

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import KFold
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
import time
```

Exploratory analysis

```
In [4]: data = pd.read_csv("C:/Users/Mounika/Downloads/data.csv")
data.head(5)

Out[4]:
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	...	texture_worst	perimeter_worst	area_worst
0	842902	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	23.41	158.80	1956.0
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0

5 rows × 33 columns

```
In [5]: print(data.shape)
(569, 33)

In [6]: data.describe()
```

	id	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	...	texture_worst	perimeter_worst	area_worst
count	5.690000e+02	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.000000	569.0000	569.000
mean	3.037183e+07	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.088799	0.048919	0.181162	...	25.677223	107.261	107.261
std	1.250206e+08	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079720	0.038803	0.027414	...	6.146258	33.602	33.602
min	8.670000e+03	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000	0.106000	...	12.020000	50.410	50.410
25%	8.692180e+05	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310	0.161900	...	21.080000	84.110	84.110
50%	9.060240e+05	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061540	0.033500	0.179200	...	25.410000	97.660	97.660
75%	8.813129e+06	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074000	0.195700	...	29.720000	125.400	125.400
max	9.113205e+08	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200	0.304000	...	49.540000	251.200	251.200

8 rows × 32 columns

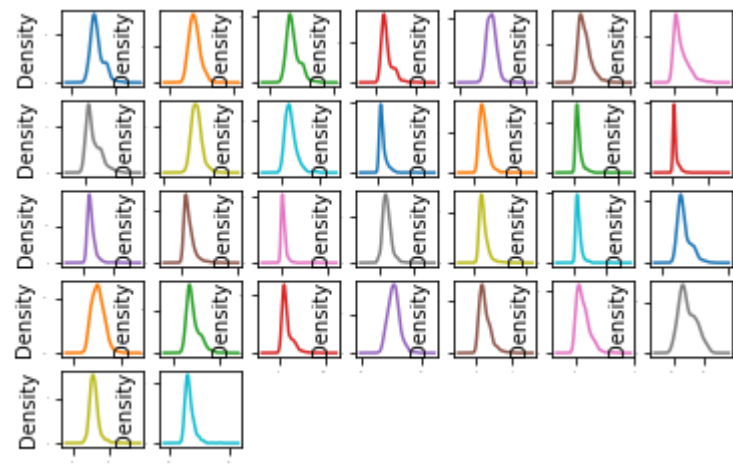
Data visualisation and pre-processing

```
In [7]: data['diagnosis'] = data['diagnosis'].apply(lambda x: '1' if x == 'M' else '0')
data = data.set_index('id')
del data['Unnamed: 32']

In [8]: print(data.groupby('diagnosis').size())

diagnosis
0    357
1    212
dtype: int64

In [9]: data.plot(kind='density', subplots=True, layout=(5,7), sharex=False, legend=False, fontsize=1)
plt.show()
```



```
In [10]: from matplotlib import cm as cm

fig = plt.figure()
ax1 = fig.add_subplot(111)
cmap = cm.get_cmap('jet', 30)
cax = ax1.imshow(data.corr(), interpolation="none", cmap=cmap)
ax1.grid(True)
plt.title('Breast Cancer Attributes Correlation')
# Add colorbar, make sure to specify tick locations to match desired ticklabels
fig.colorbar(cax, ticks=[.75,.8,.85,.9,.95,1])
plt.show()

In [11]: Y = data['diagnosis'].values
X = data.drop('diagnosis', axis=1).values

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.20, random_state=21)
```

Baseline algorithm checking

```
In [15]: models_list = []
models_list.append(('CART', DecisionTreeClassifier()))
models_list.append(('SVM', SVC()))
models_list.append(('NB', GaussianNB()))
models_list.append(('KNN', KNeighborsClassifier()))
num_folds = 10
results = []
names = []

for name, model in models_list:
    kfold = KFold(n_splits=num_folds, random_state=123)
    start = time.time()
    cv_results = cross_val_score(model, X_train, Y_train, cv=kfold, scoring='accuracy')
    end = time.time()
    results.append(cv_results)
    names.append(name)
    print("%s: %f (%f) (run time: %f)" % (name, cv_results.mean(), cv_results.std(), end-start))

C:\Users\Mounika\anaconda\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
C:\Users\Mounika\anaconda\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
C:\Users\Mounika\anaconda\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
CART: 0.921063 (0.039453) (run time: 0.108120)
SVM: 0.907681 (0.054723) (run time: 0.060003)
NB: 0.940773 (0.033921) (run time: 0.036790)
KNN: 0.927729 (0.055250) (run time: 0.076570)
C:\Users\Mounika\anaconda\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
```

```
In [16]: fig = plt.figure()
fig.suptitle('Performance Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()

Performance Comparison
```

Evaluation of algorithm on Standardised Data

```
In [17]: import warnings

# Standardize the dataset
pipelines = []

pipelines.append(('ScaledCART', Pipeline([('Scaler', StandardScaler()),('CART', DecisionTreeClassifier())]))
pipelines.append(('ScaledSVM', Pipeline([('Scaler', StandardScaler()),('SVM', SVC())]))
pipelines.append(('ScaledNB', Pipeline([('Scaler', StandardScaler()),('NB', GaussianNB())]))
pipelines.append(('ScaledKNN', Pipeline([('Scaler', StandardScaler()),('KNN', KNeighborsClassifier())]))

results = []
names = []
with warnings.catch_warnings():
    warnings.simplefilter("ignore")
    kfold = KFold(n_splits=num_folds, random_state=123)
    for name, model in pipelines:
        start = time.time()
        cv_results = cross_val_score(model, X_train, Y_train, cv=kfold, scoring='accuracy')
        end = time.time()
        results.append(cv_results)
        names.append(name)
        print("%s: %f (%f) (run time: %f)" % (name, cv_results.mean(), cv_results.std(), end-start))

ScaledCART: 0.916473 (0.035364) (run time: 0.165991)
ScaledSVM: 0.964879 (0.038621) (run time: 0.093325)
ScaledNB: 0.931932 (0.038625) (run time: 0.043218)
ScaledKNN: 0.958357 (0.038595) (run time: 0.125401)
```

```
In [18]: fig = plt.figure()
fig.suptitle('Performance Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()

Performance Comparison
```

Algorithm Tuning - Tuning SVM

```
In [19]: scaler = StandardScaler().fit(X_train)
rescaledx = scaler.transform(X_train)
c_values = [0.1, 0.3, 0.5, 0.7, 0.9, 1.0, 1.3, 1.5, 1.7, 2.0]
kernel_values = ['linear', 'poly', 'rbf', 'sigmoid']
param_grid = dict(C=c_values, kernel=kernel_values)
model = SVC()
kfold = KFold(n_splits=num_folds, random_state=21)
grid = GridSearchCV(estimator=model, param_grid=param_grid, scoring='accuracy', cv=kfold)
grid_result = grid.fit(rescaledx, Y_train)
print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
means = grid_result.cv_results_['mean_test_score']
stds = grid_result.cv_results_['std_test_score']
params = grid_result.cv_results_['params']
for mean, stdev, param in zip(means, stds, params):
    print("%f (%f) with: %r" % (mean, stdev, param))

C:\Users\Mounika\anaconda\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
Best: 0.969324 using {'C': 2.0, 'kernel': 'rbf'}
0.964976 (0.026211) with: {'C': 0.1, 'kernel': 'linear'}
0.828551 (0.054707) with: {'C': 0.1, 'kernel': 'poly'}
0.940725 (0.038380) with: {'C': 0.1, 'kernel': 'rbf'}
0.949469 (0.032899) with: {'C': 0.1, 'kernel': 'sigmoid'}
0.962754 (0.029531) with: {'C': 0.3, 'kernel': 'linear'}
0.863720 (0.050997) with: {'C': 0.3, 'kernel': 'poly'}
0.956039 (0.032908) with: {'C': 0.3, 'kernel': 'linear'}
0.960386 (0.029499) with: {'C': 0.3, 'kernel': 'sigmoid'}
0.956184 (0.030988) with: {'C': 0.5, 'kernel': 'linear'}
0.879034 (0.053507) with: {'C': 0.5, 'kernel': 'poly'}
0.964879 (0.030954) with: {'C': 0.5, 'kernel': 'rbf'}
0.956087 (0.027848) with: {'C': 0.5, 'kernel': 'sigmoid'}
0.954010 (0.031644) with: {'C': 0.7, 'kernel': 'linear'}
0.885604 (0.038275) with: {'C': 0.7, 'kernel': 'poly'}
0.967053 (0.037461) with: {'C': 0.7, 'kernel': 'rbf'}
0.949565 (0.027831) with: {'C': 0.7, 'kernel': 'sigmoid'}
0.951836 (0.028830) with: {'C': 0.9, 'kernel': 'linear'}
0.887826 (0.039039) with: {'C': 0.9, 'kernel': 'poly'}
0.967053 (0.037461) with: {'C': 0.9, 'kernel': 'rbf'}
0.947391 (0.032846) with: {'C': 0.9, 'kernel': 'sigmoid'}
0.954010 (0.026552) with: {'C': 1.0, 'kernel': 'linear'}
0.890048 (0.038398) with: {'C': 1.0, 'kernel': 'poly'}
0.967101 (0.033160) with: {'C': 1.0, 'kernel': 'rbf'}
0.947391 (0.032846) with: {'C': 1.0, 'kernel': 'sigmoid'}
0.956184 (0.025767) with: {'C': 1.3, 'kernel': 'linear'}
0.894396 (0.039509) with: {'C': 1.3, 'kernel': 'poly'}
0.967150 (0.028285) with: {'C': 1.3, 'kernel': 'rbf'}
0.947391 (0.029759) with: {'C': 1.3, 'kernel': 'sigmoid'}
0.958357 (0.024768) with: {'C': 1.5, 'kernel': 'linear'}
0.898792 (0.033329) with: {'C': 1.5, 'kernel': 'poly'}
0.967150 (0.028285) with: {'C': 1.5, 'kernel': 'rbf'}
0.940821 (0.039361) with: {'C': 1.5, 'kernel': 'sigmoid'}
0.956135 (0.021744) with: {'C': 1.7, 'kernel': 'linear'}
0.900966 (0.034861) with: {'C': 1.7, 'kernel': 'poly'}
0.967150 (0.024547) with: {'C': 1.7, 'kernel': 'rbf'}
0.940918 (0.037900) with: {'C': 1.7, 'kernel': 'sigmoid'}
0.956135 (0.021744) with: {'C': 1.7, 'kernel': 'linear'}
0.907536 (0.034327) with: {'C': 2.0, 'kernel': 'poly'}
0.969324 (0.022458) with: {'C': 2.0, 'kernel': 'rbf'}
0.932029 (0.028300) with: {'C': 2.0, 'kernel': 'sigmoid'}
```

Application of SVC on dataset

```
In [20]: # prepare the model
with warnings.catch_warnings():
    warnings.simplefilter("ignore")
    scaler = StandardScaler().fit(X_train)
    X_train_scaled = scaler.transform(X_train)
    model = SVC(C=2.0, kernel='rbf')
    start = time.time()
    model.fit(X_train_scaled, Y_train)
    end = time.time()
    print("Run Time: %f" % (end-start))

Run Time: 0.006001

In [21]: # estimate accuracy on test dataset
with warnings.catch_warnings():
    warnings.simplefilter("ignore")
    X_test_scaled = scaler.transform(X_test)
    predictions = model.predict(X_test_scaled)

In [22]: print("Accuracy score %f" % accuracy_score(Y_test, predictions))
print(classification_report(Y_test, predictions))

Accuracy score 0.991228
precision      recall  f1-score   support

0             1.00      0.99      0.99         1
1             0.97      1.00      0.99         39

accuracy      0.99      0.99      0.99      114
macro avg     0.99      0.99      0.99      114
weighted avg  0.99      0.99      0.99      114

In [23]: print(confusion_matrix(Y_test, predictions))

[[74  1]
 [ 0 39]]

In [ ]:
```