


```
[32]: #Checking X values
X
array([[ 26.,  12., 17., ...,  0.,  1.,  0.],
       [ 95.,  44.,  28., ...,  1.,  0.,  0.],
       [  0.,  16.,  12., ...,  0.,  1.,  0.],
       [112.,  39.,  19., ...,  0.,  1.,  0.],
       [ 32.,  35.,  27., ...,  0.,  1.,  0.],
       [106.,  17.,  23., ...,  1.,  0.,  0.]])
```

Scaling

Scaling is done with MinMax scaler which scales the data with respect to the minimum and maximum values, this is better than the standard scaler since it will mostly incline towards the mean of the data.

```
In [33]: # Scaling the dataset
#through histogram we see that data is not much normally distributed for some of the columns.
#Hence we use Standard Scaler to normalize and scale
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

In [34]: Classification_Scores=()
```

1) Voting Classifier

Voting is an ensemble classifier which is a combination from multiple machine learning algorithms. It helps in lowering error and lessen overfitting

```
In [35]: from sklearn.ensemble import VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
```

Hard voting : Case where the class with more votes are selected

a) Voting Classifier - Hard

Logistic regression is run with $c=1$, penalty set as l2 and liblinear solver is used.

SVM poly kernel is run with $\gamma=0.1$, $c=1$

KNN with neighbors = 9 is run

These models are used with voting to find a best model with highest votes

```
In [36]: logistic_clf = LogisticRegression(C=1, penalty='l2', solver='liblinear')
svc_clf = SVC(kernel='poly', gamma=0.1, C=1, probability=True)
knn_clf = KNeighborsClassifier(9)
```

Train data is fit into the different models without voting

```
In [37]: logistic_clf.fit(X_train, y_train)
svc_clf.fit(X_train, y_train)
#decisiontree_clf.fit(X_train, y_train)
knn_clf.fit(X_train, y_train)

Out[37]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=9, p=2,
                             weights='uniform')
```

Voting is done with the three models

```
In [38]: voting = VotingClassifier(estimators= [('l1', logistic_clf), ('svc', svc_clf), ('kn', knn_clf)], voting
                                = 'hard')
```

The scores of the training set of the models are shown:

```
In [39]: print('log_clf ', logistic_clf.score(X_train, y_train))
print('svc_clf ', svc_clf.score(X_train, y_train))
print('kn_clf ', knn_clf.score(X_train, y_train))

log_clf: 0.9461139896373058
svc_clf: 0.9471502590673575
knn_clf: 0.93860103626943
```

The train and test score for the voting classifier is obtained

```
In [40]: voting.fit(X_train, y_train)
print('vot_clf Train: (0.46)'.format(voting.score(X_train, y_train)))
print('vot_clf Test: (0.46)'.format(voting.score(X_test, y_test)))

vot_clf Train: 0.9482
vot_clf Test: 0.9380
```

Confusion matrix is printed with values of the voting classifier

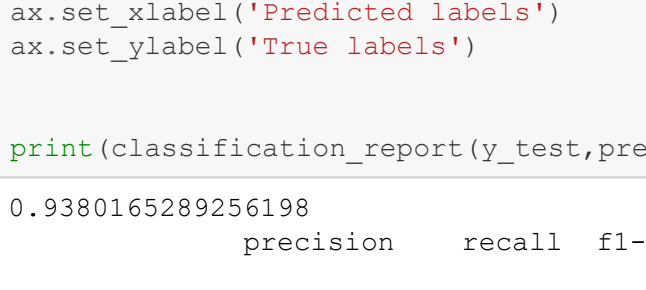
```
In [41]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_hardvotingclf = voting.predict(X_test)
print(metrics.accuracy_score(y_test, pred_hardvotingclf))

confusion = confusion_matrix(y_test, pred_hardvotingclf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
print(classification_report(y_test, pred_hardvotingclf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.92	0.91	0.92	91
accuracy			0.94	242
macro avg	0.93	0.93	0.93	242
weighted avg	0.94	0.94	0.94	242



The F1 score and Recall is obtained

```
In [42]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_hardvotingclf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9536423841059603
F1Score : 0.9504950495049505
```

```
In [43]: Classification_Scores.update({'Hard voting Classifier': [metrics.accuracy_score(y_test, pred_hardvotingclf),
                                                                recall[0], fscore[0]]})
```

b) Voting Classifier - Soft

Soft voting is done with probability vector of each class summed and averaged

```
In [44]: soft_voting = VotingClassifier(estimators= [('l1', logistic_clf), ('svc', svc_clf), ('kn', knn_clf)], voting
                                      = 'soft')
```

The train and test of the soft voting classifier is obtained

```
In [45]: soft_voting.fit(X_train, y_train)
print('vot_clf Train: (0.46)'.format(soft_voting.score(X_train, y_train)))
print('vot_clf Test: (0.46)'.format(soft_voting.score(X_test, y_test)))

vot_clf Train: 0.9451
vot_clf Test: 0.9380
```

The confusion matrix of the soft voting classifier is printed

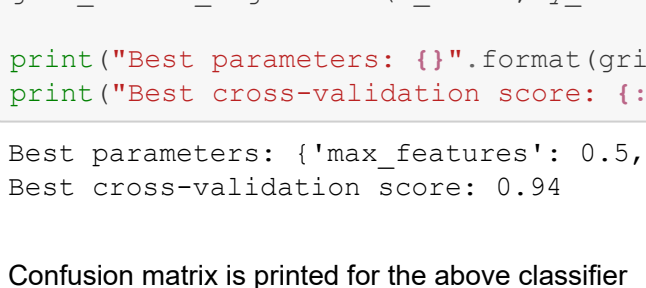
```
In [46]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_softvotingclf = soft_voting.predict(X_test)
print(metrics.accuracy_score(y_test, pred_softvotingclf))

confusion = confusion_matrix(y_test, pred_softvotingclf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
print(classification_report(y_test, pred_softvotingclf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.92	0.91	0.92	91
accuracy			0.94	242
macro avg	0.93	0.93	0.93	242
weighted avg	0.94	0.94	0.94	242



The F1 score and Recall values are obtained

```
In [47]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_softvotingclf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9536423841059603
F1Score : 0.9504950495049505
```

```
In [48]: Classification_Scores.update({'Soft voting classifier': [metrics.accuracy_score(y_test, pred_softvotingclf),
                                                                recall[0], fscore[0]]})
```

2)Bagging

Bagging is an aggregation ensemble meta-algorithm designed to improve the stability and accuracy of algorithms.

It helps to reduce overfitting. One main parameter passed is the bootstrap parameter. It uses sampling with replacement.

$N_{estimators}$ and $max_samples$ are parameters that control the number of decision trees and sample size used.

The bagging is done with KNN and decision Tree

Bagging-Decision tree

Bagging is run with decision tree with $n_{estimator}=500$, $max_features=0.5$

```
In [49]: from sklearn.ensemble import BaggingClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import BaggingClassifier
bag_clf = BaggingClassifier(DecisionTreeClassifier(), bootstrap=True, n_jobs=-1, random_state=0, max_features=0.5, n_estimators=500, max_samples=0.6)
```

```
bag_clf.fit(X_train, y_train)
print('Accuracy on training set: ({}).format(bag_clf.score(X_train, y_train))')
print('Accuracy on test set: ({}).format(bag_clf.score(X_test, y_test))')

Accuracy on training set: 0.977
Accuracy on test set: 0.942
```

Cross validation score for the bagging classifier

```
In [50]: scores = cross_val_score(bag_clf, X_train, y_train, cv = 10, scoring = 'accuracy' )
print('Cross-validation scores: ({}).format(scores)')

print('Average cross-validation score: ({}).format(scores.mean())')

Cross-validation scores: [0.93814433 0.95876289 0.88659794 0.96907216 0.96907216 0.91666667 0.94791667 0.88541667 0.97916667 0.94791667]
Average cross-validation score: 0.94
```

Grid Search is done with $max_samples$, $max_features$ and $n_{estimators}$ to find the best parameters

```
In [51]: param_grid = {'max_samples': [0.5, 0.6, 0.8],
                    'max_features': [0.4, 0.5, 0.6],
                    'n_estimators': [100, 300, 500]}

grid_search_bagclf = GridSearchCV(BaggingClassifier(random_state=0, bootstrap=True, n_jobs=-1), param_grid, cv=5, return_train_score=True)
grid_search_bagclf.fit(X_train, y_train)
```

```
print('Best parameters: ({}).format(grid_search_bagclf.best_params_)')
print('Best cross-validation score: ({}).format(grid_search_bagclf.best_score_)')

Best parameters: {'max_features': 0.5, 'max_samples': 0.6, 'n_estimators': 500}
Best cross-validation score: 0.94
```

Confusion matrix is printed for the above classifier

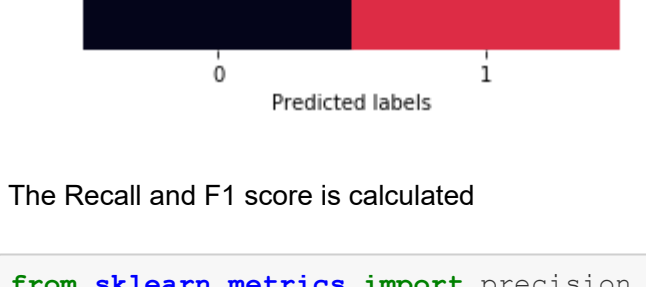
```
In [52]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_bagclf = grid_search_bagclf.predict(X_test)
print(metrics.accuracy_score(y_test, pred_bagclf))

confusion = confusion_matrix(y_test, pred_bagclf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
#ax.xaxis.set_ticklabels(['Attrition', 'No-Attrition'])
#ax.yaxis.set_ticklabels(['Attrition', 'No-Attrition'])
print(classification_report(y_test, pred_bagclf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.92	0.92	0.92	91
accuracy			0.94	242
macro avg	0.94	0.94	0.94	242
weighted avg	0.94	0.94	0.94	242



The Recall and F1 score for the decision tree with bagging is obtained

```
In [53]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_bagclf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9536423841059603
F1Score : 0.9536423841059603
```

```
In [54]: Classification_Scores.update({'Bagging-Decision Tree Classifier': [metrics.accuracy_score(y_test, pred_bagclf),
                                                                           recall[0], fscore[0]]})
```

Bagging - KNN

The KNN model is run with neighbors = 4 with bagging

```
In [55]: from sklearn.ensemble import BaggingClassifier
bagknn_clf = BaggingClassifier(KNeighborsClassifier(n_neighbors=4), bootstrap=True, n_jobs=-1, random_state=0)
bagknn_clf.fit(X_train, y_train)
print('Accuracy on training set: ({}).format(bagknn_clf.score(X_train, y_train))')
print('Accuracy on test set: ({}).format(bagknn_clf.score(X_test, y_test))')

Accuracy on training set: 0.945
Accuracy on test set: 0.938
```

The cross validation score is calculated with cv=5

```
In [56]: scores = cross_val_score(bagknn_clf, X_train, y_train, cv = 5, scoring = 'accuracy' )
print('Cross-validation scores: ({}).format(scores)')

print('Average cross-validation score: ({}).format(scores.mean())')

Cross-validation scores: [0.93264249 0.89119171 0.94818653 0.919171 0.94818653]
Average cross-validation score: 0.93
```

Grid search is used to find the best parameters for max_sample , $max_features$ and $n_{estimators}$

```
In [57]: param_grid = {'max_samples': [0.5, 0.6],
                    'max_features': [0.4, 0.5],
                    'n_estimators': [100, 500]}

grid_search_bagclfknn = GridSearchCV(BaggingClassifier(KNeighborsClassifier(n_neighbors=4), bootstrap=True, n_jobs=-1), param_grid, cv=5, return_train_score=True)
grid_search_bagclfknn.fit(X_train, y_train)
```

```
print('Best parameters: ({}).format(grid_search_bagclfknn.best_params_)')
print('Best cross-validation score: ({}).format(grid_search_bagclfknn.best_score_)')

Best parameters: {'max_features': 0.5, 'max_samples': 0.6, 'n_estimators': 500}
Best cross-validation score: 0.94
```

Confusion matrix is printed

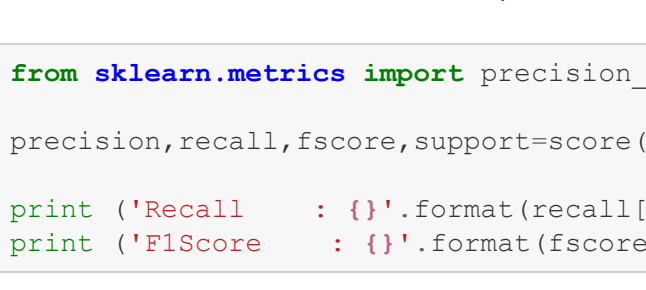
```
In [58]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_bagknn_clf = grid_search_bagclfknn.predict(X_test)
print(metrics.accuracy_score(y_test, pred_bagknn_clf))

confusion = confusion_matrix(y_test, pred_bagknn_clf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
#ax.xaxis.set_ticklabels(['Attrition', 'No-Attrition'])
#ax.yaxis.set_ticklabels(['Attrition', 'No-Attrition'])
print(classification_report(y_test, pred_bagknn_clf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.92	0.91	0.92	91
accuracy			0.94	242
macro avg	0.93	0.93	0.93	242
weighted avg	0.94	0.94	0.94	242



The Recall and F1 score is calculated

```
In [59]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_bagknn_clf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9470198675496688
F1Score : 0.9501661129568105
```

```
In [60]: Classification_Scores.update({'Bagging-KNN Classifier': [metrics.accuracy_score(y_test, pred_bagknn_clf),
                                                                recall[0], fscore[0]]})
```

3)Pasting

Pasting is an aggregation ensemble meta-algorithm designed to improve the stability and accuracy of algorithms. It helps to reduce overfitting. One main parameter passed is the bootstrap parameter which is set to false. It uses sampling without replacement.

$N_{estimators}$ and $max_samples$ are parameters that control the number of decision trees and sample size used.

Pasting is done on Decision Tree and KNN models

Pasting - Decision Tree

Decision tree is run with bootstrap=false for pasting

```
In [61]: from sklearn.ensemble import BaggingClassifier
pasteknn_clf = BaggingClassifier(DecisionTreeClassifier(random_state = 0, max_depth= 2), bootstrap=False, n_jobs=-1, random_state=0, n_estimators=150, max_samples=100)
pasteknn_clf.fit(X_train, y_train)
print('Accuracy on training set: ({}).format(pasteknn_clf.score(X_train, y_train))')
print('Accuracy on test set: ({}).format(pasteknn_clf.score(X_test, y_test))')

Accuracy on training set: 0.945
Accuracy on test set: 0.938
```

Cross validation is done on the model with $cv=10$

```
In [62]: scores = cross_val_score(pasteknn_clf, X_train, y_train, cv = 10, scoring = 'accuracy' )
print('Cross-validation scores: ({}).format(scores)')

print('Average cross-validation score: ({}).format(scores.mean())')

Cross-validation scores: [0.91752577 0.95876289 0.87628866 0.89690722 0.97938144 0.91666667 0.91666667 0.88541667 0.95833333 0.9375 ]
Average cross-validation score: 0.92
```

The gridsearch is used to find best parameters for $max_samples$, $max_features$ and $n_{estimators}$ for the decision tree

```
In [63]: param_grid = {'max_samples': [0.5, 0.6],
                    'max_features': [0.4, 0.5],
                    'n_estimators': [100, 500]}

grid_search_pasteknnclf = GridSearchCV(BaggingClassifier(DecisionTreeClassifier(max_depth= 2), random_state=0, bootstrap=False, n_jobs=-1), param_grid, cv=5, return_train_score=True)
grid_search_pasteknnclf.fit(X_train, y_train)
```

```
print('Best parameters: ({}).format(grid_search_pasteknnclf.best_params_)')
print('Best cross-validation score: ({}).format(grid_search_pasteknnclf.best_score_)')

Best parameters: {'max_features': 0.4, 'max_samples': 0.5, 'n_estimators': 500}
Best cross-validation score: 0.94
```

Confusion matrix is printed with the model outcome

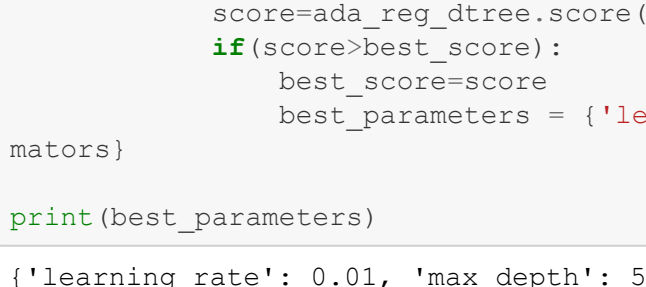
```
In [64]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_pastecf = pasteknn_clf.predict(X_test)
print(metrics.accuracy_score(y_test, pred_pastecf))

confusion = confusion_matrix(y_test, pred_pastecf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
print(classification_report(y_test, pred_pastecf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.91	0.92	0.92	91
accuracy			0.94	242
macro avg	0.93	0.93	0.93	242
weighted avg	0.94	0.94	0.94	242



The recall and F1 score of the model is printed

```
In [65]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_pastecf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9470198675496688
F1Score : 0.9501661129568105
```

```
In [66]: Classification_Scores.update({'Pasting-Decision Tree': [metrics.accuracy_score(y_test, pred_pastecf),
                                                                recall[0], fscore[0]]})
```

Pasting - KNN

KNN classifier with Pasting is modelled with bootstrap=false

```
In [67]: from sklearn.ensemble import BaggingClassifier
pasteknn_clf = BaggingClassifier(KNeighborsClassifier(n_neighbors=4), bootstrap=False, n_jobs=-1, random_state=0)
pasteknn_clf.fit(X_train, y_train)
print('Accuracy on training set: ({}).format(pasteknn_clf.score(X_train, y_train))')
print('Accuracy on test set: ({}).format(pasteknn_clf.score(X_test, y_test))')

Accuracy on training set: 0.945
Accuracy on test set: 0.938
```

Cross validation score is printed for the model

```
In [68]: scores = cross_val_score(pasteknn_clf, X_train, y_train, cv = 5, scoring = 'accuracy' )
print('Cross-validation scores: ({}).format(scores)')

print('Average cross-validation score: ({}).format(scores.mean())')

Cross-validation scores: [0.93264249 0.89119171 0.94818653 0.919171 0.94818653]
Average cross-validation score: 0.93
```

Gridsearch is used to find the best parameters for $max_samples$, $max_features$ and $n_{estimators}$ are found

```
In [69]: param_grid = {'max_samples': [0.5, 0.6],
                    'max_features': [0.4, 0.5],
                    'n_estimators': [100, 500]}

grid_search_pasteknnclf = GridSearchCV(BaggingClassifier(KNeighborsClassifier(), random_state=0, bootstrap=False, n_jobs=-1), param_grid, cv=5, return_train_score=True)
grid_search_pasteknnclf.fit(X_train, y_train)
```

```
print('Best parameters: ({}).format(grid_search_pasteknnclf.best_params_)')
print('Best cross-validation score: ({}).format(grid_search_pasteknnclf.best_score_)')

Best parameters: {'max_features': 0.4, 'max_samples': 0.5, 'n_estimators': 500}
Best cross-validation score: 0.94
```

Confusion matrix is printed for the model

```
In [70]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_pasteknn_clf = pasteknn_clf.predict(X_test)
print(metrics.accuracy_score(y_test, pred_pasteknn_clf))

confusion = confusion_matrix(y_test, pred_pasteknn_clf)
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax = plt.subplot()
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')
print(classification_report(y_test, pred_pasteknn_clf))
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	151
1	0.91	0.92	0.92	91
accuracy			0.94	242
macro avg	0.93	0.93	0.93	242
weighted avg	0.94	0.94	0.94	242

Recall and F1 score is calculated

```
In [76]: from sklearn.metrics import precision_recall_fscore_support as score
precision, recall, fscore, support = score(y_test, pred_pasteknn_clf)

print('Recall : ({}).format(recall[0])')
print('F1Score : ({}).format(fscore[0])')

Recall : 0.9470198675496688
F1Score : 0.9501661129568105
```



```
In [77]: Classification_Scores.update({'Adaboost=Decision Tree':metrics.accuracy_score(y_test,pred_adabo,recall[0],f_score[0])})
```

Adaboost with SVC Poly

Kernel=polynomial, n_estimators=200 and learning_rate=0.1 are used for SVC poly model with adaboost

```
In [78]: from sklearn.svm import SVC
svc_clf=SVC(C=1,gamma=0.1, kernel='poly',probability=True)
ada_clf_svc = AdaBoostClassifier(svc_clf, n_estimators=200, learning_rate=0.1, random_state=0)
ada_clf_svc.fit(X_train, y_train)

y_pred=ada_clf_svc.predict(X_test)
print('Train score: {:.4f} %'.format(ada_clf_svc.score(X_train, y_train)*100))
print('Test score: {:.4f} %'.format(ada_clf_svc.score(X_test, y_test)*100))

Train score: 93.8860 %
Test score: 90.9091 %

Confusion matrix is printed for the model
```

```
In [79]: pred_adabo_svc = ada_clf_svc.predict(X_test)
print(metrics.accuracy_score(y_test,pred_adabo_svc))

confusion = confusion_matrix(y_test, pred_adabo_svc)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

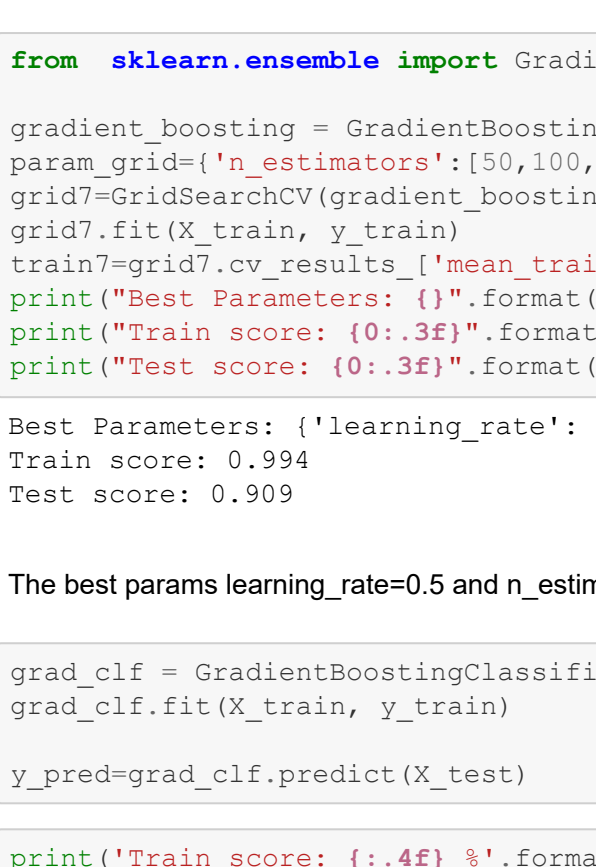
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,pred_adabo_svc))

0.9090909090909091
precision    recall    f1-score   support

0         0.91         0.95         0.93         151
1         0.91         0.85         0.87          91

accuracy          0.91          0.90          0.90         242
macro avg         0.91          0.91          0.91         242
weighted avg      0.91          0.91          0.91         242
```



```
In [80]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,pred_adabo_svc)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9470198675496688
F1Score     : 0.9285714285714285
```

```
In [81]: Classification_Scores.update({'Adaboost=SVC Poly Classifier':metrics.accuracy_score(y_test,pred_adabo_svc),recall[0],f_score[0]})
```

Gradient Boosting

Gradient boosting is a technique which minimizes the overall prediction error. The key idea is to set the target outcomes for the next best model in order to minimize the error. Gradient boosting model is primarily used with decision tree.

The key parameters used are n_estimators and learning_rate

```
In [82]: from sklearn.ensemble import GradientBoostingClassifier

gradient_boosting = GradientBoostingClassifier(max_depth=2, random_state=0)
param_grid={'n_estimators':[50,100,150], 'learning_rate':[0.5,1]}
grid=GridSearchCV(gradient_boosting,param_grid,cv=5,return_train_score=True)
grid7.fit(X_train, y_train)
train7=grid7.cv_results_['mean_train_score']
print("Best Parameters: {}".format(grid7.best_params_))
print("Train score: (0:.3f)".format(train7.mean()))
print("Test score: (0:.3f)".format(grid7.score(X_test, y_test)))

Best Parameters: {'learning_rate': 0.5, 'n_estimators': 50}
Train score: 0.994
Test score: 0.909
```

The best params learning_rate=0.5 and n_estimators=20 are used with max_depth=2 for decision tree

```
In [83]: grad_clf = GradientBoostingClassifier(max_depth=2, n_estimators=50, learning_rate=0.5, random_state=0)
grad_clf.fit(X_train, y_train)

y_pred=grad_clf.predict(X_test)
```

```
In [84]: print('Train score: {:.4f} %'.format(grad_clf.score(X_train, y_train)*100))
print('Test score: {:.4f} %'.format(grad_clf.score(X_test, y_test)*100))

Train score: 97.9275 %
Test score: 90.9091 %
```

Confusion Matrix is printed for the model

```
In [85]: grad = grad_clf.predict(X_test)
print(metrics.accuracy_score(y_test,grad))

confusion = confusion_matrix(y_test, grad)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

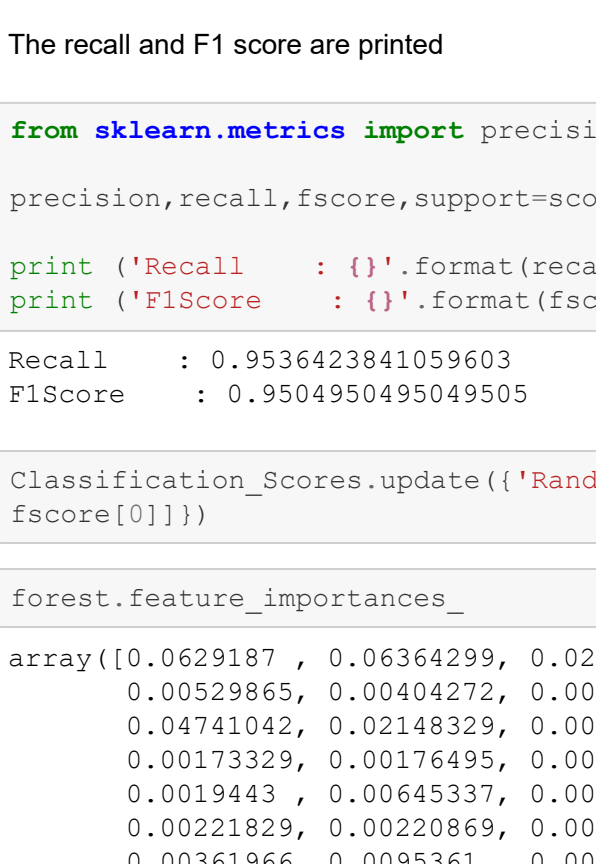
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,grad))

0.9090909090909091
precision    recall    f1-score   support

0         0.91         0.95         0.93         151
1         0.91         0.85         0.87          91

accuracy          0.91          0.90          0.91         242
macro avg         0.91          0.90          0.90         242
weighted avg      0.91          0.91          0.91         242
```



```
In [86]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,grad)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9470198675496688
F1Score     : 0.9285714285714285
```

```
In [87]: Classification_Scores.update({'Gradient boosting':metrics.accuracy_score(y_test,grad),recall[0],f_score[0]})
```

Random Forest

The random forest is a classifier consisting of many decisions trees

```
In [88]: from sklearn.ensemble import RandomForestClassifier

forest = RandomForestClassifier(random_state=0)
forest.fit(X_train, y_train)

print("Accuracy on training set: {:.3f}".format(forest.score(X_train, y_train)))
print("Accuracy on test set: {:.3f}".format(forest.score(X_test, y_test)))

Accuracy on training set: 0.997
Accuracy on test set: 0.938
```

Confusion Matrix is generated

```
In [89]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_forestclf = forest.predict(X_test)
print(metrics.accuracy_score(y_test,pred_forestclf))

confusion = confusion_matrix(y_test, pred_forestclf)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

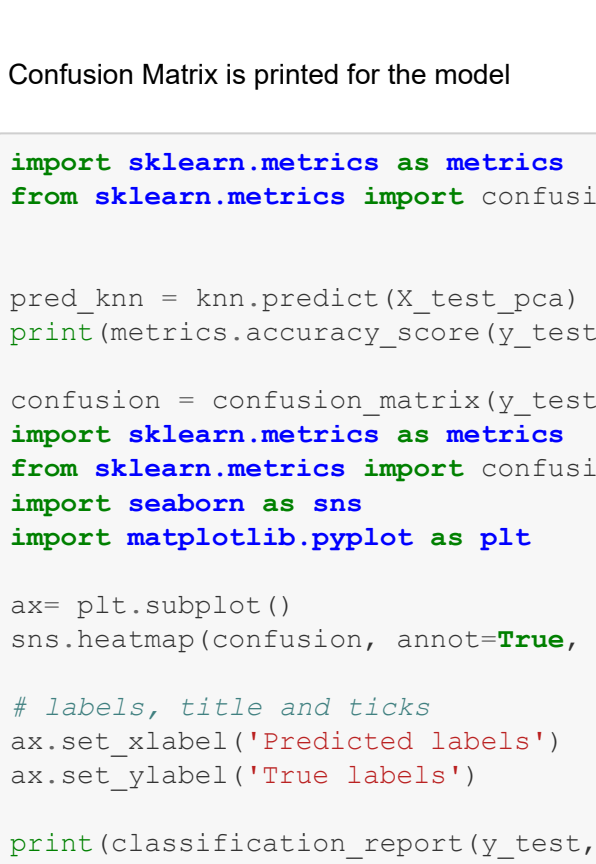
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,pred_forestclf))

0.9380165289256198
precision    recall    f1-score   support

0         0.95         0.95         0.95         151
1         0.92         0.91         0.92          91

accuracy          0.94          0.94          0.94         242
macro avg         0.93          0.93          0.93         242
weighted avg      0.94          0.94          0.94         242
```



The recall and F1 score are printed

```
In [90]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,pred_forestclf)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9536423841059603
F1Score     : 0.9504950495049505
```

```
In [91]: Classification_Scores.update({'Random Forest':metrics.accuracy_score(y_test,pred_forestclf),recall[0],f_score[0]})

In [92]: forest.feature_importances_
array([0.0629187 , 0.06364299, 0.02162273, 0.01250533, 0.02042624,
       0.00529865, 0.00404272, 0.00353806, 0.00871052, 0.00425258,
       0.04741042, 0.02148329, 0.00285393, 0.00686045, 0.0022707,
       0.00173329, 0.00176495, 0.00229307, 0.00082871, 0.00435428,
       0.0019443 , 0.00645337, 0.00466284, 0.00269751, 0.00331846,
       0.00221829, 0.00220869, 0.00113739, 0.00128925, 0.00165149,
       0.00361666, 0.0085361 , 0.00724826, 0.00876158, 0.00781325,
       0.01682447, 0.00158011, 0.00061867, 0.00340788, 0.00454487,
       0.28132018, 0.32447714, 0.00397691])
```

PCA

PCA is statistical technique using dimensionality reduction with sole basis that large number of features add to high dimensionality which essentially leads to overfitting.

PCA thus reduces the number of features taken into consideration.

```
In [93]: from sklearn.decomposition import PCA

# initialize pca and logistic regression model
pca = PCA(n_components=0.95)

In [94]: X_train_pca= pca.fit_transform(X_train)
X_test_pca = pca.transform(X_test)

print(X_train_pca.shape)
print(y_train.shape)
print(X_test_pca.shape)
print(y_test.shape)

(965, 24)
(965,)
(242, 24)
(242,)
```

PCA is run on all the models

KNN with PCA

KNN classifier is modeled with cv=5 and the cross val score is calculated

GridSearch is used to find the best neighbor that can be used

```
In [95]: k_range = list(range(1, 11))

param_grid = dict(n_neighbors=k_range)

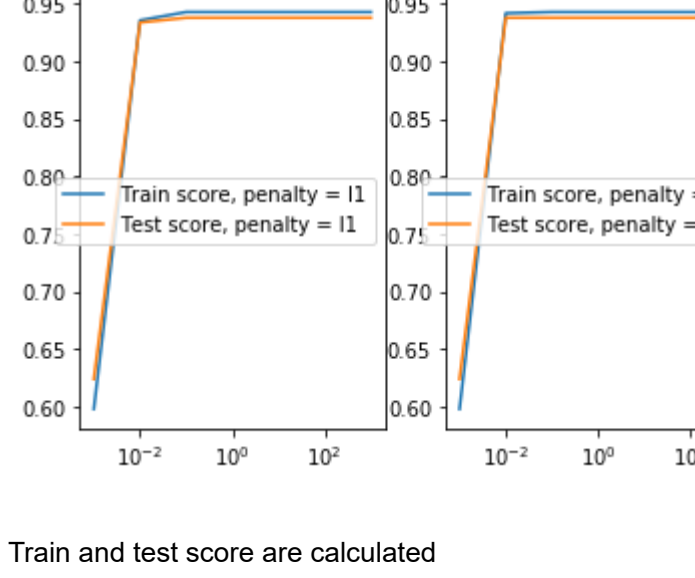
grid_search_knn = GridSearchCV(KNeighborsClassifier(), param_grid, cv=5, return_train_score=True)

grid_search_knn.fit(X_train_pca, y_train)

df = pd.DataFrame(grid_search_knn.cv_results_)
import matplotlib.pyplot as plt
ax=plt.subplot(1)
plt.plot(k_range, df.mean_train_score, c='g', label='Train Score')
plt.plot(k_range, df.mean_test_score, c='b', label='Validation Score')
plt.legend()
plt.xlabel('k')
plt.ylabel('CV Score')

print("Best parameters: {}".format(grid_search_knn.best_params_))
print("Best cross-validation score: {:.2f}".format(grid_search_knn.best_score_))

Best parameters: {'n_neighbors': 5}
Best cross-validation score: 0.93
```



Neighbour=5 is taken for the KNN classifier

```
In [96]: knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train_pca, y_train)
print("Train score on best parameters {:.4f}".format(knn.score(X_train_pca, y_train)))
print("Test score on best parameters {:.4f}".format(knn.score(X_test_pca, y_test)))

Train score on best parameters 0.9461
Test score on best parameters 0.9380
```

Cross validation score for the KNN is found

```
In [97]: from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
knn_grid = KNeighborsClassifier(5)

scores = cross_val_score(knn_grid, X_train_pca, y_train, cv=5, scoring='accuracy')
print("Cross-validation scores: {}".format(scores))

print("Average cross-validation score: {:.2f}".format(scores.mean()))

Cross-validation scores: [0.94818653 0.91709845 0.94300518 0.92227979 0.93782383]
Average cross-validation score: 0.93
```

Confusion Matrix is printed for the model

```
In [98]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report

pred_knn = knn.predict(X_test_pca)
print(metrics.accuracy_score(y_test,pred_knn))

confusion = confusion_matrix(y_test, pred_knn)
import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True, fmt='d', ax=ax); #annot=True to annotate cells

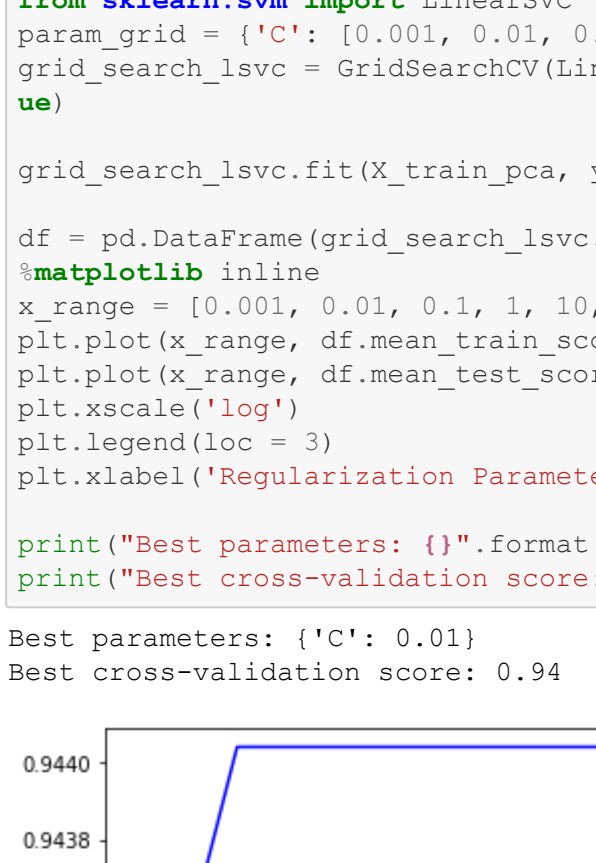
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,pred_knn))

0.9380165289256198
precision    recall    f1-score   support

0         0.95         0.95         0.95         151
1         0.91         0.92         0.92          91

accuracy          0.94          0.94          0.94         242
macro avg         0.93          0.94          0.93         242
weighted avg      0.94          0.94          0.94         242
```



Recall and F1 score is printed

```
In [99]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,pred_knn)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9470198675496688
F1Score     : 0.9501661129568105
```

```
In [100]: Classification_Scores.update({'KNN Classification with PCA':metrics.accuracy_score(y_test,pred_knn),recall[0],f_score[0]})
```

Logistic reg after PCA

Logistic regression is run after the PCA is done on the dataset

```
In [101]: # import warnings filter
from warnings import simplifyfilter
# ignore all future warnings
simplifyfilter(action='ignore', category=FutureWarning)

from sklearn.linear_model import LogisticRegression

c_range = [0.001, 0.01, 0.1, 1, 10, 100, 1000]
train_score_l1 = []
train_score_l2 = []
valid_score_l1 = []
valid_score_l2 = []

best_score = 0
l1 = ''
l2 = ''

for c in c_range:
    log_l1 = LogisticRegression(penalty = 'l1', C = c,solver='liblinear')
    log_l2 = LogisticRegression(penalty = 'l2', C = c,solver='lbfgs')
    log_l1.fit(X_train_pca, y_train)
    log_l2.fit(X_train_pca, y_train)

    train_score_l1.append(log_l1.score(X_train_pca, y_train))
    train_score_l2.append(log_l2.score(X_train_pca, y_train))

    score = log_l1.score(X_test_pca, y_test)
    valid_score_l1.append(score)
    if score > best_score:
        best_score = score
        best_parameters = {'C': c, 'penalty': 'l1'}
        best_C = c
        best_Penalty = 'l1'

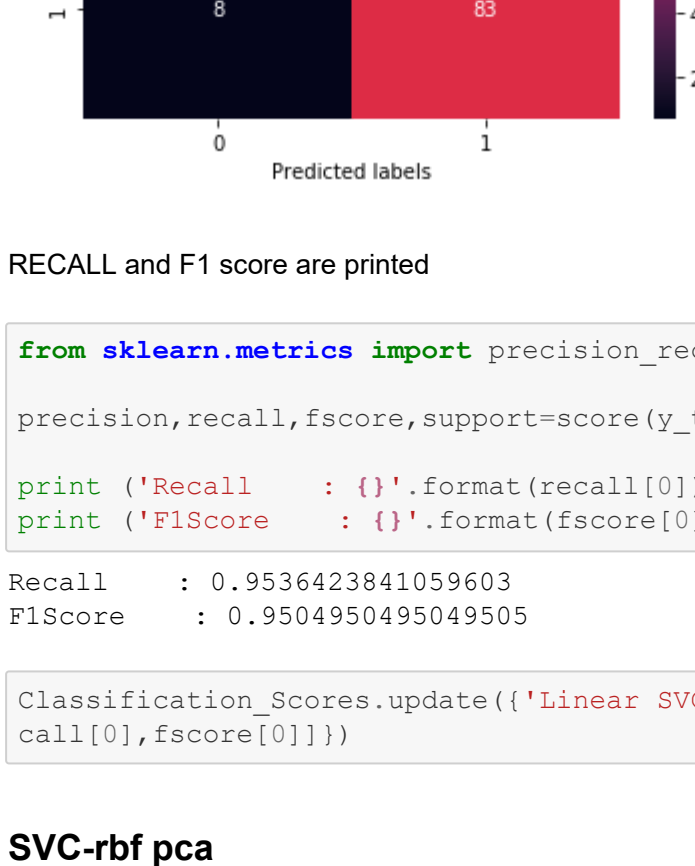
    score = log_l2.score(X_test_pca, y_test)
    valid_score_l2.append(score)
    if score > best_score:
        best_score = score
        best_parameters = {'C': c, 'penalty': 'l2'}
        best_C = c
        best_Penalty = 'l2'

plt.subplot(1,2,1)
plt.plot(c_range, train_score_l1, label = 'Train score, penalty = l1')
plt.plot(c_range, valid_score_l1, label = 'Test score, penalty = l1')
plt.xscale('log')
plt.legend()

plt.subplot(1,2,2)
plt.plot(c_range, train_score_l2, label = 'Train score, penalty = l2')
plt.plot(c_range, valid_score_l2, label = 'Test score, penalty = l2')
plt.xscale('log')
plt.legend()

print("Best score: {:.2f}".format(best_score))
print("Best parameters: {}".format(best_parameters))

Best score: 0.94
Best parameters: {'C': 0.01, 'penalty': 'l2'}
```



Train and test score are calculated

```
In [102]: lg = LogisticRegression(C=0.01,penalty='l2').fit(X_train_pca, y_train)
print("Train score on best parameters for Logistic Regression model {:.4f}".format(lg.score(X_train_pca, y_train)))
print("Test score on best parameters for Logistic Regression model {:.4f}".format(lg.score(X_test_pca, y_test)))

Train score on best parameters for Logistic Regression model 0.9420
Test score on best parameters for Logistic Regression model 0.9380
```

Logistic regression is modelled with penalty L2 and C=0.01

```
In [103]: log_grid = LogisticRegression(penalty = 'l2', C = 0.01)

scores = cross_val_score(log_grid, X_train_pca, y_train, cv=5, scoring='accuracy')
print("Cross-validation scores: {}".format(scores))

print("Average cross-validation score: {:.2f}".format(scores.mean()))

Cross-validation scores: [0.96373057 0.92227979 0.93782383 0.9119171 0.97409326]
Average cross-validation score: 0.94
```

```
In [104]: pred_log = lg.predict(X_test_pca)
print(metrics.accuracy_score(y_test,pred_log))

confusion = confusion_matrix(y_test, pred_log)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True,fmt='d', ax=ax); #annot=True to annotate cells

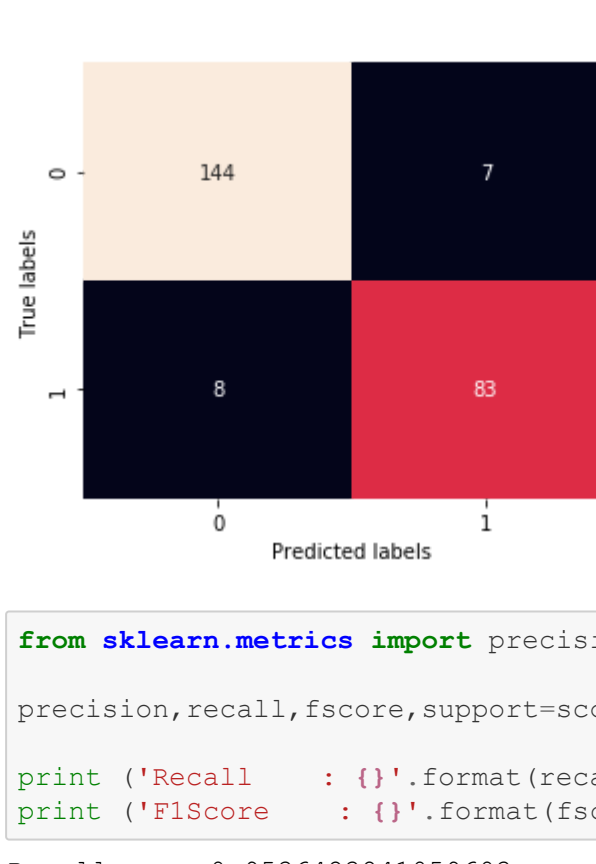
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,pred_log))

0.9380165289256198
precision    recall    f1-score   support

0         0.95         0.95         0.95         151
1         0.92         0.91         0.92          91

accuracy          0.93          0.93          0.94         242
macro avg         0.93          0.93          0.93         242
weighted avg      0.94          0.94          0.94         242
```



RECALL and F1 score are printed

```
In [111]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,pred_linear_svc)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9536423841059603
F1Score     : 0.9504950495049505
```

```
In [112]: Classification_Scores.update({'Linear SVC with PCA':metrics.accuracy_score(y_test,pred_linear_svc),recall[0],f_score[0]})
```

SVC-rbf pca

Gridsearch is used to find the best param of gamma and C

```
In [113]: param_grid = {'gamma': [0.001, 0.01, 0.1, 0.1, 1, 10, 100],
                  'C': [0.001, 0.01, 0.1, 1, 10, 100]}

grid_search_svc = GridSearchCV(LinearSVC(max_iter=100000000), param_grid, cv=5, return_train_score=True)

grid_search_svc.fit(X_train_pca, y_train)

df = pd.DataFrame(grid_search_svc.cv_results_)
print("Best parameters: {}".format(grid_search_svc.best_params_))
print("Best cross-validation score: {:.2f}".format(grid_search_svc.best_score_))

Best parameters: {'C': 0.1, 'gamma': 0.1}
Best cross-validation score: 0.93
```

The best param of C=0.1 and gamma=0.1 are used

```
In [114]: clf2 = SVC(kernel='rbf', C=0.1,gamma=0.1).fit(X_train_pca, y_train)
print("Train score on best parameters for LinearSVC {:.4f}".format(clf1.score(X_train_pca, y_train)))
print("Test score on best parameters ffor LinearSVC {:.4f}".format(clf1.score(X_test_pca, y_test)))

Train score on best parameters for LinearSVC = 0.9440
Test score on best parameters ffor LinearSVC = 0.9380
```

```
In [115]: from sklearn.model_selection import cross_val_score
svc_rbf_grid = SVC(kernel='rbf', gamma = 0.1, C = 0.1)

scores = cross_val_score(svc_rbf_grid, X_train_pca, y_train, cv=5, scoring='accuracy')
print("Cross-validation scores: {}".format(scores))

print("Average cross-validation score: {:.2f}".format(scores.mean()))

Cross-validation scores: [0.94845361 0.96807216 0.89690722 0.95876289 0.96307216 0.91666667 0.95833333 0.875 0.98598333 0.95833333]
Average cross-validation score: 0.94
```

Confusion Matrix is printed

```
In [116]: pred_rbf = clf2.predict(X_test_pca)
print(metrics.accuracy_score(y_test,pred_rbf))

confusion = confusion_matrix(y_test, pred_rbf)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot(1)
sns.heatmap(confusion, annot=True,fmt='d', ax=ax); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test,pred_rbf))

0.9380165289256198
precision    recall    f1-score   support

0         0.95         0.95         0.95         151
1         0.92         0.91         0.92          91

accuracy          0.93          0.93          0.94         242
macro avg         0.93          0.93          0.93         242
weighted avg      0.94          0.94          0.94         242
```



```
In [117]: from sklearn.metrics import precision_recall_fscore_support as score
precision,recall,f_score,support=score(y_test,pred_rbf)

print ('Recall      : {}'.format(recall[0]))
print ('F1Score     : {}'.format(fscore[0]))

Recall      : 0.9536423841059603
F1Score     : 0.9504950495049505
```

```
In [118]: Classification_Scores.update({'SVC-rbf with PCA':metrics.accuracy_score(y_test,pred_rbf),recall[0],f_score[0]})
```


SVC Poly PCA

SVC Poly is modeled with grid dearch for best param - gamma, C, degree

```
In [119]: param_grid = {'gamma': [0.01,0.1],
                  'C': [0.01, 0.1, 1,10],
                  'degree': [1,2,3,4,5]}

grid_search_svc = GridSearchCV(SVC(kernel='poly'), param_grid, cv=5, return_train_score=True)
grid_search_svc.fit(X_train_pca, y_train)

print("Best parameters: {}".format(grid_search_svc.best_params_))
print("Best cross-validation score: {:.2f}".format(grid_search_svc.best_score_))

Best parameters: {'C': 0.1, 'degree': 1, 'gamma': 0.1}
Best cross-validation score: 0.94

Train and Test score of the linear SVC

In [120]: clf = SVC(kernel='rbf', C=0.1, gamma=0.1, degree=1).fit(X_train_pca, y_train)
print("Train score on best parameters for LinearSVC - {:.4f}".format(clf.score(X_train_pca, y_train)))
print("Test score on best parameters ffor LinearSVC - {:.4f}".format(clf.score(X_test_pca, y_test)))

Train score on best parameters for LinearSVC - 0.9440
Test score on best parameters ffor LinearSVC - 0.9380

In [121]: svc_poly_grid = SVC(kernel='poly', degree = 1, C=0.01, gamma=1)

scores = cross_val_score(svc_poly_grid, X_train_pca, y_train, cv=5, scoring = 'accuracy')
print("Cross-validation scores: {}".format(scores))
print("Average cross-validation score: {:.2f}".format(scores.mean()))

Cross-validation scores: [0.96373057 0.92227979 0.94300518 0.91709845 0.97409326]
Average cross-validation score: 0.94

Confusion Matrix

In [122]: # import warnings filter
from warnings import simplefilter
# ignore all future warnings
simplefilter(action='ignore', category=FutureWarning)

pred_poly = clf3.predict(X_test_pca)
print(metrics.accuracy_score(y_test, pred_poly))

confusion = confusion_matrix(y_test, pred_poly)

import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot()
sns.heatmap(confusion, annot=True,fmt='d', ax = ax); #annot=True to annotate cells

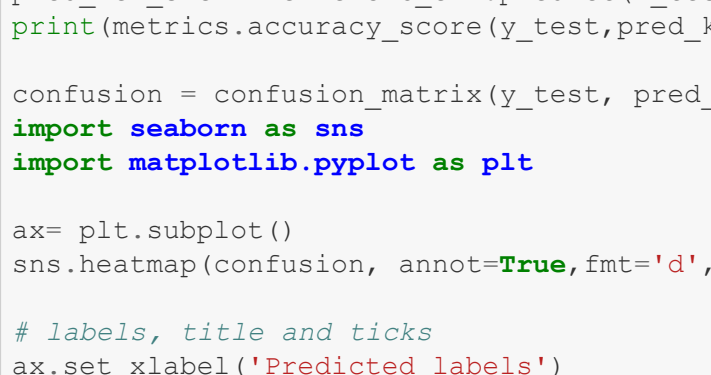
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test, pred_poly))

0.9380165289256198
              precision    recall  f1-score   support

             0         0.95         0.95         0.95         151
             1         0.92         0.91         0.92           91

 accuracy         0.93         0.93         0.94         242
 macro avg         0.93         0.93         0.93         242
weighted avg         0.94         0.94         0.94         242
```



Recall and F1 score are printed

```
In [123]: from sklearn.metrics import precision_recall_fscore_support as score

precision, recall, fscore, support=score(y_test, pred_poly)

print ('Recall      : {}'.format(recall[0]))
print ('f1score     : {}'.format(fscore[0]))

Recall      : 0.936423841059603
f1score      : 0.9304950495049505

In [124]: Classification_Scores.update({'SVC-poly with PCA':(metrics.accuracy_score(y_test, pred_poly), recall[0],
fscore[0])})
```

SVC rbf and linear with PCA

SVC RBF is modelled with gridsearch for selecting the best params for C, gamma

```
In [125]: #using Kernel SVM
from sklearn.svm import SVC
KernelSVC = SVC(max_iter=100000000)
kernelSVC_params = {'C':[0.001, 0.01, 0.1, 1], 'gamma':[1,0.1,0.001], 'kernel':['rbf', 'linear']}

In [126]: # Using Grid search to find the best parameters and fitting the model
KernelSVC_clf = GridSearchCV(KernelSVC, kernelSVC_params, cv=5)
KernelSVC_clf.fit(X_train_pca, y_train)
print("Train score on best parameters for LinearSVC - {:.4f}".format(clf1.score(X_train_pca, y_train)))
print("Test score on best parameters ffor LinearSVC - {:.4f}".format(clf1.score(X_test_pca, y_test)))

Train score on best parameters for LinearSVC - 0.9440
Test score on best parameters ffor LinearSVC - 0.9380

In [127]: #Finding the best paramter
KernelSVC_clf.best_params_

Out[127]: {'C': 0.1, 'gamma': 0.1, 'kernel': 'rbf'}
```

CONFUSION MATRIX

```
In [128]: pred_ker_svc = KernelSVC_clf.predict(X_test_pca)
print(metrics.accuracy_score(y_test, pred_ker_svc))

confusion = confusion_matrix(y_test, pred_ker_svc)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot()
sns.heatmap(confusion, annot=True,fmt='d', ax = ax); #annot=True to annotate cells

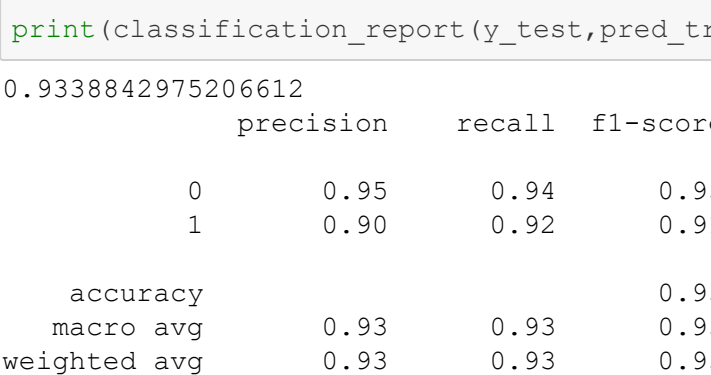
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test, pred_ker_svc))

0.9380165289256198
              precision    recall  f1-score   support

             0         0.95         0.95         0.95         151
             1         0.91         0.91         0.92           91

 accuracy         0.93         0.93         0.94         242
 macro avg         0.92         0.93         0.93         242
weighted avg         0.94         0.94         0.94         242
```



```
In [129]: from sklearn.metrics import precision_recall_fscore_support as score

precision, recall, fscore, support=score(y_test, pred_ker_svc)

print ('Recall      : {}'.format(recall[0]))
print ('f1score     : {}'.format(fscore[0]))

Recall      : 0.936423841059603
f1score      : 0.930495049505

In [130]: Classification_Scores.update({'SVC-linear with PCA':(metrics.accuracy_score(y_test, pred_ker_svc), recall
1[0], fscore[0])})
```

Decision Tree PCA

Decision tree is modeled after PCA on the dataset with grid search for best param of max_depth and min_sample_leaf

```
In [131]: from sklearn.tree import DecisionTreeClassifier

param_grid = {'criterion': ['gini', 'entropy'],
              'max_depth': [1,2,3,4,5,6,7,8,9,10],
              'min_samples_leaf': [1,2,3,4,5,6,7,8,9,10]}

grid_search_dtree = GridSearchCV(DecisionTreeClassifier(random_state=0), param_grid, cv=5, return_train_score=True)
grid_search_dtree.fit(X_train_pca, y_train)

print("Best parameters: {}".format(grid_search_dtree.best_params_))
print("Best cross-validation score: {:.2f}".format(grid_search_dtree.best_score_))

Best parameters: {'criterion': 'gini', 'max_depth': 1, 'min_samples_leaf': 1}
Best cross-validation score: 0.93

In [132]: dtree = DecisionTreeClassifier(max_depth=1, min_samples_leaf=1, criterion='gini')

dtree.fit(X_train_pca, y_train)
print("Accuracy on training set: {:.3f}".format(dtree.score(X_train_pca, y_train)))
print("Accuracy on test set: {:.3f}".format(dtree.score(X_test_pca, y_test)))

Accuracy on training set: 0.940
Accuracy on test set: 0.934

Cross validations score is obtained

In [133]: dtree_cv = DecisionTreeClassifier()
scores = cross_val_score(dtree_cv, X_train_pca, y_train, cv = 5, scoring = 'accuracy')
print("Cross-validation scores: {}".format(scores))

print("Average cross-validation score: {:.2f}".format(scores.mean()))

Cross-validation scores: [0.88601036 0.86528497 0.90673575 0.87564767 0.92227979]
Average cross-validation score: 0.89
```

CONFUSION MATRIX

```
In [134]: import sklearn.metrics as metrics
from sklearn.metrics import confusion_matrix, classification_report
pred_tree = grid_search_dtree.predict(X_test_pca)
print(metrics.accuracy_score(y_test, pred_tree))

confusion = confusion_matrix(y_test, pred_tree)
import seaborn as sns
import matplotlib.pyplot as plt

ax=plt.subplot()
sns.heatmap(confusion, annot=True,fmt='d', ax = ax); #annot=True to annotate cells

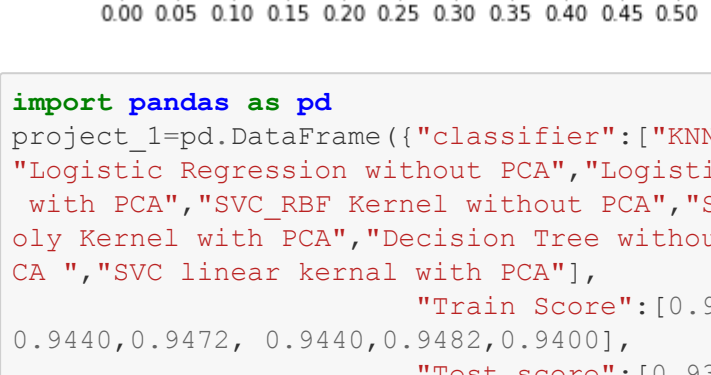
# labels, title and ticks
ax.set_xlabel('Predicted labels')
ax.set_ylabel('True labels')

print(classification_report(y_test, pred_tree))

0.9338842975206612
              precision    recall  f1-score   support

             0         0.95         0.94         0.95         151
             1         0.90         0.92         0.91           91

 accuracy         0.93         0.93         0.93         242
 macro avg         0.93         0.93         0.93         242
weighted avg         0.93         0.93         0.93         242
```



```
In [135]: from sklearn.metrics import precision_recall_fscore_support as score

precision, recall, fscore, support=score(y_test, pred_tree)

print ('Recall      : {}'.format(recall[0]))
print ('f1score     : {}'.format(fscore[0]))

Recall      : 0.940397509933775
f1score      : 0.9466666666666667

In [136]: Classification_Scores.update({'Decision Tree with PCA':(metrics.accuracy_score(y_test, pred_tree), recall
1[0], fscore[0])})
```

COMPARISON OF ALL MODELS

```
In [137]: xl=Classification_Scores.keys()
df1=pd.DataFrame(Classification_Scores, index= ['Accuracy', 'Recall', 'f1Score'])
df1

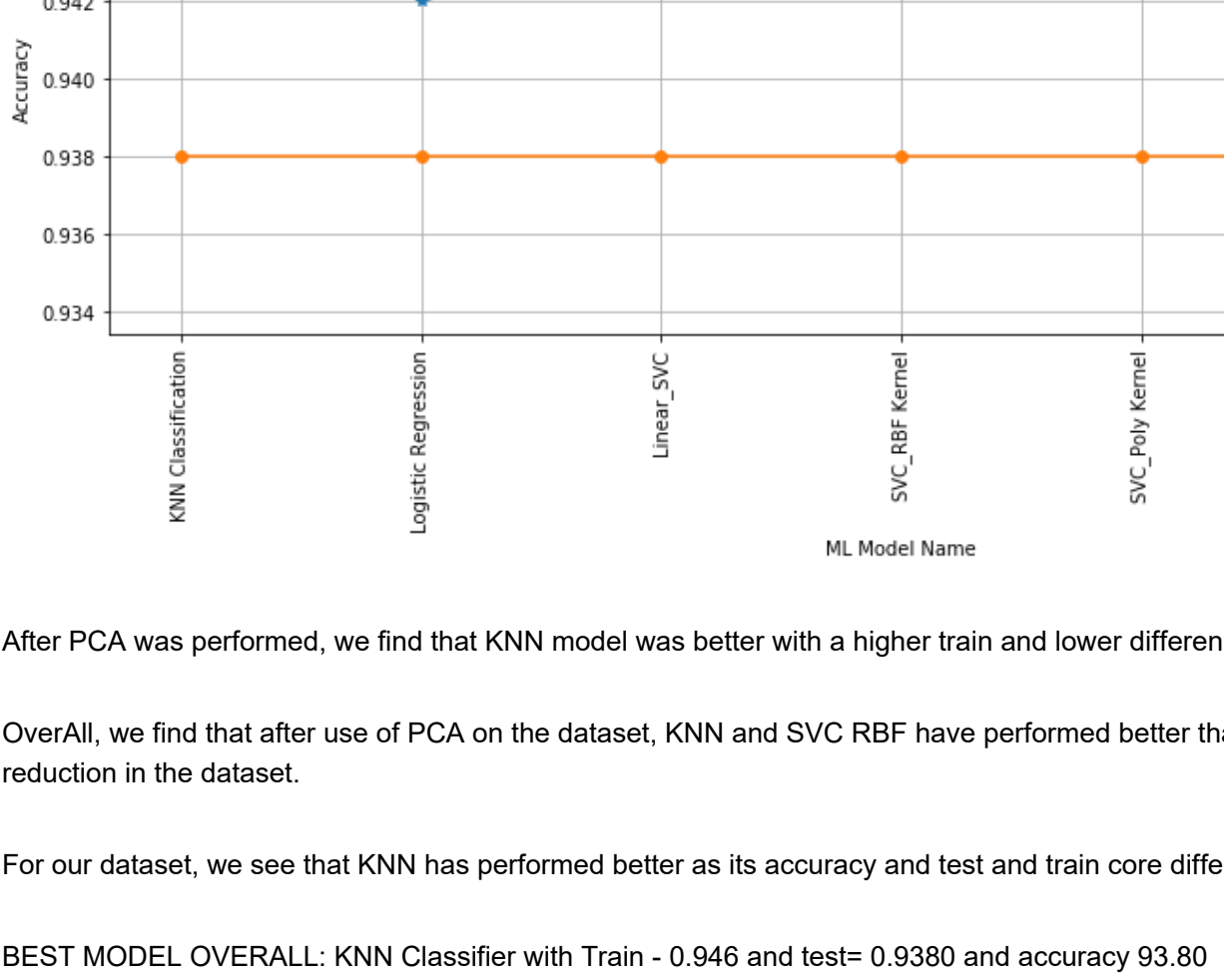
Out[137]:
```

	Hard voting Classifier	Soft voting Classifier	Bagging-Decision Tree Classifier	Bagging-KNN Classifier	Pasting-Decision Tree	Pasting-KNN Classifier	AdaBoosting-Decision Tree	AdaBoosting-SVC Poly Classifier	Gradient boosting	Random Forest	KNN Classification with PCA
Accuracy	0.938017	0.938017	0.942149	0.938017	0.938017	0.938017	0.946281	0.909091	0.909091	0.938017	0.938017
Recall	0.953642	0.953642	0.953642	0.953642	0.953642	0.947020	0.947020	0.947020	0.947020	0.953642	0.947020
F1Score	0.950495	0.950495	0.953642	0.950495	0.950495	0.950166	0.956522	0.928571	0.928571	0.950495	0.950166

The accuracy, recall and F1 score for the models are printed before and after PCA

```
In [138]: pd.DataFrame(df1).plot(kind='barh', figsize=(10,10))
plt.grid(which='both')
plt.xticks(np.arange(0,1,step=0.05))
plt.legend(loc='center', prop={'size':12})

Out[138]: <matplotlib.legend.Legend at 0x1d382ca88>
```



```
In [139]: import pandas as pd
project_1=pd.DataFrame({"Classifier":["KNN Classification without PCA", "KNN Classification with PCA",
"Logistic Regression without PCA", "Logistic Regression with PCA", "Linear SVC without PCA", "Linear SVC
with PCA", "SVC_RBF Kernel without PCA", "SVC_RBF Kernel with PCA", "SVC_Poly Kernel without PCA", "SVC_P
oly Kernel with PCA", "Decision Tree without PCA", "Decision Tree with PCA", "SVC linear kernal without P
CA ", "SVC linear kernal with PCA"]})

"Train score": [0.9389, 0.9461, 0.9461, 0.9420, 0.9461, 0.9440, 0.6070, 0.9440, 0.9472,
0.9440, 0.9472, 0.9440, 0.9469, 0.9400, 0.9400],
"Test score": [0.9380, 0.9380, 0.9380, 0.9339, 0.9380, 0.9380, 0.9380, 0.9380, 0.6240, 0.9380, 0.9380,
0.9380, 0.9380, 0.9380, 0.9380, 0.9340],
"Accuracy": [93.80, 93.80, 93.38, 93.80, 93.80, 93.80, 62.39, 93.80, 93.80, 93.80, 62.39,
93.80, 93.80, 93.80, 93.38]]

project_1

Out[139]:
```

	Classifier	Train Score	Test score	Accuracy
0	KNN Classification without PCA	0.9389	0.9380	93.80
1	KNN Classification with PCA	0.9461	0.9380	93.80
2	Logistic Regression without PCA	0.9461	0.9339	93.38
3	Logistic Regression with PCA	0.9420	0.9380	93.80
4	Linear_SVC without PCA	0.9461	0.9380	93.80
5	Linear_SVC with PCA	0.9440	0.9380	93.80
6	SVC_RBF Kernel without PCA	0.6070	0.6240	62.39
7	SVC_RBF Kernel with PCA	0.9440	0.9380	93.80
8	SVC_Poly Kernel without PCA	0.9472	0.9380	93.80
9	SVC_Poly Kernel with PCA	0.9440	0.9380	93.80
10	Decision Tree without PCA	0.9472	0.9380	62.39
11	Decision Tree with PCA	0.9440	0.9380	93.80
12	SVC linear kernal without PCA	0.9482	0.9380	93.80
13	SVC linear kernal with PCA	0.9400	0.9340	93.38

BEFORE PCA

```
In [140]: import matplotlib.pyplot as plt
btrain_score=[0.9389, 0.9461, 0.9461, 0.9420, 0.9461, 0.9440, 0.6070, 0.9440, 0.9472, 0.9440, 0.9472,
0.9440, 0.9472, 0.9440, 0.9469, 0.9400]
btest_score=[0.9380, 0.9380, 0.9380, 0.9339, 0.9380, 0.9380, 0.6240, 0.9380, 0.9380, 0.9380, 0.9380,
0.9380, 0.9380, 0.9380, 0.9380, 0.9340]
Models=["KNN Classification", "Logistic Regression", "Linear_SVC", "SVC_RBF Kernel", "SVC_Poly Kernel", "De
cision Tree", "SVC linear kernal"]

plt.figure(figsize=(15,5))
plt.plot(Models, btrain_score, marker='*', label='Train Score')
plt.plot(Models, btest_score, marker='o', label='Test Score')
plt.legend(prop={'size': 12})
plt.grid()
plt.xticks(rotation='vertical')
plt.title('Before PCA, ML Models Train Vs Test score comparison')
plt.xlabel('ML Model Name')
plt.ylabel('Accuracy')
plt.show()
```



Before PCA, we find that SVM Kernel poly was the better model with low test and train score difference.

AFTER PCA

```
In [141]: import matplotlib.pyplot as plt
btrain_score=[0.9461, 0.9420, 0.9440, 0.9440, 0.9440, 0.9440, 0.9440, 0.9400]
btest_score=[0.9380, 0.9380, 0.9380, 0.9380, 0.9380, 0.9380, 0.9380, 0.9340]
Models=["KNN Classification", "Logistic Regression", "Linear_SVC", "SVC_RBF Kernel", "SVC_Poly Kernel", "De
cision Tree", "SVC linear kernal"]

plt.figure(figsize=(15,5))
plt.plot(Models, btrain_score, marker='*', label='Train Score')
plt.plot(Models, btest_score, marker='o', label='Test Score')
plt.legend(prop={'size': 12})
plt.grid()
plt.xticks(rotation='vertical')
plt.title('Before PCA, ML Models Train Vs Test score comparison')
plt.xlabel('ML Model Name')
plt.ylabel('Accuracy')
plt.show()
```



After PCA was performed, we find that KNN model was better with a higher train and lower difference with test and train score

Overall, we find that after use of PCA on the dataset, KNN and SVC RBF have performed better than before. This is owed to the dimension reduction in the dataset.

For our dataset, we see that KNN has performed better as its accuracy and test and train core difference is less after PCA.

BEST MODEL OVERALL: KNN Classifier with Train - 0.946 and test= 0.9380 and accuracy 93.80

Deep Learning Tasks

Neural Networks

```
In [142]: import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
import numpy

# fix random seed for reproducibility
numpy.random.seed(10)

In [143]: import tensorflow as tf

In [144]: model = Sequential()
model.add(Dense(12, input_dim=43, activation='relu'))
model.add(Dense(8, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

In [145]: model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

In [146]: X_train=np.asarray(X_train)
y_train = np.asarray(y_train)
X_test=np.asarray(X_test)
y_test = np.asarray(y_test)
```



```
[165]: # Fit the model
model1.fit(X_train, y_train, epochs=100, batch_size=30)

Epoch 1/100
33/33 [=====] - 0s 1ms/step - loss: 0.6818 - accuracy: 0.5399
Epoch 2/100
33/33 [=====] - 0s 1ms/step - loss: 0.5995 - accuracy: 0.6508
Epoch 3/100
33/33 [=====] - 0s 1ms/step - loss: 0.4570 - accuracy: 0.8902
Epoch 4/100
33/33 [=====] - 0s 1ms/step - loss: 0.2822 - accuracy: 0.9347
Epoch 5/100
33/33 [=====] - 0s 1ms/step - loss: 0.2148 - accuracy: 0.9430
Epoch 6/100
33/33 [=====] - 0s 1ms/step - loss: 0.1988 - accuracy: 0.9430
Epoch 7/100
33/33 [=====] - 0s 1ms/step - loss: 0.1877 - accuracy: 0.9451
Epoch 8/100
33/33 [=====] - 0s 1ms/step - loss: 0.1799 - accuracy: 0.9461
Epoch 9/100
33/33 [=====] - 0s 1ms/step - loss: 0.1761 - accuracy: 0.9451
Epoch 10/100
33/33 [=====] - 0s 1ms/step - loss: 0.1717 - accuracy: 0.9451
Epoch 11/100
33/33 [=====] - 0s 1ms/step - loss: 0.1709 - accuracy: 0.9472
Epoch 12/100
33/33 [=====] - 0s 1ms/step - loss: 0.1644 - accuracy: 0.9513
Epoch 13/100
33/33 [=====] - 0s 1ms/step - loss: 0.1617 - accuracy: 0.9482
Epoch 14/100
33/33 [=====] - 0s 1ms/step - loss: 0.1589 - accuracy: 0.9513
Epoch 15/100
33/33 [=====] - 0s 1ms/step - loss: 0.1569 - accuracy: 0.9554
Epoch 16/100
33/33 [=====] - 0s 1ms/step - loss: 0.1550 - accuracy: 0.9523
Epoch 17/100
33/33 [=====] - 0s 1ms/step - loss: 0.1528 - accuracy: 0.9534
Epoch 18/100
33/33 [=====] - 0s 1ms/step - loss: 0.1488 - accuracy: 0.9575
Epoch 19/100
33/33 [=====] - 0s 1ms/step - loss: 0.1475 - accuracy: 0.9534
Epoch 20/100
33/33 [=====] - 0s 1ms/step - loss: 0.1452 - accuracy: 0.9534
Epoch 21/100
33/33 [=====] - 0s 1ms/step - loss: 0.1442 - accuracy: 0.9565
Epoch 22/100
33/33 [=====] - 0s 1ms/step - loss: 0.1414 - accuracy: 0.9565
Epoch 23/100
33/33 [=====] - 0s 1ms/step - loss: 0.1398 - accuracy: 0.9575
Epoch 24/100
33/33 [=====] - 0s 1ms/step - loss: 0.1381 - accuracy: 0.9554
Epoch 25/100
33/33 [=====] - 0s 1ms/step - loss: 0.1352 - accuracy: 0.9575
Epoch 26/100
33/33 [=====] - 0s 1ms/step - loss: 0.1339 - accuracy: 0.9596
Epoch 27/100
33/33 [=====] - 0s 1ms/step - loss: 0.1323 - accuracy: 0.9596
Epoch 28/100
33/33 [=====] - 0s 1ms/step - loss: 0.1307 - accuracy: 0.9565
Epoch 29/100
33/33 [=====] - 0s 1ms/step - loss: 0.1294 - accuracy: 0.9596
Epoch 30/100
33/33 [=====] - 0s 1ms/step - loss: 0.1277 - accuracy: 0.9585
Epoch 31/100
33/33 [=====] - 0s 1ms/step - loss: 0.1275 - accuracy: 0.9596
Epoch 32/100
33/33 [=====] - 0s 1ms/step - loss: 0.1232 - accuracy: 0.9596
Epoch 33/100
33/33 [=====] - 0s 1ms/step - loss: 0.1220 - accuracy: 0.9596
Epoch 34/100
33/33 [=====] - 0s 1ms/step - loss: 0.1209 - accuracy: 0.9617
Epoch 35/100
33/33 [=====] - 0s 1ms/step - loss: 0.1176 - accuracy: 0.9606
Epoch 36/100
33/33 [=====] - 0s 1ms/step - loss: 0.1159 - accuracy: 0.9617
Epoch 37/100
33/33 [=====] - 0s 1ms/step - loss: 0.1148 - accuracy: 0.9617
Epoch 38/100
33/33 [=====] - 0s 1ms/step - loss: 0.1153 - accuracy: 0.9627
Epoch 39/100
33/33 [=====] - 0s 1ms/step - loss: 0.1119 - accuracy: 0.9585
Epoch 40/100
33/33 [=====] - 0s 1ms/step - loss: 0.1118 - accuracy: 0.9606
Epoch 41/100
33/33 [=====] - 0s 1ms/step - loss: 0.1089 - accuracy: 0.9606
Epoch 42/100
33/33 [=====] - 0s 1ms/step - loss: 0.1126 - accuracy: 0.9617
Epoch 43/100
33/33 [=====] - 0s 1ms/step - loss: 0.1098 - accuracy: 0.9606
Epoch 44/100
33/33 [=====] - 0s 1ms/step - loss: 0.1098 - accuracy: 0.9627
Epoch 45/100
33/33 [=====] - 0s 1ms/step - loss: 0.1055 - accuracy: 0.9606
Epoch 46/100
33/33 [=====] - 0s 1ms/step - loss: 0.1055 - accuracy: 0.9606
Epoch 47/100
33/33 [=====] - 0s 1ms/step - loss: 0.1060 - accuracy: 0.9606
Epoch 48/100
33/33 [=====] - 0s 1ms/step - loss: 0.1050 - accuracy: 0.9627
Epoch 49/100
33/33 [=====] - 0s 1ms/step - loss: 0.1084 - accuracy: 0.9606
Epoch 50/100
33/33 [=====] - 0s 1ms/step - loss: 0.1032 - accuracy: 0.9617
Epoch 51/100
33/33 [=====] - 0s 1ms/step - loss: 0.1009 - accuracy: 0.9648
Epoch 52/100
33/33 [=====] - 0s 1ms/step - loss: 0.1054 - accuracy: 0.9617
Epoch 53/100
33/33 [=====] - 0s 1ms/step - loss: 0.1000 - accuracy: 0.9637
Epoch 54/100
33/33 [=====] - 0s 1ms/step - loss: 0.0967 - accuracy: 0.9648
Epoch 55/100
33/33 [=====] - 0s 1ms/step - loss: 0.0962 - accuracy: 0.9648
Epoch 56/100
33/33 [=====] - 0s 1ms/step - loss: 0.0966 - accuracy: 0.9637
Epoch 57/100
33/33 [=====] - 0s 1ms/step - loss: 0.0943 - accuracy: 0.9648
Epoch 58/100
33/33 [=====] - 0s 1ms/step - loss: 0.0944 - accuracy: 0.9637
Epoch 59/100
33/33 [=====] - 0s 1ms/step - loss: 0.0916 - accuracy: 0.9668
Epoch 60/100
33/33 [=====] - 0s 1ms/step - loss: 0.0935 - accuracy: 0.9668
Epoch 61/100
33/33 [=====] - 0s 1ms/step - loss: 0.0911 - accuracy: 0.9648
Epoch 62/100
33/33 [=====] - 0s 1ms/step - loss: 0.0900 - accuracy: 0.9668
Epoch 63/100
33/33 [=====] - 0s 1ms/step - loss: 0.0896 - accuracy: 0.9658
Epoch 64/100
33/33 [=====] - 0s 1ms/step - loss: 0.0955 - accuracy: 0.9648
Epoch 65/100
33/33 [=====] - 0s 1ms/step - loss: 0.0867 - accuracy: 0.9648
Epoch 66/100
33/33 [=====] - 0s 1ms/step - loss: 0.0866 - accuracy: 0.9658
Epoch 67/100
33/33 [=====] - 0s 1ms/step - loss: 0.0900 - accuracy: 0.9637
Epoch 68/100
33/33 [=====] - 0s 1ms/step - loss: 0.0872 - accuracy: 0.9637
Epoch 69/100
33/33 [=====] - 0s 1ms/step - loss: 0.0852 - accuracy: 0.9668
Epoch 70/100
33/33 [=====] - 0s 1ms/step - loss: 0.0885 - accuracy: 0.9658
Epoch 71/100
33/33 [=====] - 0s 1ms/step - loss: 0.0835 - accuracy: 0.9658
Epoch 72/100
33/33 [=====] - 0s 1ms/step - loss: 0.0867 - accuracy: 0.9637
Epoch 73/100
33/33 [=====] - 0s 1ms/step - loss: 0.0825 - accuracy: 0.9648
Epoch 74/100
33/33 [=====] - 0s 1ms/step - loss: 0.0815 - accuracy: 0.9658
Epoch 75/100
33/33 [=====] - 0s 1ms/step - loss: 0.0815 - accuracy: 0.9658
Epoch 76/100
33/33 [=====] - 0s 1ms/step - loss: 0.0821 - accuracy: 0.9658
Epoch 77/100
33/33 [=====] - 0s 1ms/step - loss: 0.0791 - accuracy: 0.9689
Epoch 78/100
33/33 [=====] - 0s 1ms/step - loss: 0.0815 - accuracy: 0.9658
Epoch 79/100
33/33 [=====] - 0s 1ms/step - loss: 0.0794 - accuracy: 0.9658
Epoch 80/100
33/33 [=====] - 0s 1ms/step - loss: 0.0802 - accuracy: 0.9689
Epoch 81/100
33/33 [=====] - 0s 1ms/step - loss: 0.0793 - accuracy: 0.9658
Epoch 82/100
33/33 [=====] - 0s 1ms/step - loss: 0.0780 - accuracy: 0.9699
Epoch 83/100
33/33 [=====] - 0s 1ms/step - loss: 0.0782 - accuracy: 0.9689
Epoch 84/100
33/33 [=====] - 0s 1ms/step - loss: 0.0750 - accuracy: 0.9648
Epoch 85/100
33/33 [=====] - 0s 1ms/step - loss: 0.0793 - accuracy: 0.9689
Epoch 86/100
33/33 [=====] - 0s 1ms/step - loss: 0.0743 - accuracy: 0.9658
Epoch 87/100
33/33 [=====] - 0s 1ms/step - loss: 0.0782 - accuracy: 0.9668
Epoch 88/100
33/33 [=====] - 0s 1ms/step - loss: 0.0765 - accuracy: 0.9668
Epoch 89/100
33/33 [=====] - 0s 1ms/step - loss: 0.0763 - accuracy: 0.9689
Epoch 90/100
33/33 [=====] - 0s 1ms/step - loss: 0.0778 - accuracy: 0.9679
Epoch 91/100
33/33 [=====] - 0s 1ms/step - loss: 0.0737 - accuracy: 0.9679
Epoch 92/100
33/33 [=====] - 0s 1ms/step - loss: 0.0778 - accuracy: 0.9668
Epoch 93/100
33/33 [=====] - 0s 1ms/step - loss: 0.0726 - accuracy: 0.9699
Epoch 94/100
33/33 [=====] - 0s 1ms/step - loss: 0.0725 - accuracy: 0.9689
Epoch 95/100
33/33 [=====] - 0s 1ms/step - loss: 0.0704 - accuracy: 0.9710
Epoch 96/100
33/33 [=====] - 0s 1ms/step - loss: 0.0728 - accuracy: 0.9679
Epoch 97/100
33/33 [=====] - 0s 1ms/step - loss: 0.0708 - accuracy: 0.9689
Epoch 98/100
33/33 [=====] - 0s 1ms/step - loss: 0.0699 - accuracy: 0.9689
Epoch 99/100
33/33 [=====] - 0s 1ms/step - loss: 0.0735 - accuracy: 0.9699
Epoch 100/100
33/33 [=====] - 0s 1ms/step - loss: 0.0708 - accuracy: 0.9699
Epoch 100/100
33/33 [=====] - 0s 1ms/step - loss: 0.0699 - accuracy: 0.9699
```

Out[165]: <tensorflow.python.keras.callbacks.History at 0x1dfc6c0c8>

Multi Layer Perceptron Evaluation Metrics

```
In [166]: model1.evaluate(X_train, y_train)

31/31 [=====] - 0s 740us/step - loss: 0.0652 - accuracy: 0.9720

Out[166]: [0.0652037039399147, 0.9720207452774048]

In [167]: model1.evaluate(X_test, y_test)

8/8 [=====] - 0s 1ms/step - loss: 0.4219 - accuracy: 0.9298

Out[167]: [0.42187103629112244, 0.9297520518302917]
```

Multi Layer Perceptron Precision Scores

```
In [168]: train_pred=model1.predict_classes(X_train)
          pred=model1.predict_classes(X_test)

In [169]: print("train",precision_score(y_train, train_pred))
          print("test",precision_score(y_test,pred))

train 0.941320293398533
test 0.8936170212765957

In [170]: from sklearn.metrics import recall_score
          print("train",recall_score(y_train, train_pred))
          print("test",recall_score(y_test,pred))

train 0.9922680412371134
test 0.9230769230769231
```