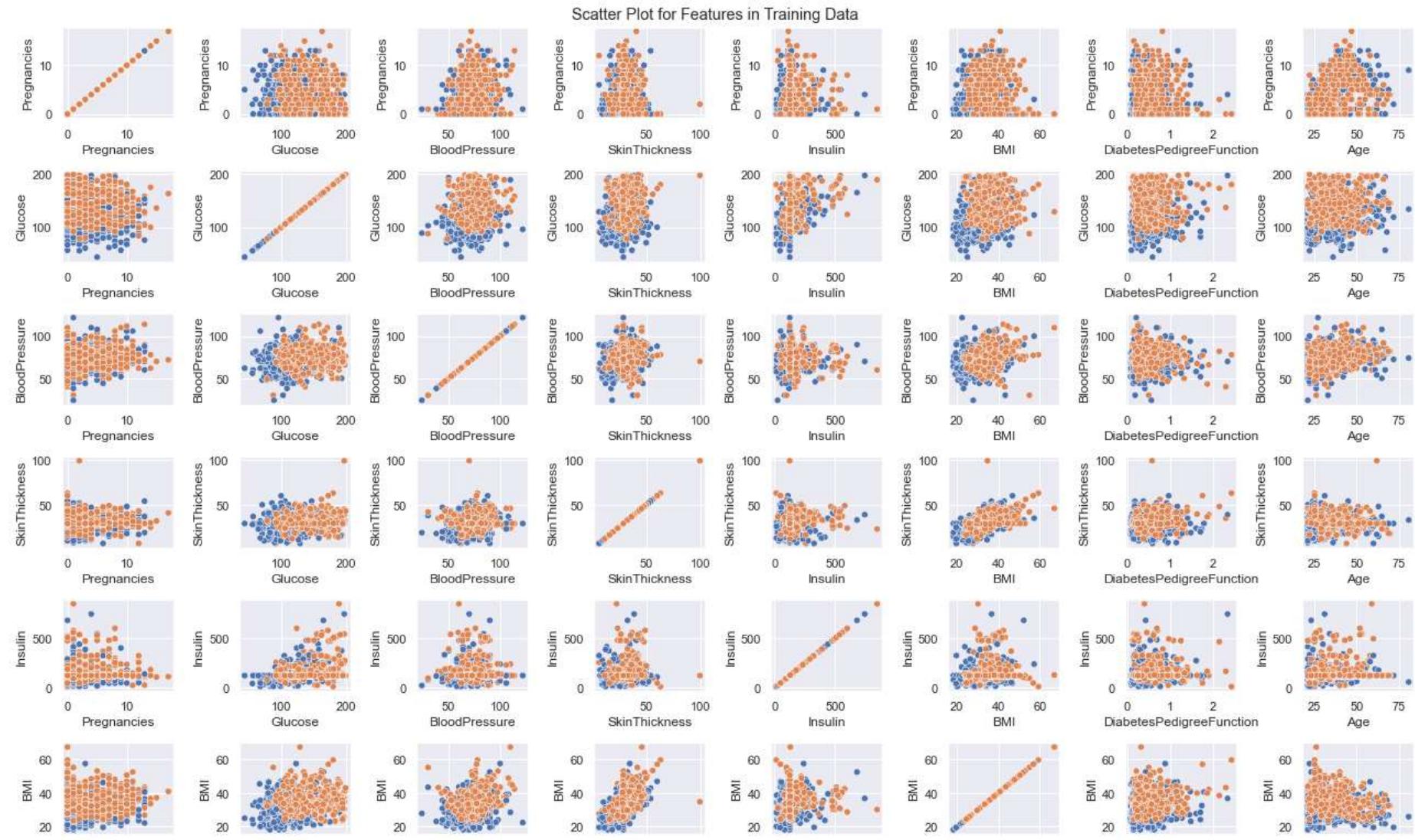
```
In [30]: sns.set(rc={'figure.figsize':(5,5)})  
sns.scatterplot(x='Pregnancies', y='Glucose', data=df_resampled, hue='Outcome');
```

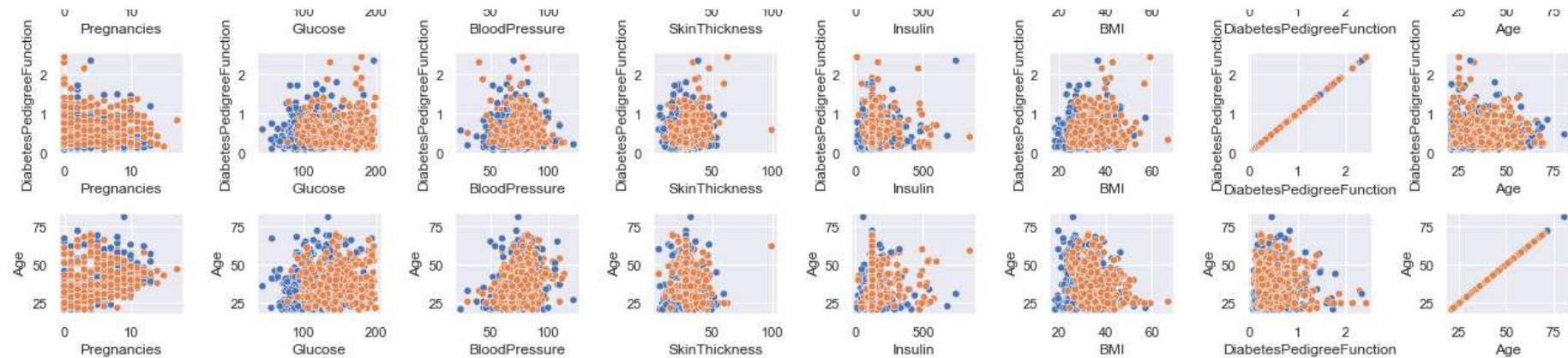


```
In [31]: fig, axes = plt.subplots(8, 8, figsize=(18, 15))
fig.suptitle('Scatter Plot for Features in Training Data')

for i, col_y in enumerate(df_X_resampled.columns):
    for j, col_x in enumerate(df_X_resampled.columns):
        sns.scatterplot(ax=axes[i, j], x=col_x, y=col_y, data=df_resampled, hue="Outcome", legend = False)

plt.tight_layout()
```





We have some interesting observations from above scatter plot of pairs of features:

- **Glucose** alone is impressively good to distinguish between the **Outcome** classes.
- **Age** alone is also able to distinguish between classes to some extent.
- It seems none of pairs in the dataset is able to clearly distinguish between the **Outcome** classes.
- We need to use combination of features to build model for prediction of classes in **Outcome**.

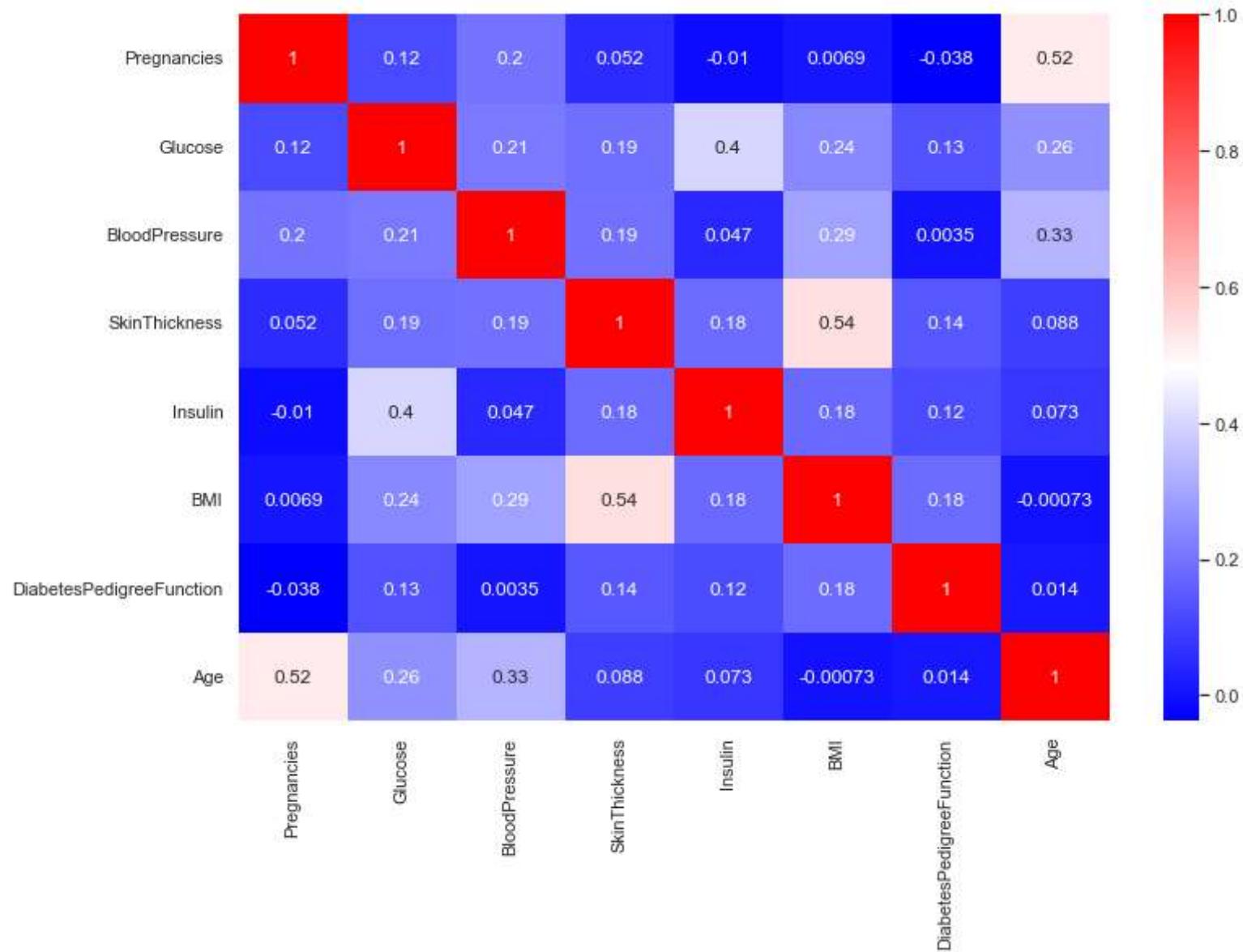
**(3) Perform correlation analysis. Visually explore it using a heat map:

In [32]: `df_X_resampled.corr()`

Out[32]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
Pregnancies	1.000000	0.115771	0.200407	0.052411	-0.010482	0.006903	-0.037803	0.518113
Glucose	0.115771	1.000000	0.210302	0.190414	0.399092	0.238361	0.129039	0.260946
BloodPressure	0.200407	0.210302	1.000000	0.194699	0.046989	0.291062	0.003523	0.330728
SkinThickness	0.052411	0.190414	0.194699	1.000000	0.183435	0.542466	0.141941	0.088340
Insulin	-0.010482	0.399092	0.046989	0.183435	1.000000	0.177066	0.116594	0.072942
BMI	0.006903	0.238361	0.291062	0.542466	0.177066	1.000000	0.181157	-0.000725
DiabetesPedigreeFunction	-0.037803	0.129039	0.003523	0.141941	0.116594	0.181157	1.000000	0.014423
Age	0.518113	0.260946	0.330728	0.088340	0.072942	-0.000725	0.014423	1.000000

```
In [33]: plt.figure(figsize=(12,8))
sns.heatmap(df_X_resampled.corr(), cmap='bwr', annot=True);
```



It appears from correlation matrix and heatmap that there exists significant correlation between some pairs such as -

- Age-Pregnancies
- BMI-SkinThickness

**Week 3:

Data Modeling:

1. Devise strategies for model building. It is important to decide the right validation framework. Express your thought process:

Answer: Since this is a classification problem, we will be building all popular classification models for our training data and then compare performance of each model on test data to accurately predict target variable (Outcome):

- 1) Logistic Regression
- 2) Decision Tree
- 3) RandomForest Classifier
- 4) K-Nearest Neighbour (KNN)
- 5) Support Vector Machine (SVM)
- 6) Naive Bayes
- 7) Ensemble Learning -> Boosting -> Adaptive Boosting
- 8) Ensemble Learning -> Boosting -> Gradient Boosting (XGBClassifier)

We will use ****GridSearchCV**** with Cross Validation (CV) = 5 for training and testing model which will give us insight about model performance on versatile data. It helps to loop through predefined hyperparameters and fit model on training set. GridSearchCV performs hyper parameter tuning which will give us optimal hyper parameters for each of the model. We will again train model with these optimized hyper parameters and then predict test data to get metrics for comparing all models.

****Performing Train - Test split on input data (To train and test model without Cross Validation and Hyper Parameter Tuning):**

In [220]:

```
from sklearn.model_selection import train_test_split, KFold, RandomizedSearchCV
from sklearn.metrics import accuracy_score, average_precision_score, f1_score, confusion_matrix, classification_report,
```

In [221]: `X_train, X_test, y_train, y_test = train_test_split(df_X_resampled, df_y_resampled, test_size=0.15, random_state =10)`

```
In [222]: X_train.shape, X_test.shape
```

```
Out[222]: ((850, 8), (150, 8))
```

2. Apply an appropriate classification algorithm to build a model. Compare various models with the results from KNN algorithm.

```
In [223]: models = []
model_accuracy = []
model_f1 = []
model_auc = []
```

1) Logistic Regression:

```
In [224]: from sklearn.linear_model import LogisticRegression
lr1 = LogisticRegression(max_iter=300)
```

```
In [225]: lr1.fit(X_train,y_train)
```

```
Out[225]: LogisticRegression
LogisticRegression(max_iter=300)
```

```
In [226]: lr1.score(X_train,y_train)
```

```
Out[226]: 0.7435294117647059
```

```
In [227]: lr1.score(X_test, y_test)
```

```
Out[227]: 0.7066666666666667
```

Performance evaluation and optimizing parameters using GridSearchCV: Logistic regression does not really have any critical hyperparameters to tune. However we will try to optimize one of its parameters 'C' with the help of GridSearchCV. So we have set this parameter as a list of values from which GridSearchCV will select the best value of parameter.

```
In [228]: from sklearn.model_selection import GridSearchCV, cross_val_score
```

```
In [229]: parameters = {'C':np.logspace(-5, 5, 50)}
```

```
In [230]: gs_lr = GridSearchCV(lr1, param_grid = parameters, cv=5, verbose=0)
gs_lr.fit(df_X_resampled, df_y_resampled)
```

```
Out[230]:
```

```
GridSearchCV
    ...
    6.55128557e-05, 1.04811313e-04, 1.67683294e-04, 2.68269580e-04,
    4.29193426e-04, 6.86648845e-04, 1.09854114e-03, 1.75751062e-03,
    2.81176870e-03, 4.49843267e-03, 7.19685673e-03, 1.15139540e-02,
    1.84206997e-02, 2.94705170e...
    7.90604321e-01, 1.26485522e+00, 2.02358965e+00, 3.23745754e+00,
    5.17947468e+00, 8.28642773e+00, 1.32571137e+01, 2.12095089e+01,
    3.39322177e+01, 5.42867544e+01, 8.68511374e+01, 1.38949549e+02,
    2.22299648e+02, 3.55648031e+02, 5.68986603e+02, 9.10298178e+02,
    1.45634848e+03, 2.32995181e+03, 3.72759372e+03, 5.96362332e+03,
    9.54095476e+03. 1.52641797e+04. 2.44205309e+04. 3.90693994e+04.
    ▶ estimator: LogisticRegression
        ▶ LogisticRegression
```

```
In [232]: gs_lr.best_params_
```

```
Out[232]: {'C': 5.1794746792312125}
```

```
In [233]: gs_lr.best_score_
```

```
Out[233]: 0.741
```

```
In [234]: lr2 = LogisticRegression(C=5.1794746792312125, max_iter=300)
```

```
In [235]: lr2.fit(X_train,y_train)
```

```
Out[235]:
```

`LogisticRegression`

`LogisticRegression(C=5.1794746792312125, max_iter=300)`

```
In [236]: lr2.score(X_train,y_train)
```

```
Out[236]: 0.7435294117647059
```

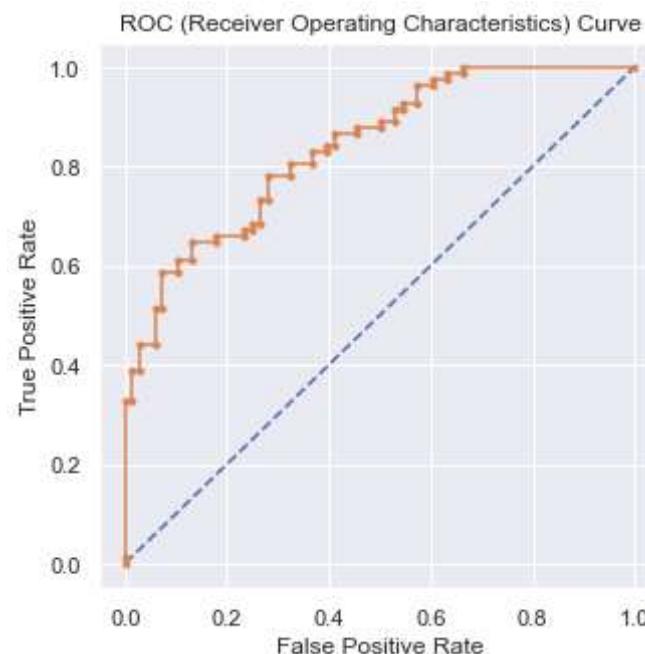
```
In [237]: lr2.score(X_test, y_test)
```

```
Out[237]: 0.7066666666666667
```

In [238]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)

```
probs = lr2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                      # keep probabilities for the positive outcome only
auc_lr = roc_auc_score(y_test, probs)     # calculate AUC
print('AUC: %.3f' %auc_lr)
fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')    # plot no skill
plt.plot(fpr, tpr, marker='.')             # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

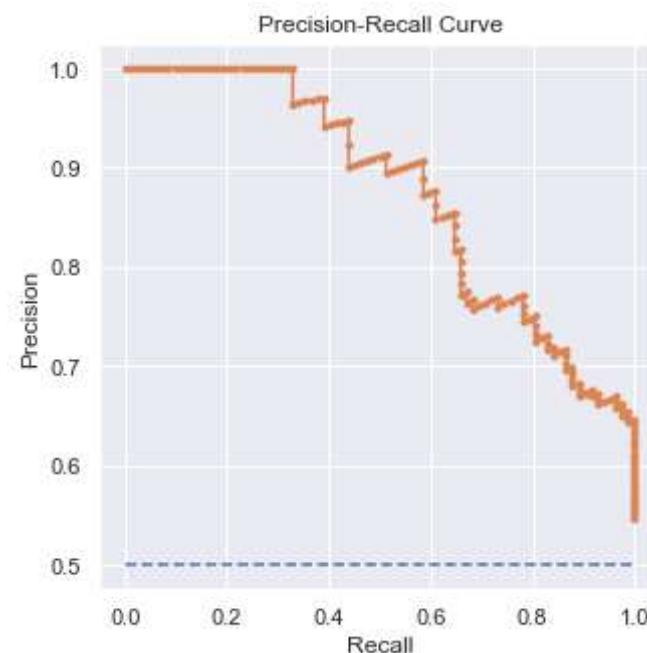
AUC: 0.839



In [239]: # Precision Recall Curve

```
pred_y_test = lr2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_lr_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_lr_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.') # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.711 auc_pr=0.874 ap=0.875



```
In [240]: models.append('LR')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_lr)
```

2) Decision Tree:

```
In [241]: from sklearn.tree import DecisionTreeClassifier
dt1 = DecisionTreeClassifier(random_state=0)
```

```
In [242]: dt1.fit(X_train,y_train)
```

```
Out[242]:
```

```
DecisionTreeClassifier
DecisionTreeClassifier(random_state=0)
```

```
In [243]: dt1.score(X_train,y_train)          # Decision Tree always 100% accuracy over train data
```

```
Out[243]: 1.0
```

```
In [244]: dt1.score(X_test, y_test)
```

```
Out[244]: 0.7933333333333333
```

**Performance evaluation and optimizing parameters using GridSearchCV:

```
In [245]: parameters = {
    'max_depth':[1,2,3,4,5,None]
}
```

```
In [246]: gs_dt = GridSearchCV(dt1, param_grid = parameters, cv=5, verbose=0)
gs_dt.fit(df_X_resampled, df_y_resampled)
```

```
Out[246]:
```

```
GridSearchCV
GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=0),
             param_grid={'max_depth': [1, 2, 3, 4, 5, None]})

  ▶ estimator: DecisionTreeClassifier
    ▶ DecisionTreeClassifier
```

```
In [247]: gs_dt.best_params_
```

```
Out[247]: {'max_depth': 5}
```

```
In [248]: gs_dt.best_score_
```

```
Out[248]: 0.7689999999999999
```

```
In [249]: X_train.columns
```

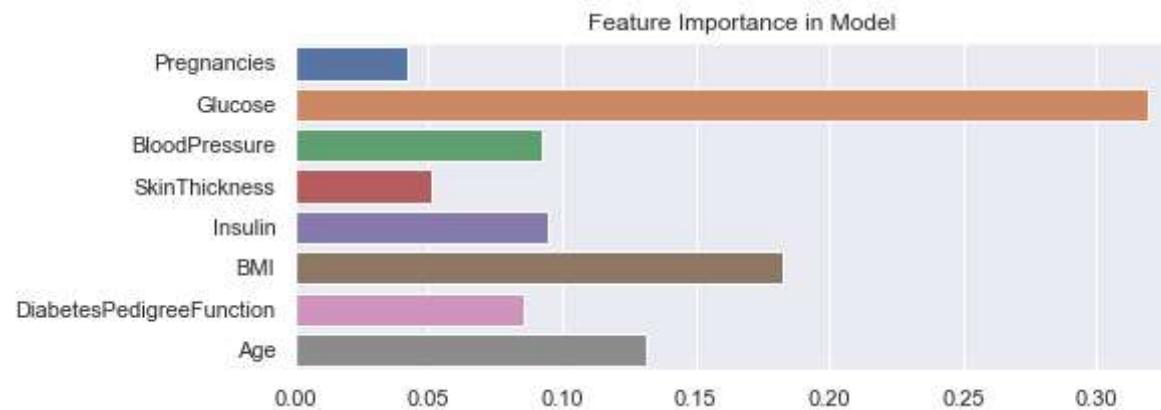
```
Out[249]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
                 'BMI', 'DiabetesPedigreeFunction', 'Age'],
                 dtype='object')
```

```
In [250]: dt1.feature_importances_
```

```
Out[250]: array([0.04238187, 0.31916028, 0.09248709, 0.05091529, 0.09483917,
                 0.18270404, 0.08596092, 0.13155132])
```

```
In [251]: import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(8,3))
sns.barplot(y=X_train.columns, x=dt1.feature_importances_)
plt.title("Feature Importance in Model");
```



```
In [253]: dt2 = DecisionTreeClassifier(max_depth=5)
```

```
In [257]: dt2.fit(X_train,y_train)
```

```
Out[257]:
```

```
▼      DecisionTreeClassifier
DecisionTreeClassifier(max_depth=5)
```

```
In [258]: dt2.score(X_train,y_train)
```

```
Out[258]: 0.8352941176470589
```

```
In [259]: dt2.score(X_test, y_test)
```

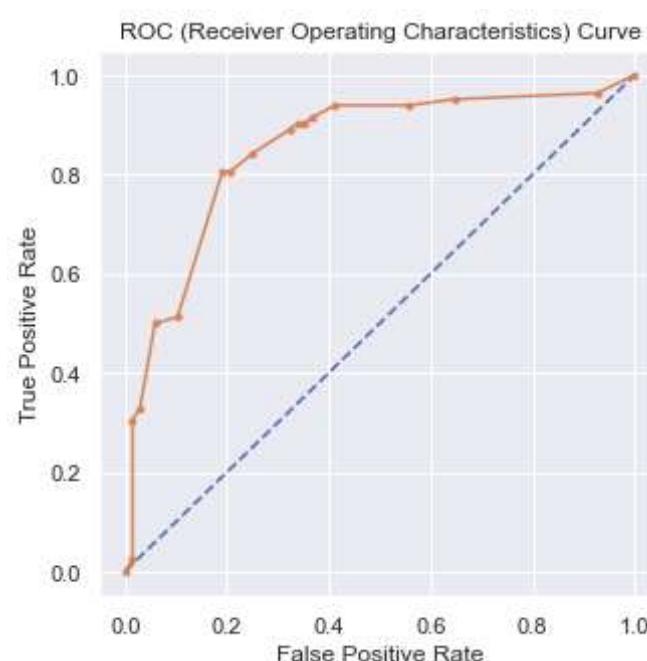
```
Out[259]: 0.8
```

```
In [260]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)
```

```
probs = dt2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                       # keep probabilities for the positive outcome only

auc_dt = roc_auc_score(y_test, probs)      # calculate AUC
print('AUC: %.3f' %auc_dt)
fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')    # plot no skill
plt.plot(fpr, tpr, marker '.')             # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

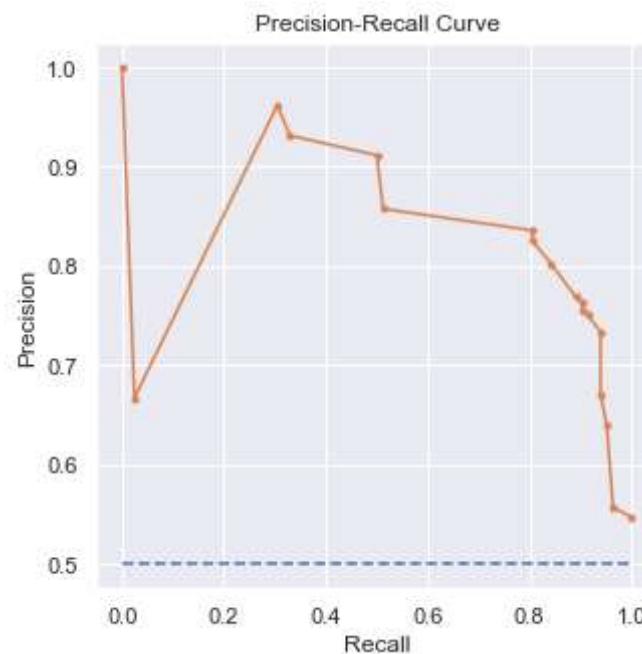
AUC: 0.851



In [261]: # Precision Recall Curve

```
pred_y_test = dt2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_dt_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_dt_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.') # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.821 auc_pr=0.828 ap=0.857



```
In [262]: models.append('DT')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_dt)
```

3) RandomForest Classifier

```
In [263]: from sklearn.ensemble import RandomForestClassifier
rf1 = RandomForestClassifier()
```

```
In [264]: rf1 = RandomForestClassifier(random_state=0)
```

```
In [265]: rf1.fit(X_train, y_train)
```

```
Out[265]: RandomForestClassifier
          RandomForestClassifier(random_state=0)
```

```
In [266]: rf1.score(X_train, y_train)      # Random Forest also 100% accuracy over train data always
```

```
Out[266]: 1.0
```

```
In [267]: rf1.score(X_test, y_test)
```

```
Out[267]: 0.8266666666666667
```

**Performance evaluation and optimizing parameters using GridSearchCV:

```
In [268]: parameters = {
    'n_estimators': [50,100,150],
    'max_depth': [None,1,3,5,7],
    'min_samples_leaf': [1,3,5]
}
```

```
In [269]: gs_dt = GridSearchCV(estimator=rf1, param_grid=parameters, cv=5, verbose=0)
gs_dt.fit(df_X_resampled, df_y_resampled)
```

```
Out[269]:
```

```
GridSearchCV
  estimator: RandomForestClassifier
    RandomForestClassifier
```

```
In [270]: gs_dt.best_params_
```

```
Out[270]: {'max_depth': None, 'min_samples_leaf': 1, 'n_estimators': 100}
```

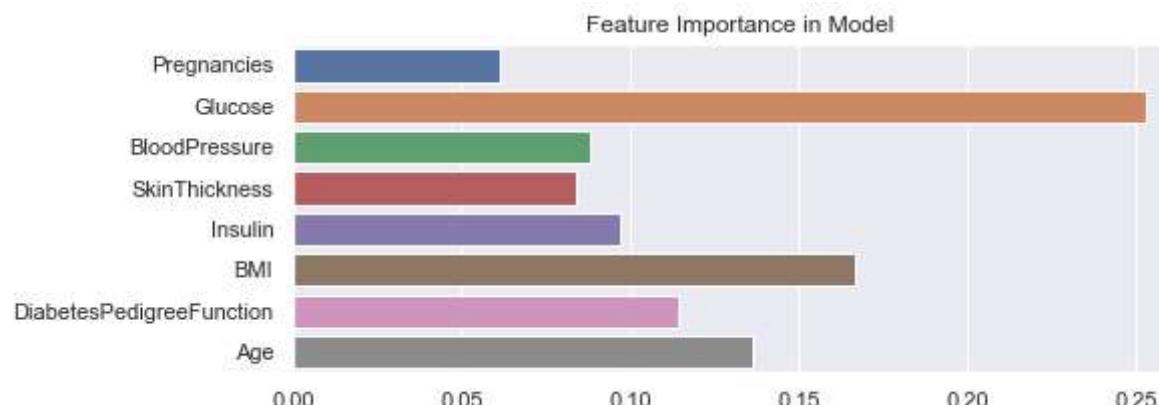
```
In [271]: gs_dt.best_score_
```

```
Out[271]: 0.819
```

```
In [272]: rf1.feature_importances_
```

```
Out[272]: array([0.0614086 , 0.25273886, 0.08832161, 0.08389693, 0.09701324,
       0.16633365, 0.11402855, 0.13625856])
```

```
In [273]: plt.figure(figsize=(8,3))
sns.barplot(y=X_train.columns, x=rf1.feature_importances_);
plt.title("Feature Importance in Model");
```



In [274]:

```
rf2 = RandomForestClassifier(max_depth=None, min_samples_leaf=1, n_estimators=100)
```

In [275]: rf2.fit(X_train,y_train)

Out[275]:

```
RandomForestClassifier  
| RandomForestClassifier()
```

In [276]: rf2.score(X_train,y_train)

Out[276]: 1.0

In [277]: rf2.score(X_test, y_test)

Out[277]: 0.8266666666666667

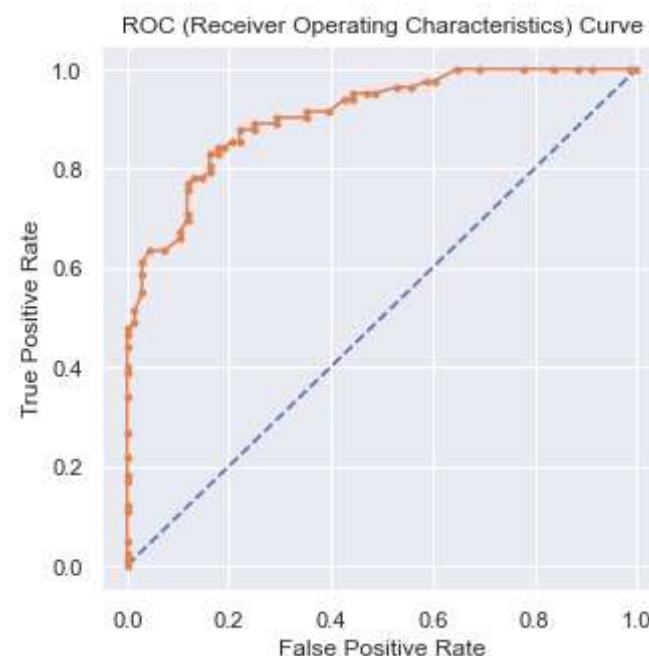
In [278]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)

```
probs = rf2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                      # keep probabilities for the positive outcome only

auc_rf = roc_auc_score(y_test, probs)      # calculate AUC
print('AUC: %.3f' %auc_rf)

fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')       # plot no skill
plt.plot(fpr, tpr, marker='.')                # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

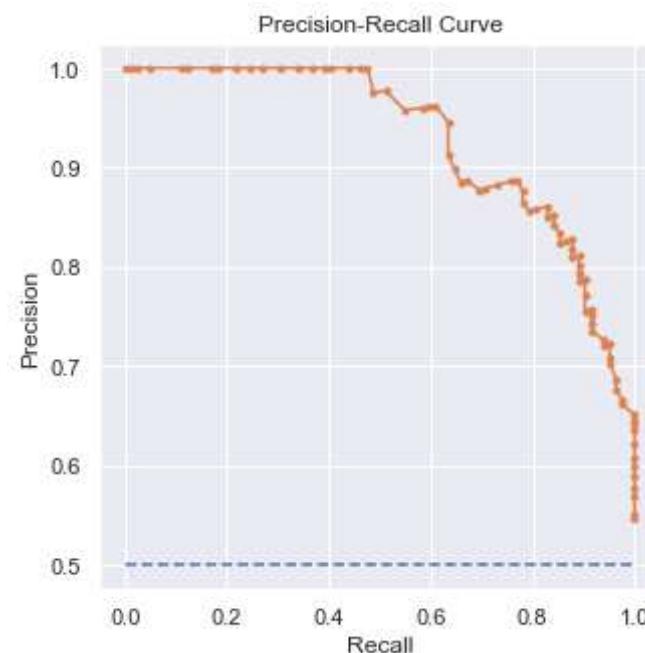
AUC: 0.907



In [279]: # Precision Recall Curve

```
pred_y_test = rf2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_rf_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_rf_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.') # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.841 auc_pr=0.929 ap=0.928



```
In [280]: models.append('RF')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_dt)
```

4) K-Nearest Neighbour (KNN) Classification:

```
In [281]: from sklearn.neighbors import KNeighborsClassifier
knn1 = KNeighborsClassifier(n_neighbors=3)
```

```
In [282]: knn1.fit(X_train, y_train)
```

```
Out[282]:
```

```
    KNeighborsClassifier
    KNeighborsClassifier(n_neighbors=3)
```

```
In [283]: knn1.score(X_train,y_train)
```

```
Out[283]: 0.8776470588235294
```

```
In [284]: knn1.score(X_test,y_test)
```

```
Out[284]: 0.7666666666666667
```

**Performance evaluation and optimizing parameters using GridSearchCV:

```
In [285]: knn_neighbors = [i for i in range(2,16)]
parameters = {
    'n_neighbors': knn_neighbors
}
```

```
In [286]: gs_knn = GridSearchCV(estimator=knn1, param_grid=parameters, cv=5, verbose=0)
gs_knn.fit(df_X_resampled, df_y_resampled)
```

```
Out[286]:
```

GridSearchCV

```
GridSearchCV(cv=5, estimator=KNeighborsClassifier(n_neighbors=3),
             param_grid={'n_neighbors': [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13,
                                         14, 15]})

    ▶ estimator: KNeighborsClassifier
        ▶ KNeighborsClassifier
```

```
In [287]: gs_knn.best_params_
```

```
Out[287]: {'n_neighbors': 3}
```

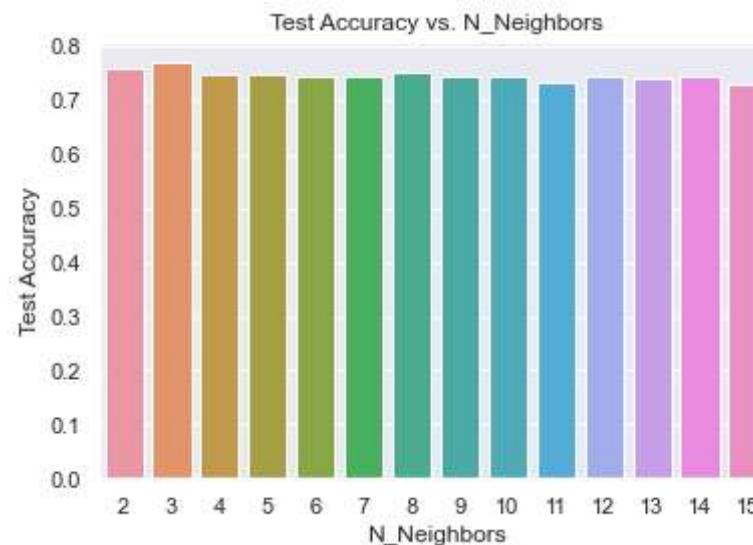
```
In [288]: gs_knn.best_score_
```

```
Out[288]: 0.7689999999999999
```

```
In [289]: # gs_knn.cv_results_
gs_knn.cv_results_['mean_test_score']
```

```
Out[289]: array([0.758, 0.769, 0.749, 0.748, 0.743, 0.745, 0.752, 0.743, 0.742,
                 0.734, 0.744, 0.74 , 0.742, 0.73 ])
```

```
In [290]: plt.figure(figsize=(6,4))
sns.barplot(x=knn_neighbors, y=gs_knn.cv_results_['mean_test_score'])
plt.xlabel("N_Neighbors")
plt.ylabel("Test Accuracy")
plt.title("Test Accuracy vs. N_Neighbors");
```



```
In [291]: knn2 = KNeighborsClassifier(n_neighbors=3)
```

```
In [292]: knn2.fit(X_train, y_train)
```

```
Out[292]: KNeighborsClassifier
```

```
KNeighborsClassifier(n_neighbors=3)
```

```
In [293]: knn2.score(X_train,y_train)
```

```
Out[293]: 0.8776470588235294
```

```
In [294]: knn2.score(X_test,y_test)
```

```
Out[294]: 0.7666666666666667
```

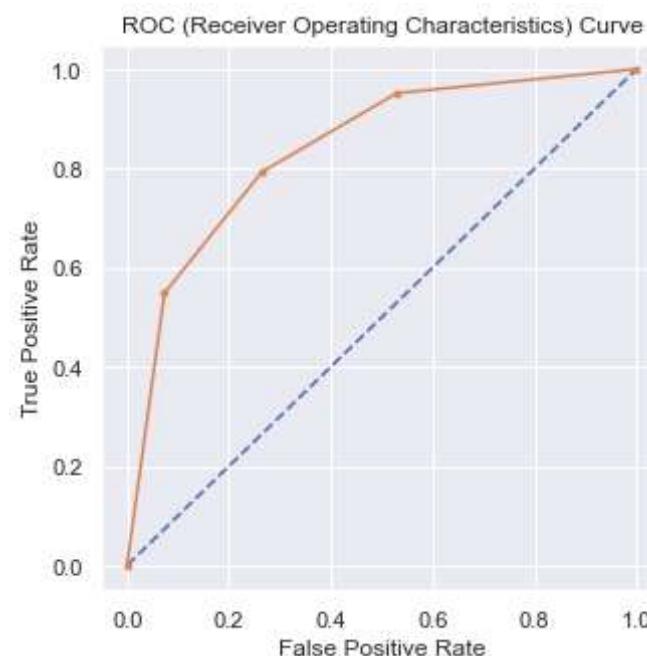
In [295]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)

```
probs = knn2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                         # keep probabilities for the positive outcome only

auc_knn = roc_auc_score(y_test, probs)       # calculate AUC
print('AUC: %.3f' %auc_knn)

fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')        # plot no skill
plt.plot(fpr, tpr, marker='.')                 # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

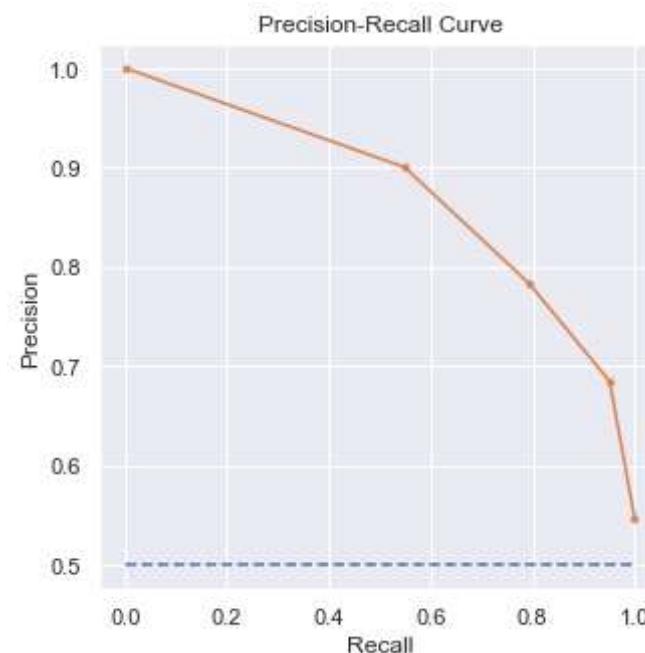
AUC: 0.838



In [296]: # Precision Recall Curve

```
pred_y_test = knn2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_knn_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_knn_pr, ap)) # plot no skill
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot the precision-recall curve for the model
plt.plot(recall, precision, marker='.')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.788 auc_pr=0.873 ap=0.820



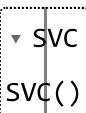
```
In [297]: models.append('KNN')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_knn)
```

5) Support Vector Machine (SVM) Algorithm:

```
In [298]: from sklearn.svm import SVC
svm1 = SVC(kernel='rbf')
```

```
In [299]: svm1.fit(X_train, y_train)
```

```
Out[299]:
```



```
SVC()
```

```
In [300]: svm1.score(X_train, y_train)
```

```
Out[300]: 0.7352941176470589
```

```
In [144]: svm1.score(X_test, y_test)
```

```
Out[144]: 0.7
```

**Performance evaluation and optimizing parameters using GridSearchCV:

```
In [301]: parameters = {
    'C':[1, 5, 10, 15, 20, 25],
    'gamma':[0.001, 0.005, 0.0001, 0.00001]
}
```

```
In [302]: gs_svm = GridSearchCV(estimator=svm1, param_grid=parameters, cv=5, verbose=0)
gs_svm.fit(df_X_resampled, df_y_resampled)
```

```
Out[302]:
```

```
GridSearchCV
GridSearchCV(cv=5, estimator=SVC(),
             param_grid={'C': [1, 5, 10, 15, 20, 25],
                         'gamma': [0.001, 0.005, 0.0001, 1e-05]})
  ▶ estimator: SVC
    ▶ SVC
```

```
In [303]: gs_svm.best_params_
```

```
Out[303]: {'C': 5, 'gamma': 0.005}
```

```
In [304]: gs_svm.best_score_
```

```
Out[304]: 0.808
```

```
In [305]: svm2 = SVC(kernel='rbf', C=5, gamma=0.005, probability=True)
```

```
In [306]: svm2.fit(X_train, y_train)
```

```
Out[306]:
```

```
SVC
SVC(C=5, gamma=0.005, probability=True)
```

```
In [307]: svm2.score(X_train, y_train)
```

```
Out[307]: 0.9741176470588235
```

```
In [308]: svm2.score(X_test, y_test)
```

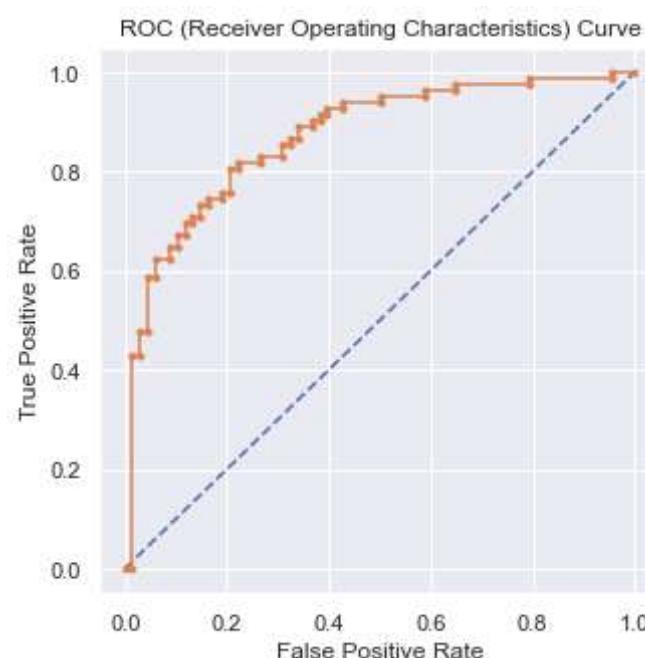
```
Out[308]: 0.7866666666666666
```

```
In [309]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)
```

```
probs = svm2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                         # keep probabilities for the positive outcome only

auc_svm = roc_auc_score(y_test, probs)       # calculate AUC
print('AUC: %.3f' %auc_svm)
fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')      # plot no skill
plt.plot(fpr, tpr, marker '.')               # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

AUC: 0.871

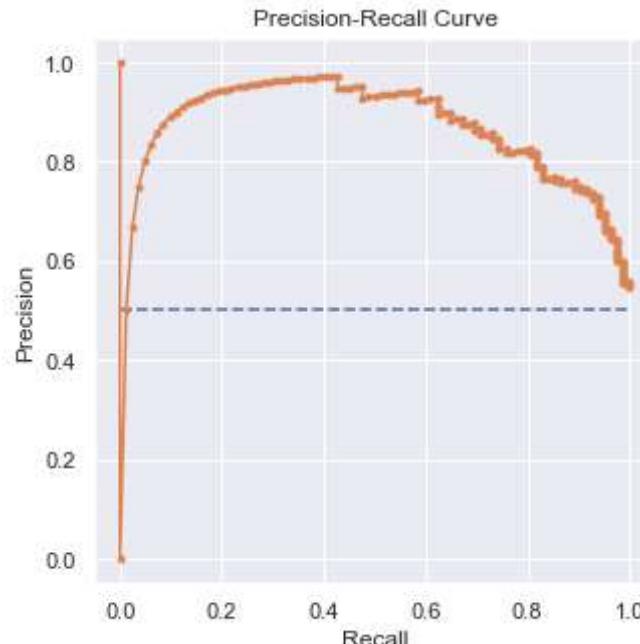


In [310]: # Precision Recall Curve

```
pred_y_test = svm2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_svm_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_svm_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.') # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
```

f1=0.800 auc_pr=0.862 ap=0.869

Out[310]: Text(0.5, 1.0, 'Precision-Recall Curve')



```
In [311]: models.append('SVM')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_svm)
```

6) Naive Bayes Algorithm:

```
In [312]: from sklearn.naive_bayes import GaussianNB, BernoulliNB, MultinomialNB
gnb = GaussianNB()
```

```
In [313]: gnb.fit(X_train, y_train)
```

```
Out[313]:
```

▼ GaussianNB
GaussianNB()

```
In [314]: gnb.score(X_train, y_train)
```

```
Out[314]: 0.7282352941176471
```

```
In [315]: gnb.score(X_test, y_test)
```

```
Out[315]: 0.72
```

****Naive Bayes has almost no hyperparameters to tune, so it usually generalizes well.**

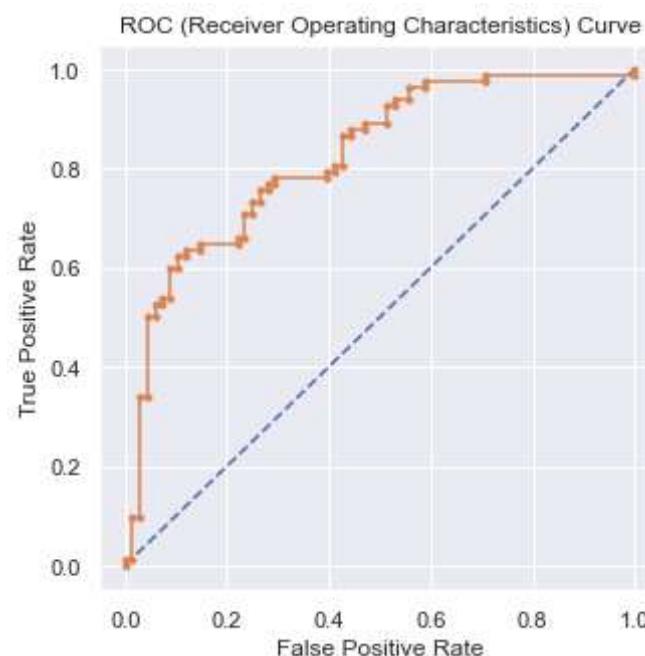
In [316]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)

```
probs = gnb.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                      # keep probabilities for the positive outcome only

auc_gnb = roc_auc_score(y_test, probs)    # calculate AUC
print('AUC: %.3f' %auc_gnb)

fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')      # plot no skill
plt.plot(fpr, tpr, marker='.')               # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

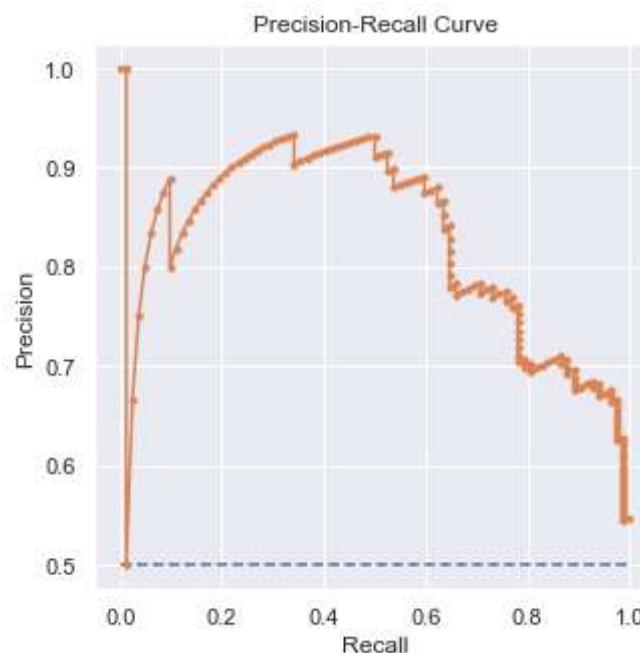
AUC: 0.824



In [317]: # Precision Recall Curve

```
pred_y_test = gnb.predict(X_test)          # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test)         # calculate F1 score
auc_gnb_pr = auc(recall, precision)       # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_gnb_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.')    # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.727 auc_pr=0.824 ap=0.828



In [318]: models.append('GNB')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_gnb)

7) Ensemble Learning --> Boosting --> Adaptive Boosting:

```
In [319]: from sklearn.ensemble import AdaBoostClassifier  
ada1 = AdaBoostClassifier(n_estimators=100)
```

```
In [320]: ada1.fit(X_train,y_train)
```

```
Out[320]:
```

```
▼       AdaBoostClassifier  
AdaBoostClassifier(n_estimators=100)
```

```
In [321]: ada1.score(X_train,y_train)
```

```
Out[321]: 0.84
```

```
In [322]: ada1.score(X_test, y_test)
```

```
Out[322]: 0.8
```

```
**Performance evaluation and optimizing parameters using cross_val_score:
```

```
In [323]: parameters = {'n_estimators': [100,200,300,400,500,700,1000]}
```

```
In [324]: gs_ada = GridSearchCV(ada1, param_grid = parameters, cv=5, verbose=0)  
gs_ada.fit(df_X_resampled, df_y_resampled)
```

```
Out[324]:
```

```
►       GridSearchCV  
  ► estimator: AdaBoostClassifier  
    ► AdaBoostClassifier
```

```
In [325]: gs_ada.best_params_
```

```
Out[325]: {'n_estimators': 200}
```

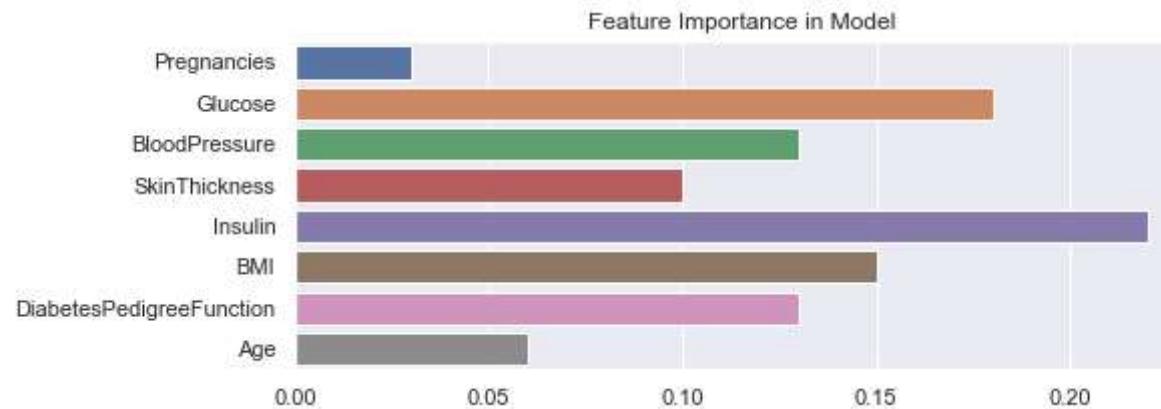
```
In [326]: gs_ada.best_score_
```

```
Out[326]: 0.766
```

```
In [327]: ada1.feature_importances_
```

```
Out[327]: array([0.03, 0.18, 0.13, 0.1 , 0.22, 0.15, 0.13, 0.06])
```

```
In [328]: plt.figure(figsize=(8,3))
sns.barplot(y=X_train.columns, x=ada1.feature_importances_)
plt.title("Feature Importance in Model");
```



```
In [329]: ada2 = AdaBoostClassifier(n_estimators=200)
```

```
In [330]: ada2.fit(X_train,y_train)
```

```
Out[330]:
```

```
AdaBoostClassifier  
AdaBoostClassifier(n_estimators=200)
```

```
In [331]: ada2.score(X_train,y_train)
```

```
Out[331]: 0.8776470588235294
```

```
In [332]: ada2.score(X_test, y_test)
```

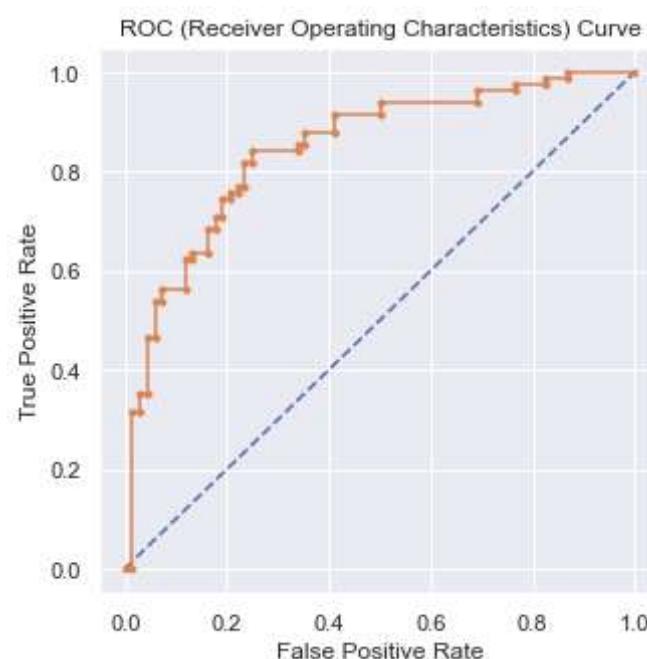
```
Out[332]: 0.7866666666666666
```

```
In [333]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)
```

```
probs = ada2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                        # keep probabilities for the positive outcome only

auc_ada = roc_auc_score(y_test, probs)      # calculate AUC
print('AUC: %.3f' %auc_ada)
fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')      # plot no skill
plt.plot(fpr, tpr, marker='.')               # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

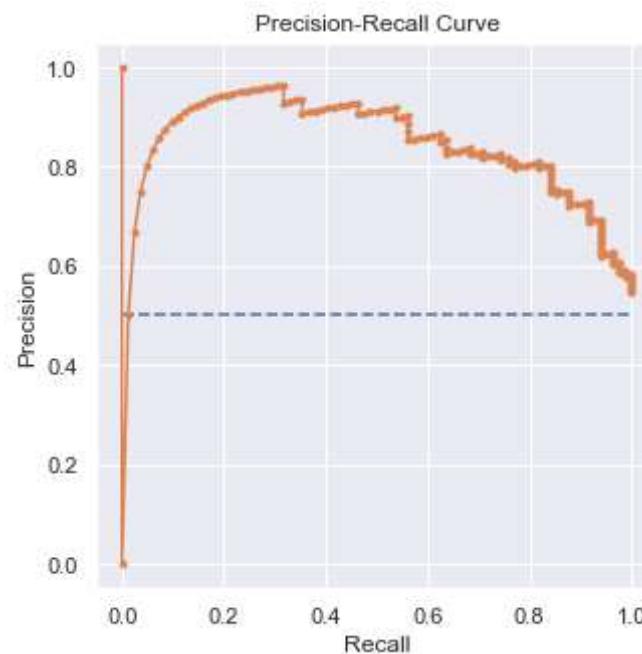
AUC: 0.846



In [334]: # Precision Recall Curve

```
pred_y_test = ada2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_ada_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_ada_pr, ap))
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot no skill
plt.plot(recall, precision, marker='.') # plot the precision-recall curve for the model
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

f1=0.805 auc_pr=0.839 ap=0.845



```
In [335]: models.append('ADA')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_ada)
```

8) Ensemble Learning --> Boosting --> Gradient Boosting (XGBClassifier):

```
In [336]: from xgboost import XGBClassifier
xgb1 = XGBClassifier(use_label_encoder=False, objective = 'binary:logistic', nthread=4, seed=10)
```

```
In [337]: xgb1.fit(X_train, y_train)
```

```
Out[337]: XGBClassifier
XGBClassifier(base_score=0.5, booster='gbtree', callbacks=None,
              colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,
              early_stopping_rounds=None, enable_categorical=False,
              eval_metric=None, gamma=0, gpu_id=-1, grow_policy='depthwise',
              importance_type=None, interaction_constraints='',
              learning_rate=0.300000012, max_bin=256, max_cat_to_onehot=4,
              max_delta_step=0, max_depth=6, max_leaves=0, min_child_weight=1,
              missing=nan, monotone_constraints='()', n_estimators=100,
              n_jobs=4, nthread=4, num_parallel_tree=1, predictor='auto',
              random_state=10, reg_alpha=0, ...)
```

```
In [338]: xgb1.score(X_train, y_train)
```

```
Out[338]: 1.0
```

```
In [339]: xgb1.score(X_test, y_test)
```

```
Out[339]: 0.8133333333333334
```

**Performance evaluation and optimizing parameters using GridSearchCV:

```
In [340]: parameters = {
    'max_depth': range(2, 10, 1),
    'n_estimators': range(60, 220, 40),
    'learning_rate': [0.1, 0.01, 0.05]
}
```

```
In [341]: gs_xgb = GridSearchCV(xgb1, param_grid = parameters, scoring = 'roc_auc', n_jobs = 10, cv=5, verbose=0)
gs_xgb.fit(df_X_resampled, df_y_resampled)
```

Out[341]:

```
GridSearchCV
GridSearchCV(cv=5,
            estimator=XGBClassifier(base_score=0.5, booster='gbtree',
                                    callbacks=None, colsample_bylevel=1,
                                    colsample_bynode=1, colsample_bytree=1,
                                    early_stopping_rounds=None,
                                    enable_categorical=False, eval_metric=None,
                                    gamma=0, gpu_id=-1,
                                    grow_policy='depthwise',
                                    importance_type=None,
                                    interaction_constraints='',
            estimator: XGBClassifier
                XGBClassifier
```

```
In [342]: gs_xgb.best_params_
```

```
Out[342]: {'learning_rate': 0.05, 'max_depth': 7, 'n_estimators': 180}
```

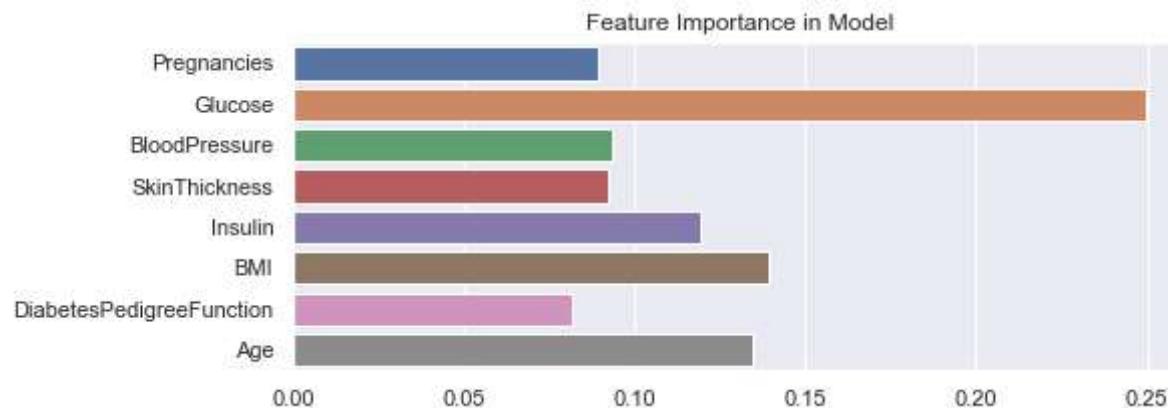
```
In [343]: gs_xgb.best_score_
```

```
Out[343]: 0.8767400000000001
```

```
In [344]: xgb1.feature_importances_
```

```
Out[344]: array([0.08949539, 0.24965179, 0.09345236, 0.09250528, 0.11924399,
       0.13939798, 0.08189643, 0.13435684], dtype=float32)
```

```
In [345]: plt.figure(figsize=(8,3))
sns.barplot(y=X_train.columns, x=xgb1.feature_importances_)
plt.title("Feature Importance in Model");
```



```
In [346]: xgb2 = XGBClassifier(use_label_encoder=False, objective = 'binary:logistic',
                           nthread=4, seed=10, learning_rate= 0.05, max_depth= 7, n_estimators= 180)
```

In [347]: `xgb2.fit(X_train,y_train)`

Out[347]:

```
XGBClassifier  
XGBClassifier(base_score=0.5, booster='gbtree', callbacks=None,  
             colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,  
             early_stopping_rounds=None, enable_categorical=False,  
             eval_metric=None, gamma=0, gpu_id=-1, grow_policy='depthwise',  
             importance_type=None, interaction_constraints='',  
             learning_rate=0.05, max_bin=256, max_cat_to_onehot=4,  
             max_delta_step=0, max_depth=7, max_leaves=0, min_child_weight=1,  
             missing=nan, monotone_constraints='()', n_estimators=180,  
             n_jobs=4, nthread=4, num_parallel_tree=1, predictor='auto',  
             random_state=10, reg_alpha=0, ...)
```

In [348]: `xgb2.score(X_train,y_train)`

Out[348]: 0.9952941176470588

In [349]: `xgb2.score(X_test, y_test)`

Out[349]: 0.84

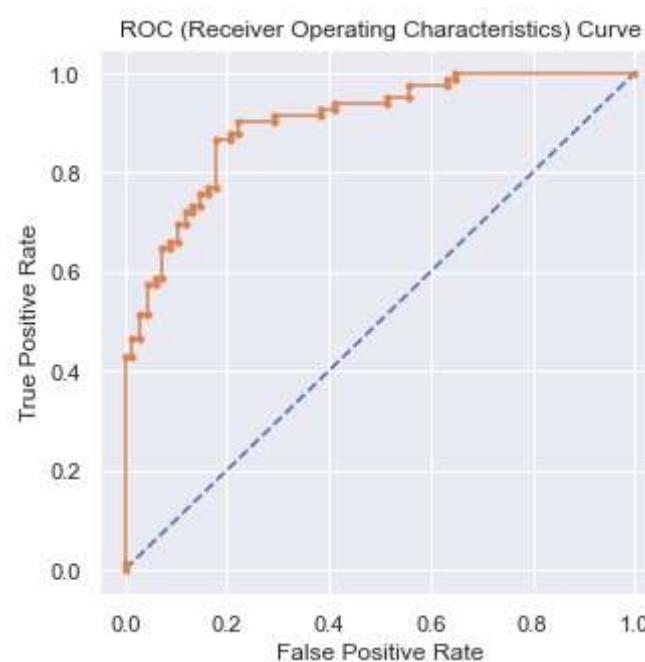
In [350]: # Preparing ROC Curve (Receiver Operating Characteristics Curve)

```
probs = xgb2.predict_proba(X_test)          # predict probabilities
probs = probs[:, 1]                         # keep probabilities for the positive outcome only

auc_xgb = roc_auc_score(y_test, probs)       # calculate AUC
print('AUC: %.3f' %auc_xgb)

fpr, tpr, thresholds = roc_curve(y_test, probs) # calculate roc curve
plt.plot([0, 1], [0, 1], linestyle='--')        # plot no skill
plt.plot(fpr, tpr, marker='.')                 # plot the roc curve for the model
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC (Receiver Operating Characteristics) Curve");
```

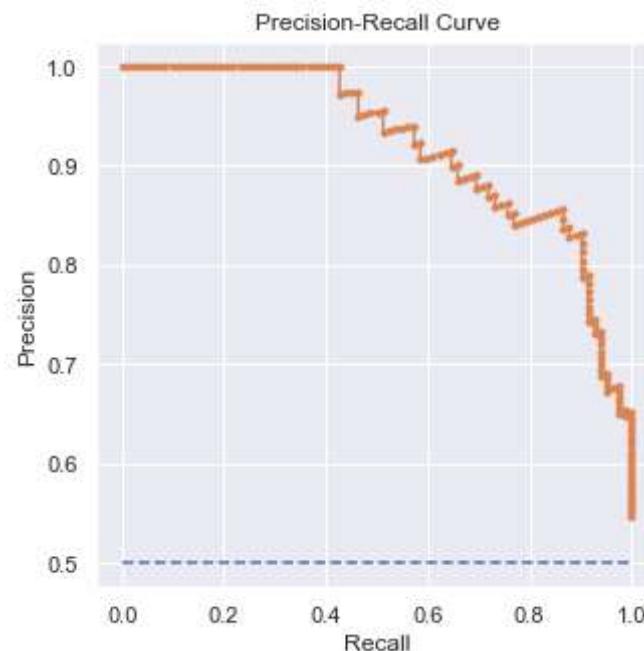
AUC: 0.901



In [351]: # Precision Recall Curve

```
pred_y_test = xgb2.predict(X_test) # predict class values
precision, recall, thresholds = precision_recall_curve(y_test, probs) # calculate precision-recall curve
f1 = f1_score(y_test, pred_y_test) # calculate F1 score
auc_xgb_pr = auc(recall, precision) # calculate precision-recall AUC
ap = average_precision_score(y_test, probs) # calculate average precision score
print('f1=% .3f auc_pr=% .3f ap=% .3f' % (f1, auc_xgb_pr, ap)) # plot no skill
plt.plot([0, 1], [0.5, 0.5], linestyle='--') # plot the precision-recall curve for the model
plt.plot(recall, precision, marker='.')
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve");
```

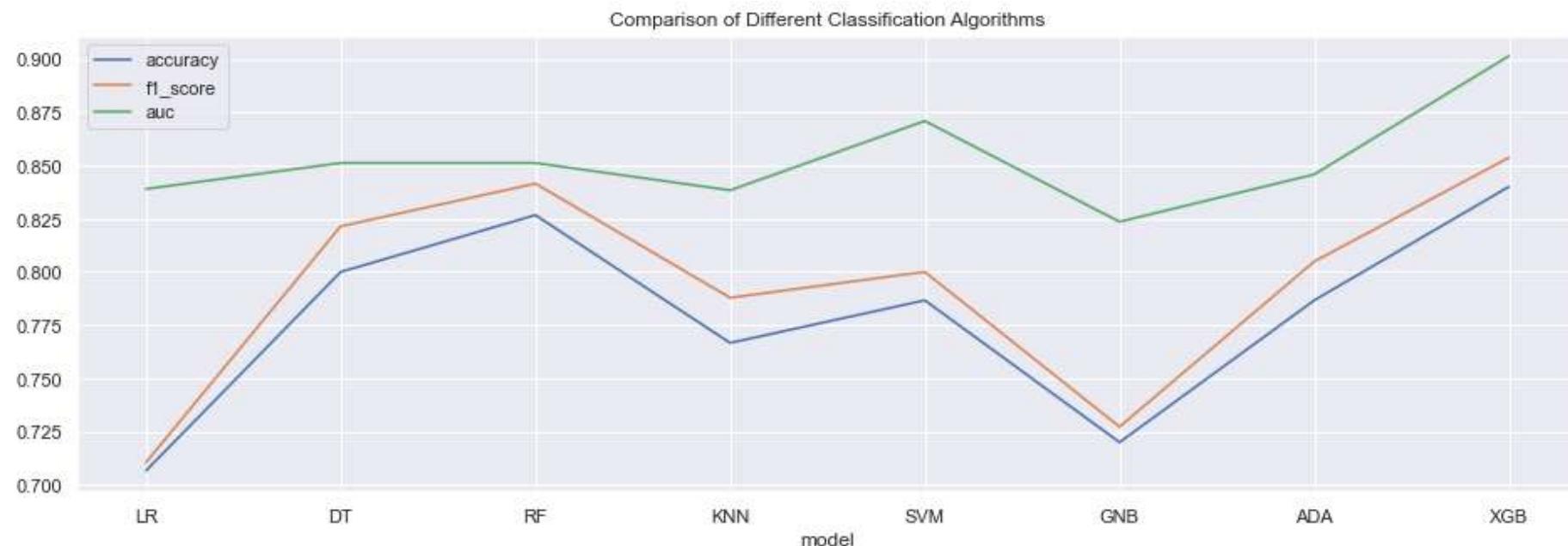
f1=0.854 auc_pr=0.921 ap=0.922



```
In [352]: models.append('XGB')
model_accuracy.append(accuracy_score(y_test, pred_y_test))
model_f1.append(f1)
model_auc.append(auc_xgb)
```

```
In [353]: model_summary = pd.DataFrame(zip(models,model_accuracy,model_f1,model_auc), columns = ['model','accuracy','f1_score','auc'])
model_summary = model_summary.set_index('model')
```

```
In [354]: model_summary.plot(figsize=(16,5))
plt.title("Comparison of Different Classification Algorithms");
```



```
In [356]: model_summary
```

Out[356]:

model	accuracy	f1_score	auc
LR	0.706667	0.710526	0.838953
DT	0.800000	0.821429	0.851148
RF	0.826667	0.841463	0.851148
KNN	0.766667	0.787879	0.838325
SVM	0.786667	0.800000	0.870875
GNB	0.720000	0.727273	0.823529
ADA	0.786667	0.804878	0.845768
XGB	0.840000	0.853659	0.901363

**Among all models, RandomForest has given best accuracy and f1_score. Therefore we will build final model using RandomForest.

**FINAL CLASSIFIER:

```
In [357]: final_model = rf2
```

**Week 4:

Data Modeling:

(1) Create a classification report by analyzing sensitivity, specificity, AUC (ROC curve), etc. Please be descriptive to explain what values of these parameter you have used:

```
In [358]: cr = classification_report(y_test, final_model.predict(X_test))
print(cr)
```

	precision	recall	f1-score	support
0	0.81	0.81	0.81	68
1	0.84	0.84	0.84	82
accuracy			0.83	150
macro avg	0.83	0.83	0.83	150
weighted avg	0.83	0.83	0.83	150

```
In [359]: confusion = confusion_matrix(y_test, final_model.predict(X_test))
print("Confusion Matrix:\n", confusion)
```

Confusion Matrix:
[[55 13]
[13 69]]

```
In [360]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives

Accuracy = (TP+TN)/(TP+TN+FP+FN)
Precision = TP/(TP+FP)
Sensitivity = TP/(TP+FN) # also called recall
Specificity = TN/(TN+FP)
```

```
In [361]: print("Accuracy: %.3f"%Accuracy)
print("Precision: %.3f"%Precision)
print("Sensitivity: %.3f"%Sensitivity)
print("Specificity: %.3f"%Specificity)
print("AUC: %.3f"%auc_rf)
```

Accuracy: 0.827
Precision: 0.841
Sensitivity: 0.841
Specificity: 0.809
AUC: 0.907

****Sensitivity and Specificity:**** By changing the threshold, target classification will be changed hence the sensitivity and specificity will also be changed. Which one of these two we should maximize? What should be ideal threshold?

Ideally we want to maximize both Sensitivity & Specificity. But this is not possible always. There is always a trade-off. Sometimes we want to be 100% sure on Predicted negatives, sometimes we want to be 100% sure on Predicted positives. Sometimes we simply don't want to compromise on sensitivity sometimes we don't want to compromise on specificity.

The threshold is set based on business problem. There are some cases where Sensitivity is important and need to be nearer to 1. There are business cases where Specificity is important and need to be nearer to 1. We need to understand the business problem and decide the importance of Sensitivity and Specificity.

Thakur

-Sourabh Singh