

CS 634 Data Mining Final term Project

Option 1: Supervised Data Mining (Classification)

Submitted by: Sowmya Kothapalli

UCID: sk272

Email: sk272@njit.edu

Date: 05/04/2021

Programming Language: Python

Classification Algorithms Implemented:

- i) Random Forest
- ii) Naïve Bayes
- iii) KNN, K-Nearest Neighbor

Dataset used: Bank Marketing Data Set

Source of Data:

<https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>

Code:

Detailed implementation of the classifiers process

Steps to Implement the process:

1. Preprocess the data

Supervised Data Mining (Classification) ¶

```
In [287]: 1 import pandas as pd
          2 import numpy as np
          3 import matplotlib.pyplot as plt
          4 import seaborn as sns
          5 %matplotlib inline
          6 import warnings
          7 warnings.filterwarnings("ignore")
```

```
In [288]: 1 data = pd.read_csv('bank.csv', header=0)
          2 data = data.dropna()
          3 print(data.shape)
          4 print(list(data.columns))
```

```
(4521, 17)
['age', 'job', 'marital', 'education', 'default', 'balance', 'housing', 'loan', 'contact', 'day', 'month', 'duration', 'campaign', 'pdays', 'previous', 'poutcome', 'y']
```

```
In [289]: 1 dataset.head()
```

```
Out[289]:
```

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
0	30	unemployed	married	primary	no	1787	no	no	cellular	19	oct	79	1	-1	0	unknown	no
1	33	services	married	secondary	no	4789	yes	yes	cellular	11	may	220	1	339	4	failure	no
2	35	management	single	tertiary	no	1350	yes	no	cellular	16	apr	185	1	330	1	failure	no
3	30	management	married	tertiary	no	1476	yes	yes	unknown	3	jun	199	4	-1	0	unknown	no
4	59	blue-collar	married	secondary	no	0	yes	no	unknown	5	may	226	1	-1	0	unknown	no

```
In [292]: 1 #checking for missing values in dataset
          2 data.isnull().sum()
```

```
Out[292]: age      0
          job      0
          marital  0
          education 0
          default  0
          balance  0
          housing  0
          loan     0
          contact  0
          day      0
          month    0
          duration 0
          campaign 0
          pdays    0
          previous 0
          poutcome 0
          y        0
          dtype: int64
```

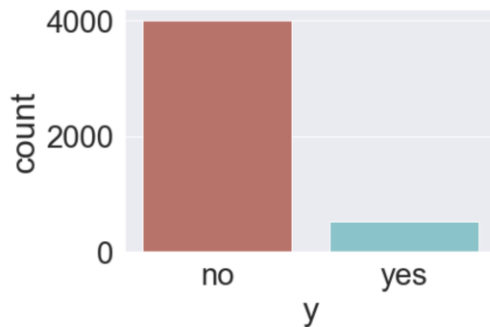
```
In [293]: 1 data['job'].unique()
```

```
Out[293]: array(['unemployed', 'services', 'management', 'blue-collar',
                  'self-employed', 'technician', 'entrepreneur', 'admin.', 'student',
                  'housemaid', 'retired', 'unknown'], dtype=object)
```

```
In [294]: 1 #Frequency of 'subscribed'
          2 data['y'].value_counts()
```

```
Out[294]: no      4000
          yes       521
          Name: y, dtype: int64
```

```
In [295]: 1 sns.countplot(x='y', data=data, palette='hls')
          2 plt.show()
          3 plt.savefig('count_plot')
```



<Figure size 432x288 with 0 Axes>

```
In [296]: 1 #Normalizing the frequency table of 'Subscribed' variable
          2 data['y'].value_counts(normalize=True)
```

```
Out[296]: no      0.88476
          yes      0.11524
          Name: y, dtype: float64
```

From the above analysis we can see that only 3,715 people out of 31,647 have subscribed which is roughly 12%.

```
In [297]: 1 data.groupby('y').mean()
```

```
Out[297]:
```

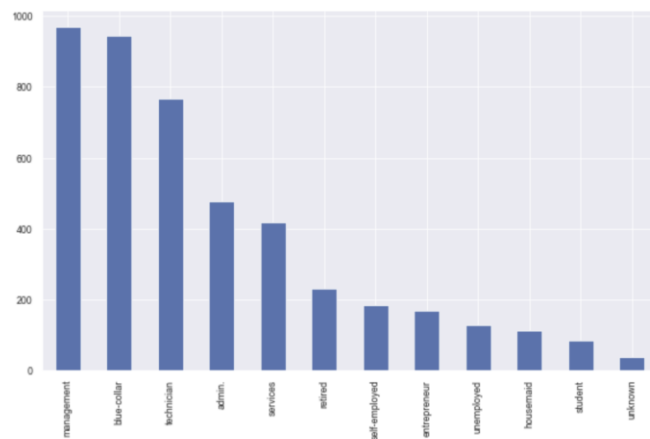
	age	balance	day	duration	campaign	pdays	previous
y							
no	40.998000	1403.211750	15.948750	226.347500	2.862250	36.006000	0.471250
yes	42.491363	1571.955854	15.658349	552.742802	2.266795	68.639155	1.090211

Analyzed the data for all attributes:

```
In [298]: 1 #Analysis 'job'
          2 #Frequency table
          3 data['job'].value_counts()
```

```
Out[298]: management      969
blue-collar      946
technician      768
admin.          478
services        417
retired         230
self-employed   183
entrepreneur    168
unemployed      128
housemaid       112
student         84
unknown         38
Name: job, dtype: int64
```

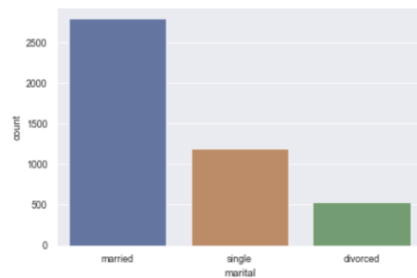
```
In [299]: 1 # Plotting the job frequency table
          2 sns.set_context('paper')
          3 data['job'].value_counts().plot(kind='bar', figsize=(10,6));
```



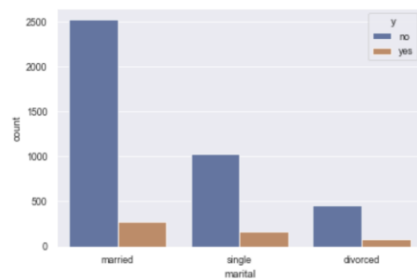
```
In [300]: 1 # Analysing 'marital'
          2 data['marital'].value_counts()
```

```
Out[300]: married    2797
          single     1196
          divorced    528
          Name: marital, dtype: int64
```

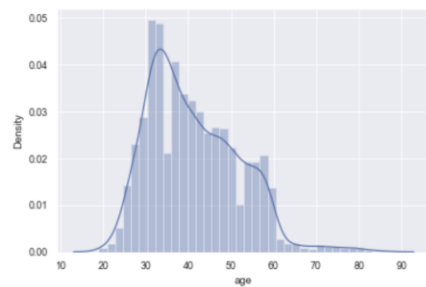
```
In [301]: 1 sns.countplot(data=data, x='marital');
```



```
In [302]: 1 sns.countplot(data=data, x='marital', hue='y');
```



```
In [303]: 1 # Analysing 'age'
          2 sns.distplot(data['age']);
```



```
In [304]: 1 #default vs subscription
          2 pd.crosstab(data['default'], data['y'])
```

```
Out[304]:
```

	y	no	yes
default			
	no	3933	512
	yes	67	9

Correlation matrix:

```
In [307]: 1 # Converting the target variables into 0s and 1s
2 data['y'].replace('no', 0,inplace=True)
3 data['y'].replace('yes', 1,inplace=True)
4 data['y']
```

```
Out[307]: 0      0
1      0
2      0
3      0
4      0
...
4516   0
4517   0
4518   0
4519   0
4520   0
Name: y, Length: 4521, dtype: int64
```

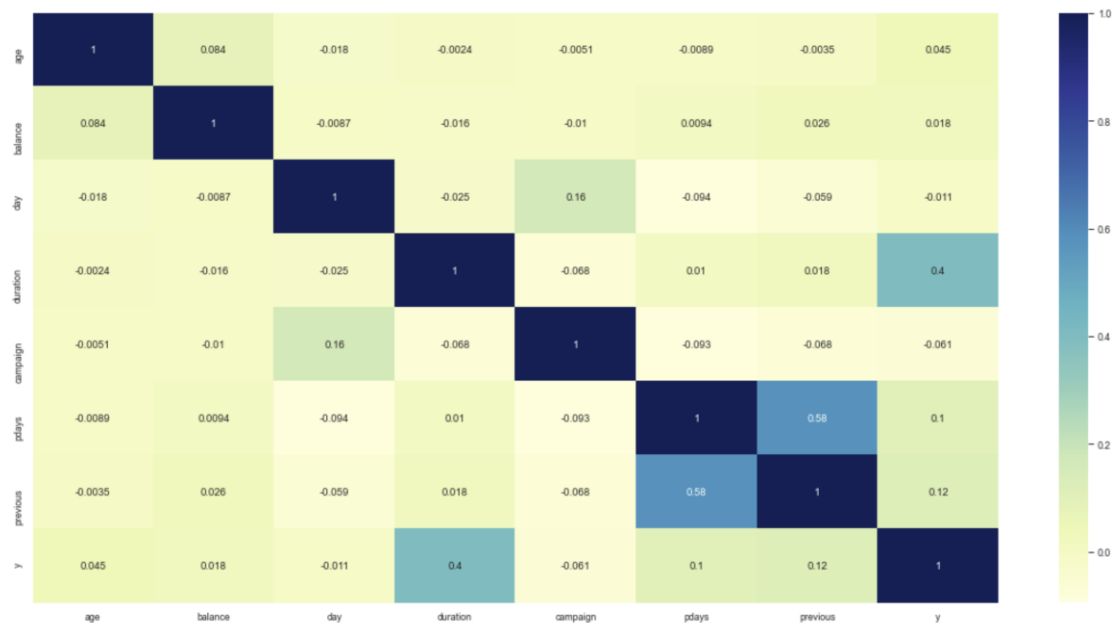
```
In [308]: 1 #Correlation matrix
2 dc = data.corr()
3 dc
```

```
Out[308]:
```

	age	balance	day	duration	campaign	pdays	previous	y
age	1.000000	0.083820	-0.017853	-0.002367	-0.005148	-0.008894	-0.003511	0.045092
balance	0.083820	1.000000	-0.008677	-0.015950	-0.009976	0.009437	0.026196	0.017905
day	-0.017853	-0.008677	1.000000	-0.024629	0.160706	-0.094352	-0.059114	-0.011244
duration	-0.002367	-0.015950	-0.024629	1.000000	-0.068382	0.010380	0.018080	0.401118
campaign	-0.005148	-0.009976	0.160706	-0.068382	1.000000	-0.093137	-0.067833	-0.061147
pdays	-0.008894	0.009437	-0.094352	0.010380	-0.093137	1.000000	0.577562	0.104087
previous	-0.003511	0.026196	-0.059114	0.018080	-0.067833	0.577562	1.000000	0.116714
y	0.045092	0.017905	-0.011244	0.401118	-0.061147	0.104087	0.116714	1.000000

```
In [309]: 1 fig,ax= plt.subplots()
2 fig.set_size_inches(20,10)
3 sns.heatmap(dc, annot=True, cmap='YlGnBu')
```

```
Out[309]: <AxesSubplot:>
```



2. Split the dataset with the help of train_test_split function

```
In [310]: 1 target = data['y']
          2 data = data.drop('y', axis=1)

In [311]: 1 #generating dummy values on the dataset
          2 data = pd.get_dummies(data)
          3 data.head()

Out[311]:
```

	age	balance	day	duration	campaign	pdays	previous	job_admin.	job_blue-collar	job_entrepreneur	...	month_jun	month_mar	month_may	month_nov	mor
0	30	1787	19	79	1	-1	0	0	0	0	...	0	0	0	0	0
1	33	4789	11	220	1	339	4	0	0	0	...	0	0	1	0	0
2	35	1350	16	185	1	330	1	0	0	0	...	0	0	0	0	0
3	30	1476	3	199	4	-1	0	0	0	0	...	1	0	0	0	0
4	59	0	5	226	1	-1	0	0	1	0	...	0	0	1	0	0

5 rows × 51 columns

```
In [312]: 1 from sklearn.model_selection import train_test_split
          2 X_train, X_val, y_train, y_val = train_test_split(data, target, test_size=0.2, random_state=12)
```

Random Forest Classifier:

Import classifier using sklearn and used confusion matrix and classification report

Random Forest Classifier

```
In [313]: 1 from sklearn.ensemble import RandomForestClassifier

In [314]: 1 #creating an object of random forest classifier
          2 rfclf = RandomForestClassifier()
          3 rfclf.fit(X_train,y_train)

Out[314]: RandomForestClassifier()

In [315]: 1 #Making predictions and Accuracy
          2 y_pred = rfclf.predict(X_val)
          3
          4 print("Accuracy of Random forest classifier on dataset:",metrics.accuracy_score(y_val, y_pred))

Accuracy of Random forest classifier on dataset: 0.8983425414364641

In [316]: 1 from sklearn.metrics import classification_report

In [317]: 1 print(classification_report(y_val, y_pred))
```

	precision	recall	f1-score	support
0	0.91	0.99	0.94	792
1	0.76	0.27	0.40	113
accuracy			0.90	905
macro avg	0.83	0.63	0.67	905
weighted avg	0.89	0.90	0.88	905

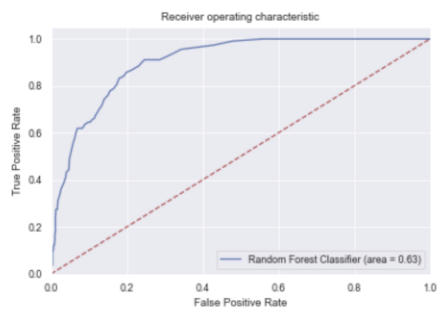

```
In [318]: 1 from sklearn.metrics import confusion_matrix
2 confusion_matrix = confusion_matrix(y_val, y_pred)
3 print(confusion_matrix)

[[782  10]
 [ 82  31]]
```

The result is telling us that we have 775+17 correct predictions and 85+28 incorrect predictions.

Roc for Random Forest classifier

```
In [320]: 1 rf_roc_auc = roc_auc_score(y_val, rfclf.predict(X_val))
2 fpr, tpr, thresholds = roc_curve(y_val, rfclf.predict_proba(X_val)[:,-1])
3 plt.figure()
4 plt.plot(fpr, tpr, label='Random Forest Classifier (area = %0.2f)' % rf_roc_auc)
5 plt.plot([0, 1], [0, 1], 'r--')
6 plt.xlim([0.0, 1.0])
7 plt.ylim([0.0, 1.05])
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver operating characteristic')
11 plt.legend(loc='lower right')
12 plt.savefig('RF_ROC')
13 plt.show()
```



Naive Bayes Classifier:

Import classifier using sklearn and used confusion matrix and classification report

Naive Bayes Classifier

```
In [321]: 1 #Import Gaussian Naive Bayes model
          2 from sklearn.naive_bayes import GaussianNB

In [322]: 1 #creating an object of Naive Bayes classifier
          2 nbclf = GaussianNB()
          3 nbclf.fit(X_train,y_train)

Out[322]: GaussianNB()

In [323]: 1 #Predict the response for dataset
          2 y_pred = nbclf.predict(X_val)
          3 print("Accuracy of Naive Bayes classifier on dataset:",metrics.accuracy_score(y_val, y_pred))

Accuracy of Naive Bayes classifier on dataset: 0.8232044198895028

In [324]: 1 from sklearn.metrics import classification_report
          2 print(classification_report(y_val, y_pred))
```

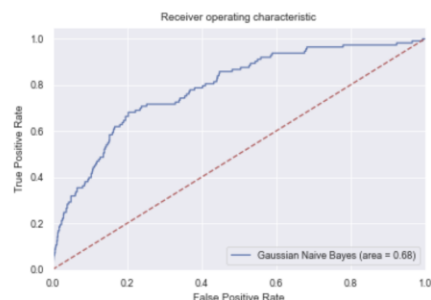
	precision	recall	f1-score	support
0	0.92	0.87	0.90	792
1	0.35	0.48	0.40	113
accuracy			0.82	905
macro avg	0.63	0.68	0.65	905
weighted avg	0.85	0.82	0.83	905

Roc for Naive Bayes Classifier

```
In [325]: 1 from sklearn.metrics import confusion_matrix
          2 confusion_matrix = confusion_matrix(y_val, y_pred)
          3 print(confusion_matrix)

[[691 101]
 [ 59  54]]

In [326]: 1 rf_roc_auc = roc_auc_score(y_val, nbclf.predict(X_val))
          2 fpr, tpr, thresholds = roc_curve(y_val, nbclf.predict_proba(X_val)[:,-1])
          3 plt.figure()
          4 plt.plot(fpr, tpr, label='Gaussian Naive Bayes (area = %0.2f)' % rf_roc_auc)
          5 plt.plot([0, 1], [0, 1], 'r--')
          6 plt.xlim([0.0, 1.0])
          7 plt.ylim([0.0, 1.05])
          8 plt.xlabel('False Positive Rate')
          9 plt.ylabel('True Positive Rate')
         10 plt.title('Receiver operating characteristic')
         11 plt.legend(loc='lower right')
         12 plt.savefig('NB_ROC')
         13 plt.show()
```



KNN, K-Nearest Neighbor Classifier:

KNN, K-Nearest Neighbor Classifier

```
In [327]: 1 from sklearn.neighbors import KNeighborsClassifier
```

```
In [328]: 1 knclf = KNeighborsClassifier(n_neighbors=3)
          2 knclf.fit(X_train,y_train)
```

```
Out[328]: KNeighborsClassifier(n_neighbors=3)
```

```
In [329]: 1 #Predict the response for dataset
          2 y_pred = knclf.predict(X_val)
          3 print("Accuracy of K-Nearest Neighbor classifier on dataset:",metrics.accuracy_score(y_val, y_pred))
```

Accuracy of K-Nearest Neighbor classifier on dataset: 0.8574585635359117

```
In [330]: 1 from sklearn.metrics import classification_report
          2 print(classification_report(y_val, y_pred))
```

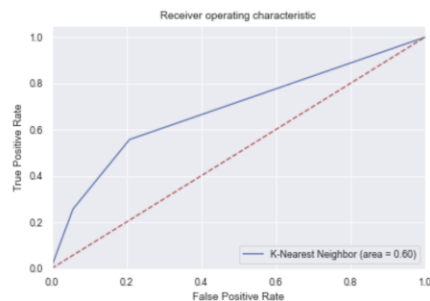
	precision	recall	f1-score	support
0	0.90	0.94	0.92	792
1	0.39	0.26	0.31	113
accuracy			0.86	905
macro avg	0.65	0.60	0.62	905
weighted avg	0.84	0.86	0.84	905

```
In [331]: 1 from sklearn.metrics import confusion_matrix
          2 confusion_matrix = confusion_matrix(y_val, y_pred)
          3 print(confusion_matrix)
```

```
[[747 45]
 [ 84 29]]
```

Roc for KNN Classifier

```
In [332]: 1 rf_roc_auc = roc_auc_score(y_val, knclf.predict(X_val))
          2 fpr, tpr, thresholds = roc_curve(y_val, knclf.predict_proba(X_val)[:,:1])
          3 plt.figure()
          4 plt.plot(fpr, tpr, label='K-Nearest Neighbor (area = %0.2f)' % rf_roc_auc)
          5 plt.plot([0, 1], [0, 1], 'r--')
          6 plt.xlim([0.0, 1.0])
          7 plt.ylim([0.0, 1.05])
          8 plt.xlabel('False Positive Rate')
          9 plt.ylabel('True Positive Rate')
         10 plt.title('Receiver operating characteristic')
         11 plt.legend(loc="lower right")
         12 plt.savefig('KNN_ROC')
         13 plt.show()
```



comparison of classification accuracies between the classification algorithms:

```
In [333]: 1 from sklearn.model_selection import train_test_split
2 from sklearn import model_selection
3 from sklearn.utils import class_weight
4 from sklearn.svm import SVC
5 X_train, X_val, y_train, y_val = train_test_split(data, target, test_size=0.25, random_state=8675309)

In [334]: 1 def run_exps(X_train: pd.DataFrame, y_train: pd.DataFrame, X_test: pd.DataFrame, y_test: pd.DataFrame) -> pd.DataFrame:
2
3     dfs = []
4     models = [
5         ('RF', RandomForestClassifier()),
6         ('KNN', KNeighborsClassifier()),
7         ('GNB', GaussianNB())
8     ]
9     results = []
10    names = []
11    scoring = ['accuracy', 'precision_weighted', 'recall_weighted', 'f1_weighted', 'roc_auc']
12    target_names = ['malignant', 'benign']
13    for name, model in models:
14        kfold = model_selection.KFold(n_splits=5, shuffle=True, random_state=90210)
15        cv_results = model_selection.cross_validate(model, X_train, y_train, cv=kfold, scoring=scoring)
16        clf = model.fit(X_train, y_train)
17        y_pred = clf.predict(X_val)
18        print(name)
19        print(classification_report(y_val, y_pred, target_names=target_names))
20    results.append(cv_results)
21    names.append(name)
22    this_df = pd.DataFrame(cv_results)
23    this_df['model'] = name
24    dfs.append(this_df)
25    final = pd.concat(dfs, ignore_index=True)
26
27    #return final
28
```

RF				
	precision	recall	f1-score	support
malignant	0.91	0.99	0.95	1001
benign	0.73	0.25	0.38	130
accuracy			0.90	1131
macro avg	0.82	0.62	0.66	1131
weighted avg	0.89	0.90	0.88	1131
KNN				
	precision	recall	f1-score	support
malignant	0.90	0.97	0.93	1001
benign	0.42	0.15	0.22	130
accuracy			0.88	1131
macro avg	0.66	0.56	0.58	1131
weighted avg	0.84	0.88	0.85	1131
GNB				
	precision	recall	f1-score	support
malignant	0.93	0.86	0.90	1001
benign	0.33	0.52	0.40	130
accuracy			0.82	1131
macro avg	0.63	0.69	0.65	1131
weighted avg	0.86	0.82	0.84	1131

```

In [335]: 1 seed = 10
          2 models = []
          3 models.append(('RF', RandomForestClassifier()))
          4 models.append(('NB', GaussianNB()))
          5 models.append(('KNN', KNeighborsClassifier()))
          6
          7
          8 results = []
          9 names = []
         10 scoring = 'accuracy'
         11 for name, model in models:
         12     kfold = model_selection.KFold(n_splits=10, random_state=seed)
         13     cv_results = model_selection.cross_val_score(model, X_train, y_train, cv=kfold, scoring=scoring)
         14     results.append(cv_results)
         15     names.append(name)
         16     msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())
         17     print(msg)

RF: 0.895280 (0.011818)
NB: 0.828614 (0.013222)
KNN: 0.875221 (0.011275)

```

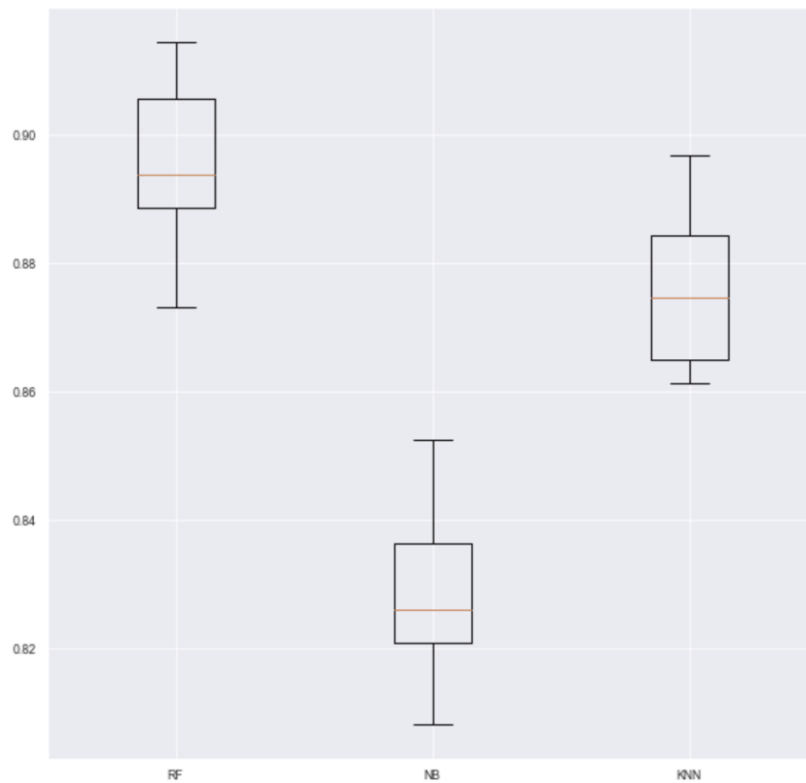
```

In [336]: 1 # boxplot algorithm comparison
          2 fig = plt.figure(figsize=(10,10))
          3 fig.suptitle('comparison of classification accuracies')
          4 ax = fig.add_subplot(111)
          5 plt.boxplot(results)
          6 ax.set_xticklabels(names)
          7 plt.show()

```

comparison of classification accuracies

comparison of classification accuracies



comparison of classification accuracies between the classification algorithms: Random Forest, KNN and Naïve Bayes:

RF				
	precision	recall	f1-score	support
malignant	0.91	0.99	0.95	1001
benign	0.73	0.25	0.38	130
accuracy			0.90	1131
macro avg	0.82	0.62	0.66	1131
weighted avg	0.89	0.90	0.88	1131

KNN				
	precision	recall	f1-score	support
malignant	0.90	0.97	0.93	1001
benign	0.42	0.15	0.22	130
accuracy			0.88	1131
macro avg	0.66	0.56	0.58	1131
weighted avg	0.84	0.88	0.85	1131

GNB				
	precision	recall	f1-score	support
malignant	0.93	0.86	0.90	1001
benign	0.33	0.52	0.40	130
accuracy			0.82	1131
macro avg	0.63	0.69	0.65	1131
weighted avg	0.86	0.82	0.84	1131