

RandomForest_Part01

March 14, 2022

```
[1]: import pandas as pd
import seaborn as sns; sns.set()
import matplotlib.pyplot as plt
import numpy as np

pd.options.display.max_columns = None
pd.set_option('display.max_rows', 200)
pd.set_option('display.float_format', lambda x: '%.3f' % x)

#import dataset
df = pd.read_csv ('WinnipegDataset.txt', sep = ",")
```

0.1 Exploratory Data Analysis

The display of the attributes and their data types was not run because when the output pdf is created the 175 columns when plotted in the pdf take up a lot of redundant space.

```
[ ]: # determine data types in dataset
df.dtypes
```

```
[ ]: #display attributes
display(df)
```

```
[ ]: # identifies columns having null values as a list
df.columns[df.isna().any()].tolist()
```

```
[ ]: # plot histogram of label column
plt.hist(df['label'], bins=np.arange(df['label'].min(), df['label'].max()+2)-0.
↪5)
```

The histograms, boxplots and descriptive stats of the attributes were not run because when the output pdf is created the 175 columns when plotted in the pdf take up a lot of redundant space.

```
[ ]: # histograms of dataset attributes
for column in df:
    plt.figure()
    df.hist([column])
```

```
[ ]: # boxplots of dataset attributes
for column in df:
    plt.figure()
    df.boxplot([column])
```

```
[ ]: # descriptive statistics of attributes
print(df.describe(include='all'))
```

It was decided not to drop any features or remove any outliers. Outliers were not removed because in remote sensing each value of attribute corresponds to a pixel in the output remote sensing image. In essence by removing that value you would not be able to classify that pixel.

0.2 Split Dataset

```
[2]: # prepare data for train/test split
X = df.drop("label", 1)
y = df['label']
```

```
[3]: # Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    ↪random_state=0)
```

```
[4]: # standardize the X dataset
from sklearn.preprocessing import StandardScaler

sc = StandardScaler()
X_train_scaled = sc.fit_transform(X_train)
X_test_scaled = sc.transform(X_test)
```

0.3 Reduce Dimensions Using Principal Component Analysis

```
[13]: import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.decomposition import PCA

pca_test = PCA()
pca_test.fit(X_train_scaled)

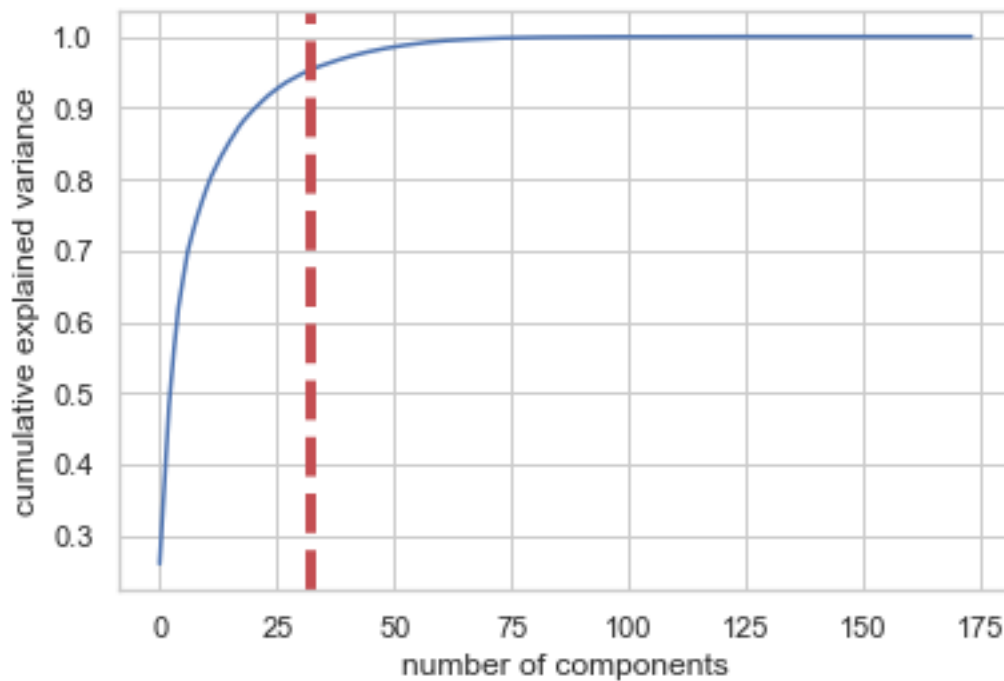
sns.set(style='whitegrid')
plt.plot(np.cumsum(pca_test.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance')
plt.axvline(linewidth=4, color='r', linestyle = '--', x=32, ymin=0, ymax=1)
display(plt.show())
```

```

evr = pca_test.explained_variance_ratio_
cvr = np.cumsum(pca_test.explained_variance_ratio_)

pca_df = pd.DataFrame()
pca_df['Cumulative Variance Ratio'] = cvr
pca_df['Explained Variance Ratio'] = evr
display(pca_df.head(40))

```



None

	Cumulative Variance Ratio	Explained Variance Ratio
0	0.260	0.260
1	0.374	0.114
2	0.481	0.107
3	0.557	0.076
4	0.619	0.062
5	0.660	0.041
6	0.699	0.039
7	0.724	0.025
8	0.747	0.023
9	0.768	0.021
10	0.787	0.019
11	0.804	0.017
12	0.818	0.014
13	0.831	0.013

14	0.842	0.012
15	0.854	0.011
16	0.864	0.010
17	0.874	0.010
18	0.883	0.009
19	0.890	0.007
20	0.897	0.007
21	0.904	0.007
22	0.910	0.006
23	0.916	0.006
24	0.922	0.005
25	0.927	0.005
26	0.931	0.005
27	0.936	0.004
28	0.940	0.004
29	0.943	0.004
30	0.947	0.003
31	0.950	0.003
32	0.953	0.003
33	0.956	0.003
34	0.958	0.003
35	0.960	0.002
36	0.963	0.002
37	0.965	0.002
38	0.967	0.002
39	0.969	0.002

```
[14]: pca = PCA(n_components=32)
      pca.fit(X_train_scaled)

      X_train_scaled_pca = pca.transform(X_train_scaled)
      X_test_scaled_pca = pca.transform(X_test_scaled)
```

0.4 Fine Tune Hyperparameters of Random Forest Classifier

```
[18]: # next sections uses random grid and search to fine tune hyperparameters of
      ↳ RandomForest classifier

      n_estimators = [5,20,50,100] # number of trees in the random forest
      max_features = ['auto', 'sqrt'] # number of features in consideration at every
      ↳ split
      max_depth = [int(x) for x in np.linspace(10, 120, num = 12)] # maximum number
      ↳ of levels allowed in each decision tree
      min_samples_split = [2, 6, 10] # minimum sample number to split a node
      min_samples_leaf = [1, 3, 4] # minimum sample number that can be stored in a
      ↳ leaf node
      bootstrap = [True, False] # method used to sample data points
```

```

random_grid = {'n_estimators': n_estimators,

               'max_features': max_features,

               'max_depth': max_depth,

               'min_samples_split': min_samples_split,

               'min_samples_leaf': min_samples_leaf,

               'bootstrap': bootstrap}

```

```

[19]: # create RandomForest classifier
from sklearn.ensemble import RandomForestClassifier

classifier = RandomForestClassifier()

```

```

[20]: # create RandomizedSearch object with parameters
from sklearn.model_selection import RandomizedSearchCV
classifier_random = RandomizedSearchCV(estimator =_
    ↪ classifier, param_distributions = random_grid,
    n_iter = 100, cv = 5, verbose=2, random_state=35, n_jobs = -1)

```

```

[21]: # run the RandomizedSearch classifier on training data to determine best_
    ↪ parameters for RandomForest classifier run
classifier_random.fit(X_train_scaled_pca, y_train)

```

Fitting 5 folds for each of 100 candidates, totalling 500 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 33 tasks      | elapsed: 4.6min
[Parallel(n_jobs=-1)]: Done 154 tasks     | elapsed: 95.3min
[Parallel(n_jobs=-1)]: Done 357 tasks     | elapsed: 191.4min
[Parallel(n_jobs=-1)]: Done 500 out of 500 | elapsed: 267.0min finished

```

```

[21]: RandomizedSearchCV(cv=5, estimator=RandomForestClassifier(), n_iter=100,
    n_jobs=-1,
    param_distributions={'bootstrap': [True, False],
        'max_depth': [10, 20, 30, 40, 50, 60,
            70, 80, 90, 100, 110,
            120],
        'max_features': ['auto', 'sqrt'],
        'min_samples_leaf': [1, 3, 4],
        'min_samples_split': [2, 6, 10],
        'n_estimators': [5, 20, 50, 100]},
    random_state=35, verbose=2)

```

```
[22]: # print random grid parameters
print ('Random grid: ', random_grid, '\n')
# print the best parameters
print ('Best Parameters: ', classifier_random.best_params_, ' \n')
```

```
Random grid: {'n_estimators': [5, 20, 50, 100], 'max_features': ['auto',
'sqrt'], 'max_depth': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, 120],
'min_samples_split': [2, 6, 10], 'min_samples_leaf': [1, 3, 4], 'bootstrap':
[True, False]}
```

```
Best Parameters: {'n_estimators': 50, 'min_samples_split': 6,
'min_samples_leaf': 1, 'max_features': 'auto', 'max_depth': 40, 'bootstrap':
False}
```

0.5 Random Forest Classifier Model

```
[23]: # run the RandomForest classifier using fine tuned settings
from sklearn.ensemble import RandomForestClassifier

classifier_finetune = RandomForestClassifier(n_estimators=50,
→min_samples_split=6, min_samples_leaf=1, max_features='auto', max_depth=40,
→bootstrap=False, random_state=0)
classifier_finetune.fit(X_train_scaled_pca, y_train)

# Predicting the Test set results
y_pred = classifier_finetune.predict(X_test_scaled_pca)
```

```
[24]: # create confusion matrix and calculate accuracy and cohen's kappa
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.metrics import cohen_kappa_score

cm = confusion_matrix(y_test, y_pred)
print(cm, ' \n')
print(classification_report(y_test, y_pred))
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Cohen's Kappa:", cohen_kappa_score(y_test, y_pred))
```

```
[[11806    0     7    29     0    10     4]
 [   1 1068     0     1     0     0     0]
 [   3     1 22799    20     2     9     0]
 [  20     1   14 22244    22    35     0]
 [  13     0   16   51 13671   217     5]
 [  12     0   13   30   74 25206     2]
 [   4     0     1     0     0     3  337]]
```

	precision	recall	f1-score	support
1	1.00	1.00	1.00	11856
2	1.00	1.00	1.00	1070
3	1.00	1.00	1.00	22834
4	0.99	1.00	1.00	22336
5	0.99	0.98	0.99	13973
6	0.99	0.99	0.99	25337
7	0.97	0.98	0.97	345
accuracy			0.99	97751
macro avg	0.99	0.99	0.99	97751
weighted avg	0.99	0.99	0.99	97751

Accuracy: 0.9936573538889628

Cohen's Kappa: 0.9919768959631383

[]: