

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import pyplot as plt
from sklearn.preprocessing import LabelEncoder
from sklearn.decomposition import PCA
from scipy.stats import norm
import seaborn as sns
from sklearn.model_selection import train_test_split
df = pd.read_csv('C:\Users\guruc\Downloads\cancer.csv')
df.head()
df.to_csv('data/data_clean_id.csv')
df.drop('id', axis=1, inplace=True)
plt.rcParams['figure.figsize'] = (15,8)
plt.rcParams['axes.titlesize'] = 'large'
df = pd.read_csv('data/data_clean_id.csv', index_col=False)
df.drop('Unnamed: 0',axis=1, inplace=True)
df.head(3)
df.describe()
df.head(3)
df.shape
df.info()
df.isnull().any()
df.diagnosis.unique()

plt.style.use('fivethirtyeight')
sns.set_style("white")

plt.rcParams['figure.figsize'] = (8,4)
```

```
df = pd.read_csv('data/data_clean.csv', index_col=False)
```

```
df.drop('Unnamed: 0',axis=1, inplace=True)
```

```
df.head(3)
```

```
array = df.values
```

```
X = array[:,1:31]
```

```
y = array[:,0]
```

```
X
```

```
le = LabelEncoder()
```

```
y = le.fit_transform(y)
```

```
y
```

```
df.diagnosis.unique()
```

```
array(['M', 'B'], dtype=object)
```

```
diag_gr = df.groupby('diagnosis', axis=0)
```

```
pd.DataFrame(diag_gr.size(), columns=['# of observations'])
```

```
sns.set_style("white")
```

```
sns.set_context({"figure.figsize": (10, 8)})
```

```
sns.countplot(df['diagnosis'],label='Count',palette="Set3")
```

```
plt.style.use('fivethirtyeight')
```

```
sns.set_style("white")
```

```
df = pd.read_csv('data/data_clean.csv', index_col=False)
```

```
df.drop('Unnamed: 0',axis=1, inplace=True)
```

```
corr = df_mean.corr()
```

```
mask = np.zeros_like(corr, dtype=np.bool)
```

```
mask[np.triu_indices_from(mask)] = True
```

```
df, ax = plt.subplots(figsize=(8, 8))
```

```
plt.title('Breast Cancer Feature Correlation')
```

```
cmap = sns.diverging_palette(260, 10, as_cmap=True)
```

```
sns.heatmap(corr, vmax=1.2, square='square', cmap=cmap, mask=mask,
```

```
ax=ax,annot=True, fmt='.2g',linewidths=2)
```

```
plt.style.use('fivethirtyeight')
```

```
sns.set_style("white")
```

```
df = pd.read_csv('data/data_clean.csv', index_col=False)
```

```
g = sns.PairGrid(df[[df.columns[1],df.columns[2], df.columns[3],
```

```
df.columns[4], df.columns[5], df.columns[6]]], hue='diagnosis')
```

```
g = g.map_diag(plt.hist)
```

```
g = g.map_offdiag(plt.scatter, s = 3)
```

```
X_train, X_test, y_train, y_test = train_test_split( X, y, test_size=0.25, random_state=7)
```

```
X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

```
pca = PCA(n_components=10)
```

```
fit = pca.fit(Xs)
```

```
X_pca = pca.transform(Xs)
```

```
PCA_df = pd.DataFrame()
```

```
PCA_df['PCA_1'] = X_pca[:,0]
```

```
PCA_df['PCA_2'] = X_pca[:,1]
```

```
plt.plot(PCA_df['PCA_1'][df.diagnosis == 'M'],PCA_df['PCA_2'][df.diagnosis == 'M'],'o', alpha = 0.7,  
color = 'r')
```

```
plt.plot(PCA_df['PCA_1'][df.diagnosis == 'B'],PCA_df['PCA_2'][df.diagnosis == 'B'],'o', alpha = 0.7,  
color = 'b')
```

```
plt.xlabel('PCA_1')
plt.ylabel('PCA_2')
plt.legend(['Malignant','Benign'])
plt.show()
```

```
var= pca.explained_variance_ratio_
```

```
plt.plot(var)
plt.title('Scree Plot')
plt.xlabel('Principal Component')
plt.ylabel('Eigenvalue')
```

```
leg = plt.legend(['Eigenvalues from PCA'], loc='best', borderpad=0.3, shadow=False, markerscale=0.4)
leg.get_frame().set_alpha(0.4)
leg.set_draggable(state=True)
plt.show()
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.model_selection import cross_val_score
from sklearn.pipeline import make_pipeline
from sklearn.metrics import confusion_matrix
from sklearn import metrics, preprocessing
from sklearn.metrics import classification_report
```

```
import seaborn as sns
```

```
plt.style.use('fivethirtyeight')
sns.set_style("white")
```

```
plt.rcParams['figure.figsize'] = (8,4)
```

```
df = pd.read_csv('data/data_clean.csv', index_col=False)
```

```
df.drop('Unnamed: 0',axis=1, inplace=True)
```

```
df.head(3)
```

```
array = df.values
```

```
X = array[:,1:31]
```

```
y = array[:,0]
```

```
le = LabelEncoder()
```

```
y = le.fit_transform(y)
```

```
scaler = StandardScaler()
```

```
Xs = scaler.fit_transform(X)
```

```
X_train, X_test, y_train, y_test = train_test_split(Xs, y, test_size=0.3, random_state=2, stratify=y)
```

```
clf = SVC(probability=True)
```

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

```
classifier_score = clf.score(X_test, y_test)
```

```
print ('\n→ The classifier accuracy score is {:.2f}\n'.format(classifier_score))
```

```
n_folds = 3
```

```
cv_error = np.average(cross_val_score(SVC(), Xs, y, cv=n_folds))
```

```
print ('\n→ The {}-fold cross-validation accuracy score for this classifier is {:.2f}\n'.format(n_folds, cv_error))
```

```
from sklearn.feature_selection import SelectKBest, f_regression
```

```
clf2 = make_pipeline(SelectKBest(f_regression, k=3),SVC(probability=True))
```

```
scores = cross_val_score(clf2, Xs, y, cv=3)
```

```

n_folds = 3

cv_error = np.average(cross_val_score(SVC(), Xs, y, cv=n_folds))

print('\n→ The {}-fold cross-validation accuracy score for this classifier is {:.2f}\n'.format(n_folds,
cv_error))


print (scores)

avg = (100*np.mean(scores), 100*np.std(scores)/np.sqrt(scores.shape[0]))

print ("→ Average score and uncertainty: (%.2f +- %.3f)%%"%avg)

from IPython.display import Image, display

y_pred = clf.fit(X_train, y_train).predict(X_test)

cm = metrics.confusion_matrix(y_test, y_pred)

fig, ax = plt.subplots(figsize=(5, 5))

ax.matshow(cm, cmap=plt.cm.Reds, alpha=0.3)

for i in range(cm.shape[0]):
    for j in range(cm.shape[1]):
        ax.text(x=j, y=i,
                s=cm[i, j],
                va='center', ha='center')

plt.xlabel('Predicted Values', )
plt.ylabel('Actual Values')

plt.show()

print(classification_report(y_test, y_pred ))


plt.figure(figsize=(10,8))

probas_ = clf.predict_proba(X_test)

fpr, tpr, thresholds = roc_curve(y_test, probas_[ :, 1])

roc_auc = auc(fpr, tpr)

plt.plot(fpr, tpr, lw=1, label='ROC fold (area = %0.2f)' % (roc_auc))

plt.plot([0, 1], [0, 1], '--', color=(0.6, 0.6, 0.6), label='Random')

```

```
plt.xlim([-0.05, 1.05])
plt.ylim([-0.05, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.axes().set_aspect(1)
```

```
plt.style.use('fivethirtyeight')
sns.set_style("white")
```

```
plt.rcParams['figure.figsize'] = (8,4)
df = pd.read_csv('data/data_clean.csv', index_col=False)
df.drop('Unnamed: 0',axis=1, inplace=True)
```

```
array = df.values
X = array[:,1:31]
y = array[:,0]
```

```
le = LabelEncoder()
y = le.fit_transform(y)
```

```
scaler =StandardScaler()
Xs = scaler.fit_transform(X)
```

```
from sklearn.decomposition import PCA
pca = PCA(n_components=10)
fit = pca.fit(Xs)
X_pca = pca.transform(Xs)
```

```
X_train, X_test, y_train, y_test = train_test_split(X_pca, y, test_size=0.3, random_state=2, stratify=y)
```

```
clf = SVC(probability=True)
```

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

```
classifier_score = clf.score(X_test, y_test)
```

```
print ('\n\nThe classifier accuracy score is {:.03.2f}\n'.format(classifier_score))
```

```
clf2 = make_pipeline(SelectKBest(f_regression, k=3),SVC(probability=True))
```

```
scores = cross_val_score(clf2, X_pca, y, cv=3)
```

```
n_folds = 5
```

```
cv_error = np.average(cross_val_score(SVC(), X_pca, y, cv=n_folds))
```

```
y_pred = clf.fit(X_train, y_train).predict(X_test)
```

```
cm = metrics.confusion_matrix(y_test, y_pred)
```

```
print(classification_report(y_test, y_pred ))
```

```
fig, ax = plt.subplots(figsize=(5, 5))
```

```
ax.matshow(cm, cmap=plt.cm.Red, alpha=0.3)
```

```
for i in range(cm.shape[0]):
```

```
    for j in range(cm.shape[1]):
```

```
        ax.text(x=j, y=i,
```

```
                s=cm[i, j],
```

```
                va='center', ha='center')
```

```
plt.xlabel('Predicted Values', )
```

```
plt.ylabel('Actual Values')
```

```
plt.show()
```



```

kernel_values = [ 'linear' , 'poly' , 'rbf' , 'sigmoid' ]
param_grid = {'C': np.logspace(-3, 2, 6), 'gamma': np.logspace(-3, 2, 6), 'kernel': kernel_values}

grid = GridSearchCV(SVC(), param_grid=param_grid, cv=5)
grid.fit(X_train, y_train)
print("The best parameters are %s with a score of %0.2f"
      % (grid.best_params_, grid.best_score_))
grid.best_estimator_.probability = True
clf = grid.best_estimator_
y_pred = clf.fit(X_train, y_train).predict(X_test)
cm = metrics.confusion_matrix(y_test, y_pred)

print(classification_report(y_test, y_pred ))

fig, ax = plt.subplots(figsize=(5, 5))
ax.matshow(cm, cmap=plt.cm.Reds, alpha=0.3)
for i in range(cm.shape[0]):
    for j in range(cm.shape[1]):
        ax.text(x=j, y=i,
                s=cm[i, j],
                va='center', ha='center')
plt.xlabel('Predicted Values', )
plt.ylabel('Actual Values')
plt.show()

def decision_plot(X_train, y_train, n_neighbors, weights):
    h = .02

Xtrain = X_train[:, :2]

```

```
cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])
```

```
cmap_bold = ListedColormap(['#FF0000', '#00FF00', '#0000FF'])
```

```
C = 1.0
```

```
svm = SVC(kernel='linear', random_state=0, gamma=0.1, C=C).fit(Xtrain, y_train)
```

```
rbf_svc = SVC(kernel='rbf', gamma=0.7, C=C).fit(Xtrain, y_train)
```

```
poly_svc = SVC(kernel='poly', degree=3, C=C).fit(Xtrain, y_train)
```

```
%matplotlib inline
```

```
plt.rcParams['figure.figsize'] = (15, 9)
```

```
plt.rcParams['axes.titlesize'] = 'large'
```

```
x_min, x_max = Xtrain[:, 0].min() - 1, Xtrain[:, 0].max() + 1
```

```
y_min, y_max = Xtrain[:, 1].min() - 1, Xtrain[:, 1].max() + 1
```

```
xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),  
                     np.arange(y_min, y_max, 0.1))
```

```
titles = ['SVC with linear kernel',
```

```
         'SVC with RBF kernel',
```

```
         'SVC with polynomial (degree 3) kernel']
```

```
for i, clf in enumerate((svm, rbf_svc, poly_svc)):
```

```
    plt.subplot(2, 2, i + 1)
```

```
    plt.subplots_adjust(wspace=0.4, hspace=0.4)
```

```
Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
```

```
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
```

```
plt.scatter(Xtrain[:, 0], Xtrain[:, 1], c=y_train, cmap=plt.cm.coolwarm)
plt.xlabel('radius_mean')
plt.ylabel('texture_mean')
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
plt.xticks(())
plt.yticks(())
plt.title(titles[i])
```

```
plt.show()
```

```
#Automation
```

```
%matplotlib inline
```

```
import matplotlib.pyplot as plt
```

```
# Create a pipeline that standardizes the data then creates a model
```

```
#Load libraries for data processing
```

```
import pandas as pd #data processing, CSV file I/O (e.g. pd.read_csv)
```

```
import numpy as np
```

```
from scipy.stats import norm
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.model_selection import cross_val_score
```

```
from sklearn.model_selection import KFold
```

```
from sklearn.preprocessing import LabelEncoder
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report

# visualization
import seaborn as sns

plt.style.use('fivethirtyeight')
sns.set_style("white")

plt.rcParams['figure.figsize'] = (8,4)
#plt.rcParams['axes.titlesize'] = 'large'

#load data
df = pd.read_csv('data/data_clean.csv', index_col=False)
df.drop('Unnamed: 0',axis=1, inplace=True)

# Split-out validation dataset
array = df.values
X = array[:,1:31]
```

```
y = array[:,0]
```

```
# Divide records in training and testing sets.
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=7)
```

```
#transform the class labels from their original string representation (M and B) into integers
```

```
le = LabelEncoder()
```

```
y = le.fit_transform(y)
```

```
Evaluate Algorithms: Baseline
```

```
# Spot-Check Algorithms
```

```
models = []
```

```
models.append(('LR', LogisticRegression()))
```

```
models.append(('LDA', LinearDiscriminantAnalysis()))
```

```
models.append(('KNN', KNeighborsClassifier()))
```

```
models.append(('CART', DecisionTreeClassifier()))
```

```
models.append(('NB', GaussianNB()))
```

```
models.append(('SVM', SVC()))
```

```
# Test options and evaluation metric
```

```
num_folds = 10
```

```
num_instances = len(X_train)
```

```
seed = 7
```

```
scoring = 'accuracy'
```

```
# Test options and evaluation metric
```

```
num_folds = 10
```

```
num_instances = len(X_train)
```

```
seed = 7
```

```
scoring = 'accuracy'
```

```
results = []
```

```
names = []
```

```

for name, model in models:

    kfold = KFold(n_splits=num_folds, random_state=seed, shuffle=True)

    cv_results = cross_val_score(model, X_train, y_train, cv=kfold, scoring=scoring)

    results.append(cv_results)

    names.append(name)

    msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())

    print(msg)

print('\n→ 10-Fold cross-validation accuray score for the training data for six classifiers')

len(X_train)

# Compare Algorithms

fig = plt.figure()
fig.suptitle('Algorithm Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()

pipelines = []

pipelines.append(('ScaledLR', Pipeline([('Scaler', StandardScaler()), ('LR', LogisticRegression())]))

pipelines.append(('ScaledLDA', Pipeline([('Scaler', StandardScaler()), ('LDA',
LinearDiscriminantAnalysis())]))

pipelines.append(('ScaledKNN', Pipeline([('Scaler', StandardScaler()), ('KNN',
KNeighborsClassifier())]))

pipelines.append(('ScaledCART', Pipeline([('Scaler', StandardScaler()), ('CART',
DecisionTreeClassifier())]))

pipelines.append(('ScaledNB', Pipeline([('Scaler', StandardScaler()), ('NB', GaussianNB())]))

pipelines.append(('ScaledSVM', Pipeline([('Scaler', StandardScaler()), ('SVM', SVC())]))

results = []

names = []

for name, model in pipelines:

```

```

kfold = KFold(n_splits=num_folds, random_state=seed, shuffle=True)
cv_results = cross_val_score(model, X_train, y_train, cv=kfold,
    scoring=scoring)
results.append(cv_results)
names.append(name)
msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())
print(msg)
fig = plt.figure()
fig.suptitle('Scaled Algorithm Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()

#Make Support Vector Classifier Pipeline
pipe_svc = Pipeline([('scl', StandardScaler()),
    ('pca', PCA(n_components=2)),
    ('clf', SVC(probability=True, verbose=False))])

#Fit Pipeline to training Data
pipe_svc.fit(X_train, y_train)

#print('→ Fitted Pipeline to training Data')

scores = cross_val_score(estimator=pipe_svc, X=X_train, y=y_train, cv=10, n_jobs=1, verbose=0)
print('→ Model Training Accuracy: %.3f +/- %.3f' % (np.mean(scores), np.std(scores)))

#Tune Hyperparameters
param_range = [0.0001, 0.001, 0.01, 0.1, 1.0, 10.0, 100.0, 1000.0]
param_grid = [{'clf__C': param_range, 'clf__kernel': ['linear']},

```

```

        {'clf__C': param_range, 'clf__gamma': param_range,
         'clf__kernel': ['rbf']}]

gs = GridSearchCV(estimator=pipe_svc,
                  param_grid=param_grid,
                  scoring='accuracy',
                  cv=10,
                  n_jobs=1)

gs = gs.fit(X_train, y_train)

print('→ Tuned Parameters Best Score: ', gs.best_score_)

print('→ Best Parameters: \n', gs.best_params_)

from sklearn.neighbors import KNeighborsClassifier as KNN

pipe_knn = Pipeline([('scl', StandardScaler()),
                     ('pca', PCA(n_components=2)),
                     ('clf', KNeighborsClassifier())])

#Fit Pipeline to training Data
pipe_knn.fit(X_train, y_train)

scores = cross_val_score(estimator=pipe_knn,
                          X=X_train,
                          y=y_train,
                          cv=10,
                          n_jobs=1)

print('→ Model Training Accuracy: %.3f +/- %.3f' %(np.mean(scores), np.std(scores)))

#Tune Hyperparameters

param_range = range(1, 31)

param_grid = [{'clf__n_neighbors': param_range}]

# instantiate the grid

gs = GridSearchCV(estimator=pipe_knn,

```



```

        param_grid=param_grid,
        cv=10,
        scoring='accuracy')

gs = gs.fit(X_train, y_train)

print('→ Tuned Parameters Best Score: ',gs.best_score_)

print('→ Best Parameters: \n',gs.best_params_)


#Use best parameters

clf_svc = gs.best_estimator_


#Get Final Scores

clf_svc.fit(X_train, y_train)

scores = cross_val_score(estimator=clf_svc,
                          X=X_train,
                          y=y_train,
                          cv=10,
                          n_jobs=1)


print('→ Final Model Training Accuracy: %.3f +/- %.3f' %(np.mean(scores), np.std(scores)))

print('→ Final Accuracy on Test set: %.5f' % clf_svc.score(X_test,y_test))

clf_svc.fit(X_train, y_train)

y_pred = clf_svc.predict(X_test)


print(accuracy_score(y_test, y_pred))

print(confusion_matrix(y_test, y_pred))

print(classification_report(y_test, y_pred))

```