



Final Artificial intelligent

NAME: MD. NAYEEM MOLLA

ID:20183290381

MAJOR: COMPUTER SCIENCE &TECHNOLOGY

SCHOOL: School of Information Science & Engineering

Subject: Artificial intelligent

Date: 6/5/2022

Final Project

Introduction

This is my final project for the Artificial intelligent course. I did complete all homework and attend each and every class it was really a wonderful learning course for me because I want to do my masters also in Artificial intelligence. In this project, I will do classification for the iris dataset from chapter 1 then I will use SVM and random forest respectively, and last, I will use matrices to check the accuracy to evaluate my model.

Use the Iris dataset in Chapter 1 (from sklearn. datasets import load_iris) to do the classification

The Iris Dataset contains four features (length and width of sepals and petals) of 50 samples of three species of Iris (Iris setosa, Iris virginica, and Iris versicolor). These measures were used to create a linear discriminant model to classify the species. Classification is the process of categorizing a set of data into classes. It can be done on both structured and unstructured data. Predicting the class of provided data points is the first step in the procedure. The classes are also known as the target, label, or categories. Approximating the mapping function from discrete input variables to discrete output variables is the problem of classification predictive modeling. The basic purpose is to determine which class/category the new data belongs to. Now I will classify iris dataset:

Name: MD NAYEEM MOLLA

ID;20183290381

```
[1] 1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from pandas.plotting import parallel_coordinates
7 from sklearn.tree import DecisionTreeClassifier, plot_tree
8 from sklearn import metrics
9 from sklearn.naive_bayes import GaussianNB
10 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.svm import SVC
13 from sklearn.linear_model import LogisticRegression
```

```
[2] 1 # load through url
2 url = 'http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data'
3 attributes = ["sepal_length", "sepal_width", "petal_length", "petal_width", "class"]
4 dataset = pd.read_csv(url, names = attributes)
5 dataset.columns = attributes
```

```
[2] 1 # load through url
2 url = 'http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data'
3 attributes = ["sepal_length", "sepal_width", "petal_length", "petal_width", "class"]
4 dataset = pd.read_csv(url, names = attributes)
5 dataset.columns = attributes
```

1 dataset

| | | | | | |
|-----|-----|-----|-----|-----|----------------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

1 data.head(5)



| | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

[6] 1 data.dtypes

sepal_length float64
sepal_width float64
petal_length float64
petal_width float64
species object
dtype: object

dtype: object

[7] 1 # number of instances in each class
2 data.groupby('species').size()

species
setosa 50
versicolor 50
virginica 50
dtype: int64

[8] 1 # Take out a test set
2 train, test = train_test_split(data, test_size = 0.4, stratify = data['species'], random_state = 42)

[9] 1 # number of instances in each class in training data
2 train.groupby('species').size()

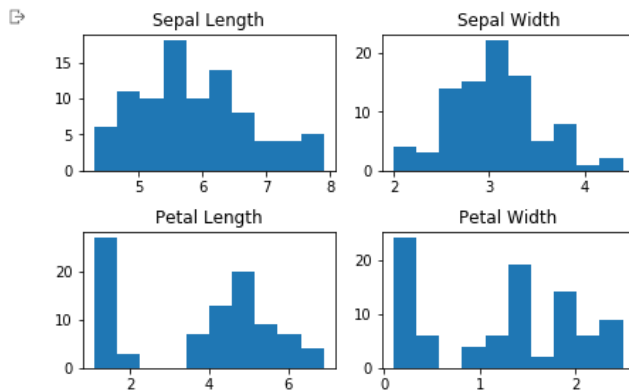
species
setosa 30
versicolor 30
virginica 30
dtype: int64

nt

```

1  # histograms
2  n_bins = 10
3  fig, axs = plt.subplots(2, 2)
4  axs[0,0].hist(train['sepal_length'], bins = n_bins);
5  axs[0,0].set_title('Sepal Length');
6  axs[0,1].hist(train['sepal_width'], bins = n_bins);
7  axs[0,1].set_title('Sepal Width');
8  axs[1,0].hist(train['petal_length'], bins = n_bins);
9  axs[1,0].set_title('Petal Length');
10 axs[1,1].hist(train['petal_width'], bins = n_bins);
11 axs[1,1].set_title('Petal Width');
12
13 # add some spacing between subplots
14 fig.tight_layout(pad=1.0);

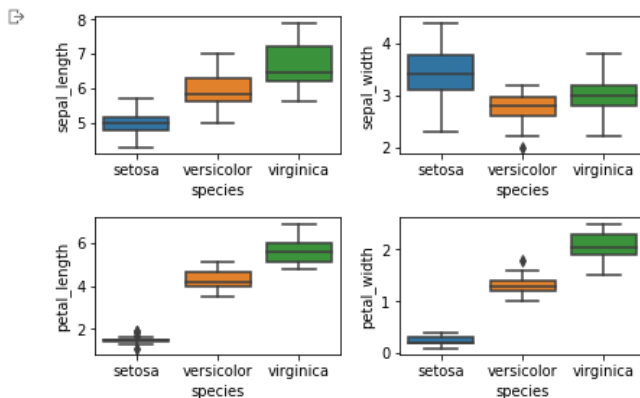
```



```

1  # boxplots using seaborn
2  fig, axs = plt.subplots(2, 2)
3  fn = ["sepal_length", "sepal_width", "petal_length", "petal_width"]
4  cn = ['setosa', 'versicolor', 'virginica']
5  sns.boxplot(x = 'species', y = 'sepal_length', data = train, order = cn, ax = axs[0,0]);
6  sns.boxplot(x = 'species', y = 'sepal_width', data = train, order = cn, ax = axs[0,1]);
7  sns.boxplot(x = 'species', y = 'petal_length', data = train, order = cn, ax = axs[1,0]);
8  sns.boxplot(x = 'species', y = 'petal_width', data = train, order = cn, ax = axs[1,1]);
9  # add some spacing between subplots
10 fig.tight_layout(pad=1.0);

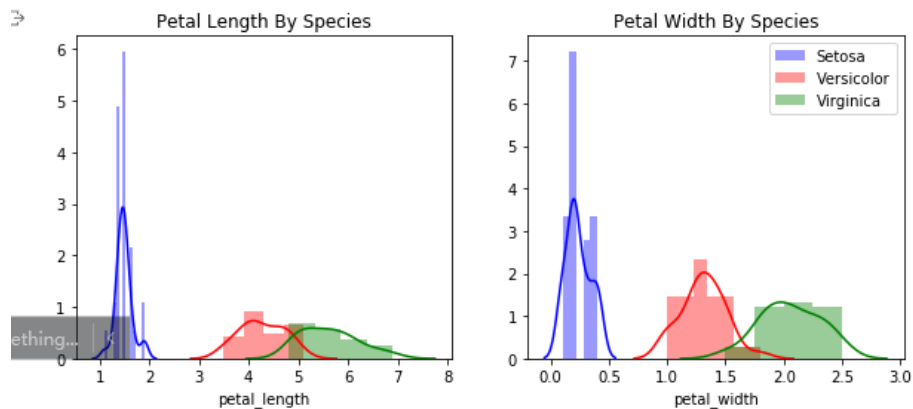
```



```

1
2 # right off the bat, we see that petal length/width can separate setosa from the others
3 # histogram by species
4 setosa_pl = train.loc[train.species=='setosa', 'petal_length']
5 versicolor_pl = train.loc[train.species=='versicolor', 'petal_length']
6 virginica_pl = train.loc[train.species=='virginica', 'petal_length']
7 setosa_pw = train.loc[train.species=='setosa', 'petal_width']
8 versicolor_pw = train.loc[train.species=='versicolor', 'petal_width']
9 virginica_pw = train.loc[train.species=='virginica', 'petal_width']
10
11 fig, axs = plt.subplots(1, 2)
12 # set figure size
13 fig.set_size_inches(10,4)
14 ax1 = sns.distplot(setosa_pl, color="blue", label="Setosa", ax = axs[0]);
15 ax1.set_title('Petal Length By Species')
16 ax1 = sns.distplot(versicolor_pl, color="red", label="Versicolor", ax = axs[0]);
17 ax1 = sns.distplot(virginica_pl, color="green", label="Virginica", ax = axs[0]);
18
19 ax2 = sns.distplot(setosa_pw, color="blue", label="Setosa", ax = axs[1]);
20 ax2.set_title('Petal Width By Species')
21 ax2 = sns.distplot(versicolor_pw, color="red", label="Versicolor", ax = axs[1]);
22 ax2 = sns.distplot(virginica_pw, color="green", label="Virginica", ax = axs[1]);
23
24 plt.legend();

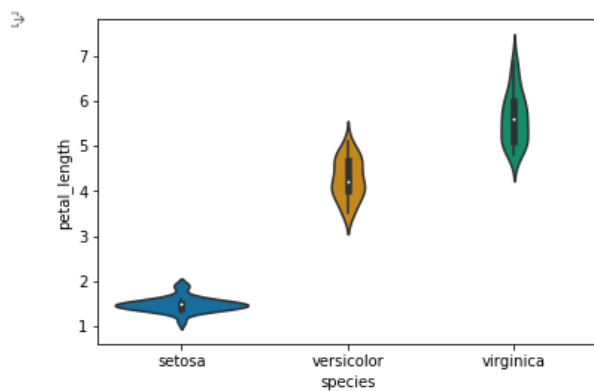
```



```

|3] 1 sns.violinplot(x="species", y="petal_length", data=train, size=5, order = cn, palette = 'colorbli ^

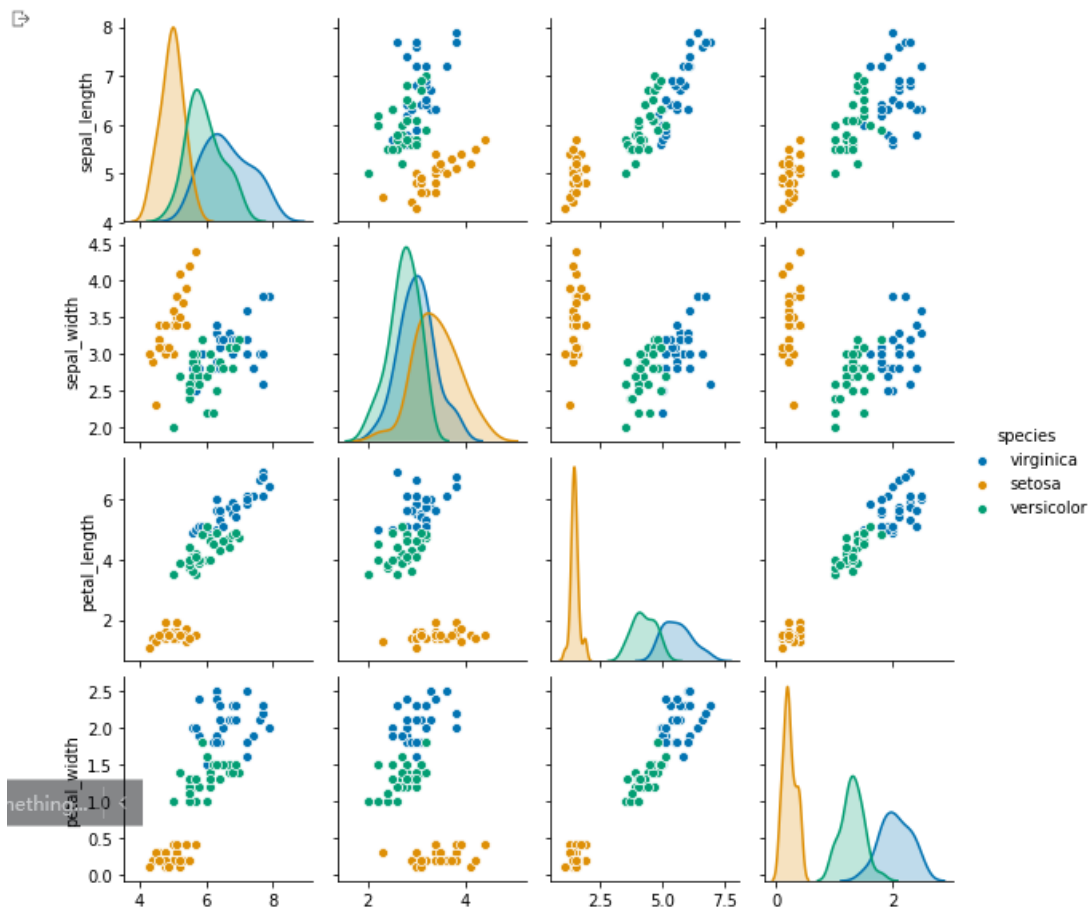
```



```

1 # bivariate relationship
2 # scatterplot matrix
3 sns.pairplot(train, hue="species", height = 2, palette = 'colorblind');

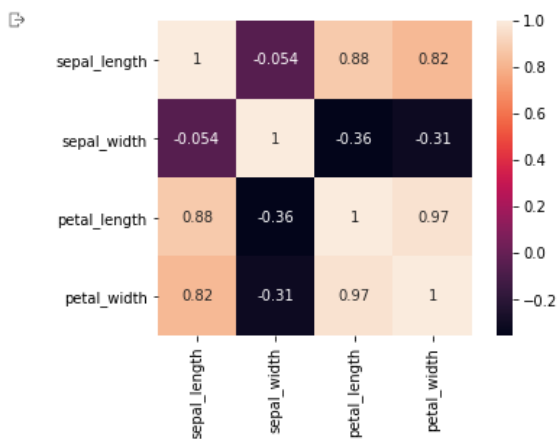
```



```

1 # correlation matrix
2 corrmat = train.corr()
3 sns.heatmap(corrmat, annot = True, square = True);

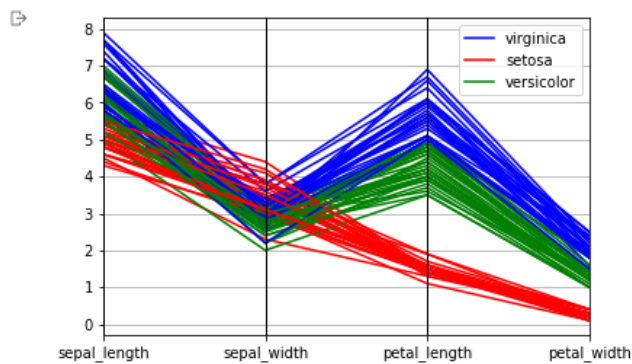
```



```

1 # parallel coordinates
2 parallel_coordinates(train, "species", color = ['blue', 'red', 'green']);

```



```

1 # first try decision tree
2 mod_dt = DecisionTreeClassifier(max_depth = 3, random_state = 1)
3 mod_dt.fit(X_train,y_train)
4 prediction=mod_dt.predict(X_test)
5 print('The accuracy of the Decision Tree is',"{:0.3f}".format(metrics.accuracy_score(prediction,y_

```

The accuracy of the Decision Tree is 0.983

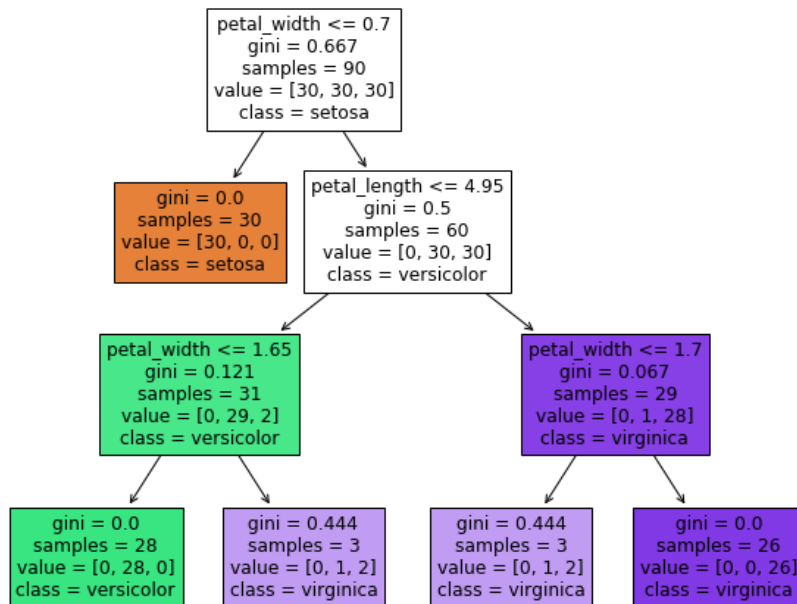
```

[19] 1 mod_dt.feature_importances_

```

array([0. , 0. , 0.42430866, 0.57569134])


```
1 # set figure size
2 plt.figure(figsize = (10,8))
3 plot_tree(mod_dt, feature_names = fn, class_names = cn, filled = True);
```

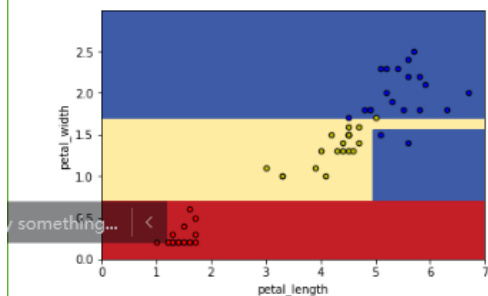


```

1 # plot decision boundary for petal width vs petal length
2 plot_step = 0.01
3 plot_colors = "ryb"
4 xx, yy = np.meshgrid(np.arange(0, 7, plot_step), np.arange(0, 3, plot_step))
5 plt.tight_layout(h_pad=1, w_pad=1, pad=2.5)
6
7 selected_predictors = ["petal_length", "petal_width"]
8 mod_dt_1 = DecisionTreeClassifier(max_depth = 3, random_state = 1)
9 y_train_en = y_train.replace({'setosa':0, 'versicolor':1, 'virginica':2}).copy()
10 mod_dt_1.fit(X_train[selected_predictors], y_train_en)
11
12 pred_all = mod_dt_1.predict(np.c_[xx.ravel(), yy.ravel()])
13 pred_all = pred_all.reshape(xx.shape)
14
15 graph = plt.contourf(xx, yy, pred_all, cmap=plt.cm.RdYlBu)
16
17 plt.xlabel(selected_predictors[0])
18 plt.ylabel(selected_predictors[1])
19
20 # plot test data points
21 n_class = 3
22 for i, color in zip(cn, plot_colors):
23     temp = np.where(y_test == i)
24     idx = [elem for elems in temp for elem in elems]
25     plt.scatter(X_test.iloc[idx, 2], X_test.iloc[idx, 3], c=color,
26               label=y_test, cmap=plt.cm.RdYlBu, edgecolor='black', s=20)
27
28 plt.suptitle("Decision Boundary Shown in 2D with Test Data")
29 plt.axis("tight");

```

Decision Boundary Shown in 2D with Test Data

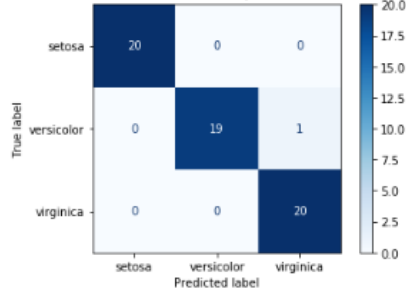


```

1 # confusion matrix
2 # one versicolor misclassified
3 disp = metrics.plot_confusion_matrix(mod_dt, X_test, y_test,
4                                     display_labels=cn,
5                                     cmap=plt.cm.Blues,
6                                     normalize=None)
7 disp.ax_.set_title('Decision Tree Confusion matrix, without normalization');

```

Decision Tree Confusion matrix, without normalization



```
[23] 1 # Gaussian Naive Bayes Classifier
2 mod_gnb_all = GaussianNB()
3 y_pred = mod_gnb_all.fit(X_train, y_train).predict(X_test)
4 print('The accuracy of the Gaussian Naive Bayes Classifier on test data is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

↳ The accuracy of the Gaussian Naive Bayes Classifier on test data is 0.933

```
[24] 1 # Gaussian Naive Bayes Classifier with two predictors
2 mod_gnb = GaussianNB()
3 y_pred = mod_gnb.fit(X_train[selected_predictors], y_train).predict(X_test[selected_predictors])
4 print('The accuracy of the Gaussian Naive Bayes Classifier with 2 predictors on test data is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

↳ The accuracy of the Gaussian Naive Bayes Classifier with 2 predictors on test data is 0.950

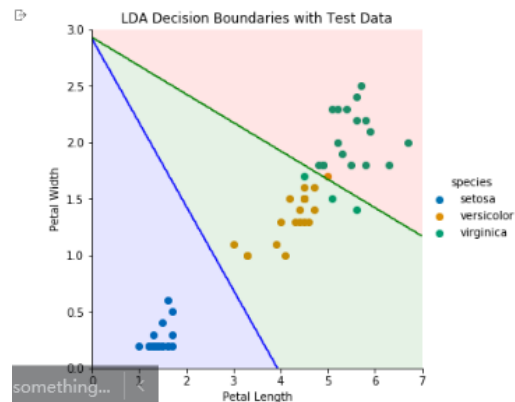
```
[25] 1 # LDA Classifier
2 mod_lda_all = LinearDiscriminantAnalysis()
3 y_pred = mod_lda_all.fit(X_train, y_train).predict(X_test)
4 print('The accuracy of the LDA Classifier on test data is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

↳ The accuracy of the LDA Classifier on test data is 0.983

```
[26] 1 # LDA Classifier with two predictors
2 mod_lda = LinearDiscriminantAnalysis()
3 y_pred = mod_lda.fit(X_train[selected_predictors], y_train).predict(X_test[selected_predictors])
4 print('The accuracy of the LDA Classifier with two predictors on test data is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

↳ The accuracy of the LDA Classifier with two predictors on test data is 0.933

```
27] 1 # LDA with 2 predictors
2 mod_lda_1 = LinearDiscriminantAnalysis()
3 y_pred = mod_lda_1.fit(X_train[selected_predictors], y_train_en).predict(X_test[selected_predictors])
4
5 N = 300
6 X = np.linspace(0, 7, N)
7 Y = np.linspace(0, 3, N)
8 X, Y = np.meshgrid(X, Y)
9
10 g = sns.FacetGrid(test, hue="species", height=5, palette = 'colorblind').map(plt.scatter, "petal_length", "petal_width", ).add_legend
11 my_ax = g.ax
12
13 zz = np.array([mod_lda_1.predict(np.array([[xx, yy]])) for xx, yy in zip(np.ravel(X), np.ravel(Y)) ])
14 Z = zz.reshape(X.shape)
15
16 #Plot the filled and boundary contours
17 my_ax.contourf( X, Y, Z, 2, alpha = .1, colors = ('blue', 'green', 'red'))
18 my_ax.contour( X, Y, Z, 2, alpha = 1, colors = ('blue', 'green', 'red'))
19
20 # Add axis and title
21 my_ax.set_xlabel('Petal Length')
22 my_ax.set_ylabel('Petal Width')
23 my_ax.set_title('LDA Decision Boundaries with Test Data');
```



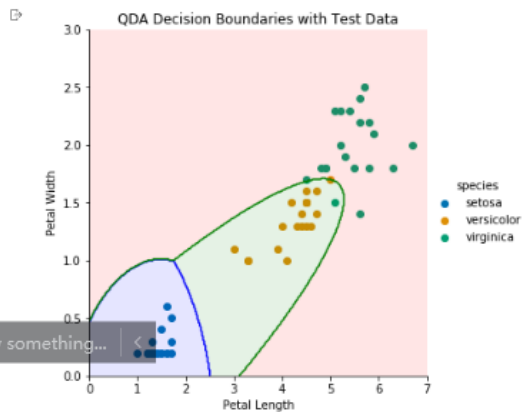
```
[28] 1 # QDA Classifier
2 mod_qda_all = QuadraticDiscriminantAnalysis()
3 y_pred = mod_qda_all.fit(X_train, y_train).predict(X_test)
4 print('The accuracy of the QDA Classifier is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

The accuracy of the QDA Classifier is 0.983

```
[29] 1 # QDA Classifier with two predictors
2 mod_qda = QuadraticDiscriminantAnalysis()
3 y_pred = mod_qda.fit(X_train[selected_predictors], y_train).predict(X_test[selected_predictors])
4 print('The accuracy of the QDA Classifier with two predictors is', "{:.3f}".format(metrics.accuracy_score(y_pred, y_test)))
```

The accuracy of the QDA Classifier with two predictors is 0.967

```
1 # QDA with 2 predictors
2 mod_qda_1 = QuadraticDiscriminantAnalysis()
3 y_pred = mod_qda_1.fit(X_train.iloc[:, 2:4], y_train_en).predict(X_test.iloc[:, 2:4])
4
5 N = 300
6 X = np.linspace(0, 7, N)
7 Y = np.linspace(0, 3, N)
8 X, Y = np.meshgrid(X, Y)
9
10 g = sns.FacetGrid(test, hue="species", height=5, palette = 'colorblind').map(plt.scatter, "petal_length", "petal_width", ).add_legend
11 my_ax = g.ax
12
13 zz = np.array([mod_qda_1.predict(np.array([[xx, yy]])) for xx, yy in zip(np.ravel(X), np.ravel(Y)) ])
14 Z = zz.reshape(X.shape)
15
16 #Plot the filled and boundary contours
17 my_ax.contourf( X, Y, Z, 2, alpha = .1, colors = ('blue', 'green', 'red'))
18 my_ax.contour( X, Y, Z, 2, alpha = 1, colors = ('blue', 'green', 'red'))
19
20 # Addd axis and title
21 my_ax.set_xlabel('Petal Length')
22 my_ax.set_ylabel('Petal Width')
23 my_ax.set_title('QDA Decision Boundaries with Test Data');
```



```
[31] 1 # KNN, first try 5
2 mod_5nn=KNeighborsClassifier(n_neighbors=5)
3 mod_5nn.fit(X_train,y_train)
4 prediction=mod_5nn.predict(X_test)
5 print('The accuracy of the 5NN Classifier is',"{:.3f}".format(metrics.accuracy_score(prediction,y_test)))
```

↳ The accuracy of the 5NN Classifier is 0.933

```
> 1 # try different k
2 acc_s = pd.Series(dtype = 'float')
3 for i in list(range(1,11)):
4     mod_knn=KNeighborsClassifier(n_neighbors=i)
5     mod_knn.fit(X_train,y_train)
6     prediction=mod_knn.predict(X_test)
7     acc_s = acc_s.append(pd.Series(metrics.accuracy_score(prediction,y_test)))
8
9 plt.plot(list(range(1,11)), acc_s)
10 plt.suptitle("Test Accuracy vs K")
11 plt.xticks(list(range(1,11)))
12 plt.ylim(0.9,0.98);
```

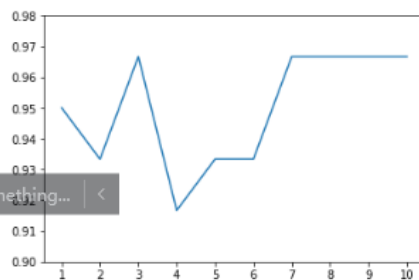
↳ /opt/conda/envs/python35-paddle120-env/lib/python3.7/site-packages/matplotlib/cbook/_init_.py:2064: FutureWarning: Support for multi-dimensional indexing (e.g. 'obj[:, None]') is deprecated and will be removed in a future version. Convert to a numpy array before indexing instead.

x[:, None]

/opt/conda/envs/python35-paddle120-env/lib/python3.7/site-packages/matplotlib/axes/_base.py:250: FutureWarning: Support for multi-dimensional indexing (e.g. 'obj[:, None]') is deprecated and will be removed in a future version. Convert to a numpy array before indexing instead.

y = y[:, np.newaxis]

Test Accuracy vs K



say something... <

```
[33] 1 # SVC with linear kernel
2 # for SVC, may be impractical beyond tens of thousands of samples
3 linear_svc = SVC(kernel='linear').fit(X_train, y_train)
4 prediction=linear_svc.predict(X_test)
5 print('The accuracy of the linear SVC is',"{:.3f}".format(metrics.accuracy_score(prediction,y_test)))
```

↳ The accuracy of the linear SVC is 1.000

```
> 1 # SVC with polynomial kernel
2 poly_svc = SVC(kernel='poly', degree = 4).fit(X_train, y_train)
3 prediction=poly_svc.predict(X_test)
4 print('The accuracy of the Poly SVC is',"{:.3f}".format(metrics.accuracy_score(prediction,y_test)))
```

↳ The accuracy of the Poly SVC is 0.933

```
[35] 1 # Logistic regression
2 mod_lr = LogisticRegression(solver = 'newton-cg').fit(X_train, y_train)
3 prediction=mod_lr.predict(X_test)
4 print('The accuracy of the Logistic Regression is',"{:.3f}".format(metrics.accuracy_score(prediction,y_test)))
```

↳ The accuracy of the Logistic Regression is 0.950

Use SVM and Random Forest, respectively.

SVM: The "Support Vector Machine" (SVM) is a supervised machine learning technique that can be used for classification and regression tasks. It is, however, largely employed in categorization difficulties. We depict each data item as a point in n-dimensional space (where n is the number of features you have) in the SVM method, with the value of each feature being the value of a certain coordinate.

Important Features of Random Forest

1. Diversity- Not all attributes/variables/features are considered while making an individual tree, each tree is different.
2. Immune to the curse of dimensionality- Since each tree does not consider all the features, the feature space is reduced.
3. Parallelization-Each tree is created independently out of different data and attributes. This means that we can make full use of the CPU to build random forests.
4. Train-Test split- In a random forest we don't have to segregate the data for train and test as there will always be 30% of the data which is not seen by the decision tree.
5. Stability- Stability arises because the result is based on majority voting/ averaging.

2、 Use SVM and Random Forest, respectively.

```
[12]: 1 #Import scikit-learn dataset library
      2 from sklearn import datasets
      3
      4 #Load dataset
      5 iris = datasets.load_iris()
```

```
[13]: 1 # print the label species (setosa, versicolor, virginica)
      2 print(iris.target_names)
      3
      4 # print the names of the four features
      5 print(iris.feature_names)
```

```
Out[13]: ['setosa' 'versicolor' 'virginica']
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

[illegible]

```

1 # Creating a DataFrame of given iris dataset.
2 import pandas as pd
3 data=pd.DataFrame({
4     'sepal length':iris.data[:,0],
5     'sepal width':iris.data[:,1],
6     'petal length':iris.data[:,2],
7     'petal width':iris.data[:,3],
8     'species':iris.target
9 })
10 data.head()

```

| | sepal length | sepal width | petal length | petal width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | 0 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | 0 |

```

[16] 1 # Import train_test_split function
2 from sklearn.model_selection import train_test_split
3
4 X=data[['sepal length', 'sepal width', 'petal length', 'petal width']] # Features
5 y=data['species'] # Labels
6
7 # Split dataset into training set and test set
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3) # 70% training and 30% t

```

```

[17] 1 #Import Random Forest Model
2 from sklearn.ensemble import RandomForestClassifier
3
4 #Create a Gaussian Classifier
5 clf=RandomForestClassifier(n_estimators=100)
6
7 #Train the model using the training sets y_pred=clf.predict(X_test)
8 clf.fit(X_train,y_train)
9
10 y_pred=clf.predict(X_test)

```

```

> 1 #Import scikit-learn metrics module for accuracy calculation
2 from sklearn import metrics
3 # Model Accuracy, how often is the classifier correct?
4 print("Accuracy:",metrics.accuracy_score(y_test, y_pred))

```

Accuracy: 0.9555555555555556

```
[19] 1 #Import scikit-learn metrics module for accuracy calculation
    2 from sklearn import metrics
    3 # Model Accuracy, how often is the classifier correct?
    4 print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

➤ Accuracy: 0.9555555555555556

```
▷ 1 from sklearn.ensemble import RandomForestClassifier
    2
    3 #Create a Gaussian Classifier
    4 clf=RandomForestClassifier(n_estimators=100)
    5
    6 #Train the model using the training sets y_pred=clf.predict(X_test)
    7 clf.fit(X_train,y_train)
```

➤ RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None, oob_score=False, random_state=None, verbose=0, warm_start=False)

```
[21] 1 RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
    2 | | | max_depth=None, max_features='auto', max_leaf_nodes=None,
    3 | | | min_impurity_decrease=0.0, min_impurity_split=None,
    4 | | | min_samples_leaf=1, min_samples_split=2,
    5 | | | min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=1,
    6 | | | oob_score=False, random_state=None, verbose=0,
    7 | | | warm_start=False)
```

➤ RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=1, oob_score=False, random_state=None, verbose=0, warm_start=False)


```

1  from sklearn.ensemble import RandomForestClassifier
2
3  #Create a Gaussian Classifier
4  clf=RandomForestClassifier(n_estimators=100)
5
6  #Train the model using the training sets y_pred=clf.predict(X_test)
7  clf.fit(X_train,y_train)
8
9  # prediction on test set
10 y_pred=clf.predict(X_test)
11
12 #Import scikit-learn metrics module for accuracy calculation
13 from sklearn import metrics
14 # Model Accuracy, how often is the classifier correct?
15 print("Accuracy:",metrics.accuracy_score(y_test, y_pred))

```

Accuracy: 0.9555555555555556

```

1  iris = load_iris()
2  df_data = pd.DataFrame(data= np.c_[iris['data'], iris['target']],
3  | | | | | columns= ['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm', 'Spe
4  df_data

```

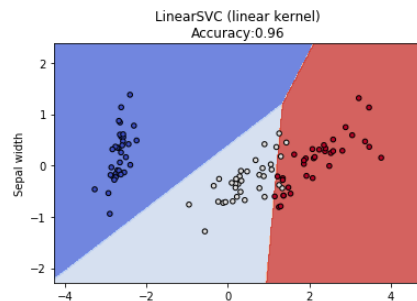
| | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|-----|---------------|--------------|---------------|--------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0.0 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | 0.0 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0.0 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0.0 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | 0.0 |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | 2.0 |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | 2.0 |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | 2.0 |
| 148 | 6.2 | 2.4 | 5.4 | 2.2 | 2.0 |

```
[38] 1 from sklearn.model_selection import train_test_split
2     X = df_data.drop(labels=['Species'],axis=1).values
3     y = df_data['Species'].values
4     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, stratif
5
6     print('train shape:', X_train.shape)
7     print('test shape:', X_test.shape)
```

```
↳ train shape: (105, 4)
   test shape: (45, 4)
```

```
> 1 from sklearn import svm
2
3
4 linearSvcModel=svm.LinearSVC(C=1, max_iter=10000)
5
6 linearSvcModel.fit(train_reduced, y_train)
7
8 predicted=linearSvcModel.predict(train_reduced)
9
10 accuracy = linearSvcModel.score(train_reduced, y_train)
11
12 X0, X1 = train_reduced[:, 0], train_reduced[:, 1]
13 xx, yy = make_meshgrid(X0, X1)
14 plot_contours(plt, linearSvcModel, xx, yy,
15               cmap=plt.cm.coolwarm, alpha=0.8)
16 plt.scatter(X0, X1, c=y_train, cmap=plt.cm.coolwarm, s=20, edgecolors='k')
17 plt.xlabel('Sepal length')
18 plt.ylabel('Sepal width')
19 plt.title('LinearSVC (linear kernel)') + '\n' + 'Accuracy:%.2f'%accuracy)
```

```
↳ Text(0.5,1,'LinearSVC (linear kernel)\nAccuracy:0.96')
```



```

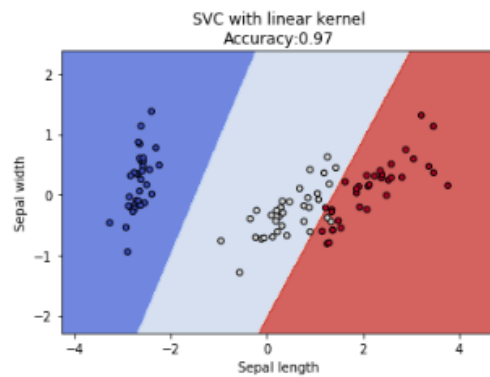
1  from sklearn import svm
2
3
4  svcModel=svm.SVC(kernel='linear', C=1)
5
6  svcModel.fit(train_reduced, y_train)
7
8  predicted=svcModel.predict(train_reduced)
9
10 accuracy = svcModel.score(train_reduced, y_train)
11
12 X0, X1 = train_reduced[:, 0], train_reduced[:, 1]
13 xx, yy = make_meshgrid(X0, X1)
14 plot_contours(plt, svcModel, xx, yy,
15               cmap=plt.cm.coolwarm, alpha=0.8)
16 plt.scatter(X0, X1, c=y_train, cmap=plt.cm.coolwarm, s=20, edgecolors='k')
17 plt.xlabel('Sepal length')
18 plt.ylabel('Sepal width')
19 plt.title('SVC with linear kernel' + '\n' + 'Accuracy: %.2f'%accuracy)

```

```

, Text(0.5,1,'SVC with linear kernel\nAccuracy:0.97')

```

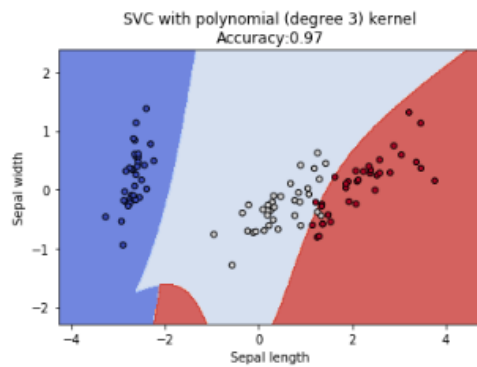


```

1  from sklearn import svm
2
3
4  polyModel=svm.SVC(kernel='poly', degree=3, gamma='auto', C=1)
5
6  polyModel.fit(train_reduced, y_train)
7
8  predicted=polyModel.predict(train_reduced)
9
10 accuracy = polyModel.score(train_reduced, y_train)
11
12 X0, X1 = train_reduced[:, 0], train_reduced[:, 1]
13 xx, yy = make_meshgrid(X0, X1)
14 plot_contours(plt, polyModel, xx, yy,
15               cmap=plt.cm.coolwarm, alpha=0.8)
16 plt.scatter(X0, X1, c=y_train, cmap=plt.cm.coolwarm, s=20, edgecolors='k')
17 plt.xlabel('Sepal length')
18 plt.ylabel('Sepal width')
19 plt.title('SVC with polynomial (degree 3) kernel'+ '\n' + 'Accuracy:0.97')

```

Text(0.5,1,'SVC with polynomial (degree 3) kernel\nAccuracy:0.97')



```

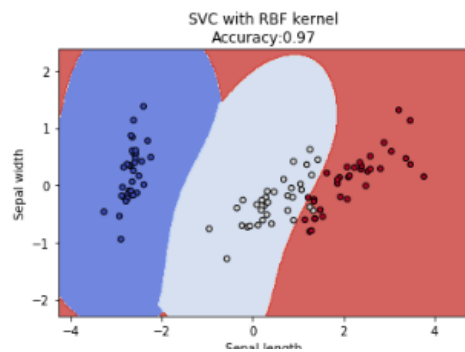
1  from sklearn import svm
2
3
4  rbfModel=svm.SVC(kernel='rbf', gamma=0.7, C=1)
5
6  rbfModel.fit(train_reduced, y_train)
7
8  predicted=rbfModel.predict(train_reduced)
9
10 accuracy = rbfModel.score(train_reduced, y_train)
11
12 X0, X1 = train_reduced[:, 0], train_reduced[:, 1]
13 xx, yy = make_meshgrid(X0, X1)
14 plot_contours(plt, rbfModel, xx, yy,
15               cmap=plt.cm.coolwarm, alpha=0.8)
16 plt.scatter(X0, X1, c=y_train, cmap=plt.cm.coolwarm, s=20, edgecolors='k')
17 plt.xlabel('Sepal length')
18 plt.ylabel('Sepal width')
19 plt.title('SVC with RBF kernel' + '\n' + 'Accuracy:%.2f'%accuracy)

```

```

20 Text(0.5,1,'SVC with RBF kernel\nAccuracy:0.97')

```



Random Forest,

```
[48] 1 from sklearn import datasets
```

```
[49] 1 iris = datasets.load_iris()
```

```
➤ 1 print(iris)
```

```
➤ [6. , 3.4, 4.5, 1.6],  
[6.7, 3.1, 4.7, 1.5],  
[6.3, 2.3, 4.4, 1.3],  
[5.6, 3. , 4.1, 1.3],  
[5.5, 2.5, 4. , 1.3],  
[5.5, 2.6, 4.4, 1.2],  
[6.1, 3. , 4.6, 1.4],  
[5.8, 2.6, 4. , 1.2],  
[5. , 2.3, 3.3, 1. ],  
[5.6, 2.7, 4.2, 1.3],  
[5.7, 3. , 4.2, 1.2],  
[5.7, 2.9, 4.2, 1.3],  
[6.2, 2.9, 4.3, 1.3],  
[5.1, 2.5, 3. , 1.1],  
[5.7, 2.8, 4.1, 1.3],  
[6.3, 3.3, 6. , 2.5],  
[5.8, 2.7, 5.1, 1.9],  
[7.1, 3. , 5.9, 2.1],
```

```
[52] 1 import pandas as pd  
2 set = 'data/data.csv'  
3 iris_csv = pd.read_csv(set)
```

```
[53] 1 iris_csv.head()
```

```
➤
```

| | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

ting

```
> 1 iris_csv.head()
```

↳

| | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

```
[54] 1 iris_csv.tail()
```

↳

| | sepal_length | sepal_width | petal_length | petal_width | species |
|-----|--------------|-------------|--------------|-------------|-----------|
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | virginica |
| | | | | | |

```
1 data = pd.DataFrame({'sepal length':iris.data[:,0],
2                       'sepal width':iris.data[:,1],
3                       'petal length':iris.data[:,2],
4                       'petal width':iris.data[:,3],
5                       'species':iris.target})
6 data.head()
```

↳

| | sepal length | sepal width | petal length | petal width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | 0 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | 0 |

```
[56] 1 from sklearn.model_selection import train_test_split
2     X = data[['sepal length', 'sepal width', 'petal length', 'petal width']] # Features
3     y = data['species'] # Labels
4     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3) # 70% training and 30% test
```

```
[57] 1 X_train.head()
```

↳

| | sepal length | sepal width | petal length | petal width |
|---|--------------|-------------|--------------|-------------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 |


```
[58] 1 from sklearn.ensemble import RandomForestClassifier
2
3 #Create a Gaussian Classifier
4 clf = RandomForestClassifier(n_estimators=100)
```

```
>> 1 clf.fit(X_train, y_train)
```

```
>> RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                           criterion='gini', max_depth=None, max_features='auto',
                           max_leaf_nodes=None, max_samples=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=1, min_samples_split=2,
                           min_weight_fraction_leaf=0.0, n_estimators=100,
                           n_jobs=None, oob_score=False, random_state=None,
                           verbose=0, warm_start=False)
```

```
[61] 1 y_pred = clf.predict(X_test)
```

```
[62] 1 from sklearn import metrics
2 print('Accuracy:', metrics.accuracy_score(y_test, y_pred))
```

```
>> Accuracy: 0.9555555555555556
```

```
[63] 1 # Prediction for a single item
2 clf.predict([[2, 3, 4, 2]])
```

```
>> array([1])
```

```
[64] 1 feature_importance = pd.Series(clf.feature_importances_, index=iris.feature_names).sort_values(ascending=False)
2 feature_importance
```

```
>> petal length (cm)    0.449624
    petal width (cm)    0.425776
    sepal length (cm)   0.100659
    sepal width (cm)    0.023941
    dtype: float64
```

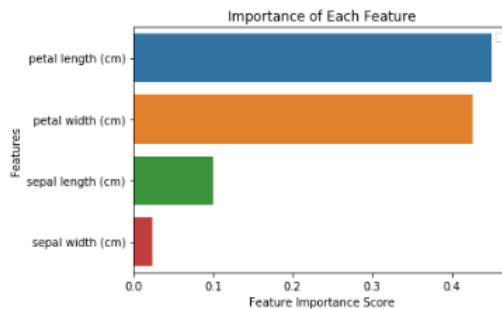
```
[65] 1 iris.feature_names
```

```
['sepal length (cm)',  
'sepal width (cm)',  
'petal length (cm)',  
'petal width (cm)']
```

```
[66] 1 import matplotlib.pyplot as plt  
2 import seaborn as sns
```

```
> 1 %matplotlib inline  
2 sns.barplot(x=feature_importance, y=feature_importance.index)  
3 plt.xlabel('Feature Importance Score')  
4 plt.ylabel('Features')  
5 plt.title('Importance of Each Feature')  
6 plt.legend()  
7 plt.show()
```

No handles with labels found to put in legend.



```
[68] 1 # let's remove the feature that's almost irrelevant for the prediction  
2 X = data[['sepal length', 'petal length', 'petal width']] # Features  
3 y = data['species'] # Labels
```

```
[69] 1 clf.fit(X_train, y_train)  
2 y_pred = clf.predict(X_test)  
3 from sklearn import metrics  
4 print('Accuracy', metrics.accuracy_score(y_test, y_pred))
```

```
> Accuracy 0.9555555555555556
```

Use Grid Search with Cross-Validation to select the best parameters.

Grid search: Grid search is the simplest algorithm for hyperparameter tuning. Basically, we divide the domain of the hyperparameters into a discrete grid. Then, try every combination of values of this grid, calculating some performance metrics using cross-validation. The point of the grid that maximizes the average value in cross-validation, is the optimal combination of values for the hyperparameters.

Cross-Validation: Cross-validation is a technique for evaluating ML models by training several ML models on subsets of the available input data and evaluating them on the complementary subset of the data. Use cross-validation to detect overfitting, failing to generalize a pattern

```
Corss_validation

1  from sklearn.model_selection import train_test_split
2  from sklearn.datasets import load_iris
3  from sklearn.neighbors import KNeighborsClassifier
4  from sklearn.metrics import accuracy_score
5
6  # Configures k-NN with 1 neighbor
7  model = KNeighborsClassifier(n_neighbors=1)
8
9  iris = load_iris()
10 X = iris.data
11 y = iris.target
12
13 # divide the original dataset in 2 equal parts
14 X1, X2, y1, y2 = train_test_split(X, y, random_state=0, train_size=0.5)
15
16 # ajusta e avalia DOIS modelos
17
18 y2_model = model.fit(X1, y1).predict(X2)
19
20 y1_model = model.fit(X2, y2).predict(X1)
21
22 accuracy_score(y1, y1_model), accuracy_score(y2, y2_model)
```

(0.96, 0.9066666666666666)

```
[64] 1 from sklearn.model_selection import cross_val_score
      2 scores = cross_val_score(model, X, y, cv=10)
      3 print('Scores in each fold: ', scores)
      4 print('Average score: ', scores.mean())
      5 print('Std deviation of scores: ', scores.std())
```

```
↳ Scores in each fold: [1.  0.93 1.  0.93 0.87 1.  0.87 1.  1.  1. ]
Average score: 0.96
Std deviation of scores: 0.05333333333333332
```

```
▷ 1 from sklearn.model_selection import LeaveOneOut
   2
   3 scores = cross_val_score(model, X, y, cv=LeaveOneOut())
   4 scores
   5 print('Average score: ', scores.mean())
   6 print('Std deviation of scores: ', scores.std())
```

```
↳ Average score: 0.96
Std deviation of scores: 0.19595917942265423
```

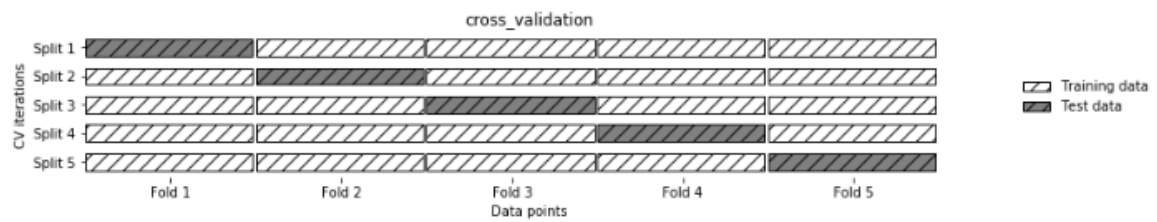
```
[4] 1 from sklearn.datasets import make_blobs
      2 from sklearn.linear_model import LogisticRegression
      3 from sklearn.model_selection import train_test_split
      4
      5 # create a synthetic dataset
      6 X, y = make_blobs(random_state=0)
      7 # split data and labels into a training and a test set
      8 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
      9 # instantiate a model and fit it to the training set
     10 logreg = LogisticRegression().fit(X_train, y_train)
     11 # evaluate the model on the test set
     12 print("Test set score: {:.2f}".format(logreg.score(X_test, y_test)))
     13
```

```
↳ Test set score: 0.88
```

```

1 import mglearn
2 mglearn.plots.plot_cross_validation()

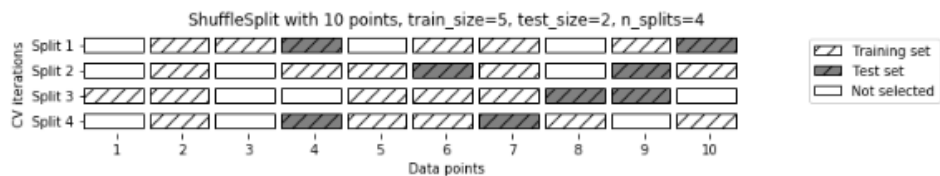
```



```

1 mglearn.plots.plot_shuffle_split()

```



grid search

```

[85] 1 import numpy as np
      2 import pandas as pd
      3 import seaborn as sns
      4 import matplotlib.pyplot as plt
      5 from sklearn import svm, datasets
      6 data = pd.read_csv("data/data.csv")
      7 data.head()

```



| | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |

9

```
# Column Non-Null Count Dtype
0 sepal_length 150 non-null float64
1 sepal_width 150 non-null float64
2 petal_length 150 non-null float64
3 petal_width 150 non-null float64
4 species 150 non-null object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

```
[90] 1 from sklearn.model_selection import GridSearchCV
```

```
[91] 1 param_grid = {'C':[0.1,1,10,100], 'gamma':[1,0.1,0.01,0.001]}
```

```
[92] 1 grid = GridSearchCV(SVC(), param_grid, refit = True, verbose=3)
      2 grid.fit(X_train, y_train)
```

```

[4] Fitting 5 folds for each of 16 candidates, totalling 80 fits
[CV] C=0.1, gamma=1 .....
[CV] ..... C=0.1, gamma=1, score=1.000, total= 0.0s
[CV] C=0.1, gamma=1 .....
[CV] ..... C=0.1, gamma=1, score=0.913, total= 0.0s
[CV] C=0.1, gamma=1 .....
[CV] ..... C=0.1, gamma=1, score=1.000, total= 0.0s
[CV] C=0.1, gamma=1 .....
[CV] ..... C=0.1, gamma=1, score=0.909, total= 0.0s
[CV] C=0.1, gamma=1 .....
[CV] ..... C=0.1, gamma=1, score=0.955, total= 0.0s
[CV] C=0.1, gamma=0.1 .....
[CV] ..... C=0.1, gamma=0.1, score=0.913, total= 0.0s

```

```
[93]      1      pred_grid = grid.predict(X_test)
```

```
[94] 1 print(confusion_matrix(y_test, pred_grid))
```

$$\begin{bmatrix} 13 & 0 & 0 \\ 0 & 15 & 1 \\ 0 & 0 & 9 \end{bmatrix}$$

```
[95] 1 print(classification_report(y_test, pred_grid))
```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 13 |
| 1 | 1.00 | 0.94 | 0.97 | 16 |
| 2 | 0.90 | 1.00 | 0.95 | 9 |
| accuracy | | | 0.97 | 38 |
| macro avg | 0.97 | 0.98 | 0.97 | 38 |
| weighted avg | 0.98 | 0.97 | 0.97 | 38 |

```

> 1 # naive grid search implementation
2 from sklearn.svm import SVC
3 from sklearn.datasets import load_iris
4
5 X_train, X_test, y_train, y_test = train_test_split(
6     iris.data, iris.target, random_state=0)
7 print("Size of training set: {}    size of test set: {}".format(
8     X_train.shape[0], X_test.shape[0]))
9
10 best_score = 0
11
12 for gamma in [0.001, 0.01, 0.1, 1, 10, 100]:
13     for C in [0.001, 0.01, 0.1, 1, 10, 100]:
14         # for each combination of parameters, train an SVC
15         svm = SVC(gamma=gamma, C=C)
16         svm.fit(X_train, y_train)
17         # evaluate the SVC on the test set
18         score = svm.score(X_test, y_test)
19         # if we got a better score, store the score and parameters
20         if score > best_score:
21             best_score = score
22             best_parameters = {'C': C, 'gamma': gamma}
23
24 print("Best score: {:.2f}".format(best_score))
25 print("Best parameters: {}".format(best_parameters))

```

```

> Size of training set: 112    size of test set: 38
Best score: 0.97
Best parameters: {'C': 100, 'gamma': 0.001}

```

Grid Search with Cross-Validation

```
1  from sklearn.model_selection import cross_val_score
2  from sklearn.svm import SVC
3  from sklearn.datasets import load_iris
4  import numpy as np
5  # split data into train+validation set and test set
6  X_trainval, X_test, y_trainval, y_test = train_test_split(
7      iris.data, iris.target, random_state=0)
8  # split train+validation set into training and validation sets
9  X_train, X_valid, y_train, y_valid = train_test_split(
10     X_trainval, y_trainval, random_state=1)
11  print("Size of training set: {}    size of validation set: {}    size of test set:"
12        " {} \n".format(X_train.shape[0], X_valid.shape[0], X_test.shape[0]))
13
14  best_score = 0
15
16  for gamma in [0.001, 0.01, 0.1, 1, 10, 100]:
17      for C in [0.001, 0.01, 0.1, 1, 10, 100]:
18          # for each combination of parameters, train an SVC
19          svm = SVC(gamma=gamma, C=C)
20          svm.fit(X_train, y_train)
21          # evaluate the SVC on the validation set
22          score = svm.score(X_valid, y_valid)
23          # if we got a better score, store the score and parameters
24          if score > best_score:
25              best_score = score
26              best_parameters = {'C': C, 'gamma': gamma}
27
28  # rebuild a model on the combined training and validation set,
29  # and evaluate it on the test set
30  svm = SVC(**best_parameters)
31  svm.fit(X_trainval, y_trainval)
32  test_score = svm.score(X_test, y_test)
33  print("Best score on validation set: {:.2f}".format(best_score))
34  print("Best parameters: ", best_parameters)
35  print("Test set score with best parameters: {:.2f}".format(test_score))
```

```
Size of training set: 84    size of validation set: 28    size of test set: 38
```

```
Best score on validation set: 0.96
```

```
Best parameters:  {'C': 10, 'gamma': 0.001}
```

```
Test set score with best parameters: 0.92
```

```

1  from sklearn.datasets import load_iris
2  for gamma in [0.001, 0.01, 0.1, 1, 10, 100]:
3      for C in [0.001, 0.01, 0.1, 1, 10, 100]:
4          # for each combination of parameters,
5          # train an SVC
6          svm = SVC(gamma=gamma, C=C)
7          # perform cross-validation
8          scores = cross_val_score(svm, X_trainval, y_trainval, cv=5)
9          # compute mean cross-validation accuracy
10         score = np.mean(scores)
11         # if we got a better score, store the score and parameters
12         if score > best_score:
13             best_score = score
14             best_parameters = {'C': C, 'gamma': gamma}
15 # rebuild a model on the combined training and validation set
16 svm = SVC(**best_parameters)
17 svm.fit(X_trainval, y_trainval)

```

```

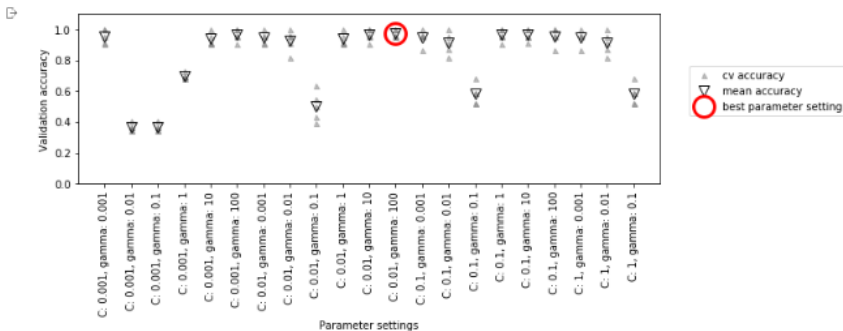
SVC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma=0.1, kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```

```

> 1  mglearn.plots.plot_cross_val_selection()

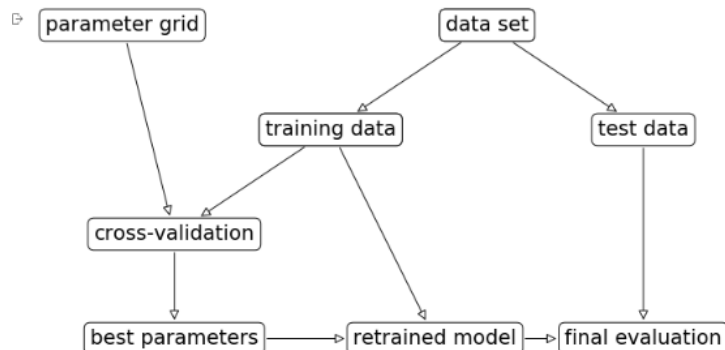
```



```

77] 1  mglearn.plots.plot_grid_search_overview()

```



```
[78] 1 param_grid = {'C': [0.001, 0.01, 0.1, 1, 10, 100],
2           |       |       | 'gamma': [0.001, 0.01, 0.1, 1, 10, 100]}
3     print("Parameter grid:\n{}".format(param_grid))
```

Parameter grid:

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

```
> 1 from sklearn.model_selection import GridSearchCV
2   from sklearn.datasets import load_iris
3   from sklearn.svm import SVC
4   grid_search = GridSearchCV(SVC(), param_grid, cv=5,
5   |                           | return_train_score=True)
6   X_train, X_test, y_train, y_test = train_test_split(
7   | iris.data, iris.target, random_state=0)
8   grid_search.fit(X_train, y_train)
9   print("Test set score: {:.2f}".format(grid_search.score(X_test, y_test)))
10  print("Best parameters: {}".format(grid_search.best_params_))
11  print("Best cross-validation score: {:.2f}".format(grid_search.best_score_))
12  print("Best estimator:\n{}".format(grid_search.best_estimator_))
13  import pandas as pd
14  # convert to DataFrame
15  results = pd.DataFrame(grid_search.cv_results_)
16  # show the first 5 rows
17  display(results.head())
```

Test set score: 0.97

Best parameters: {'C': 10, 'gamma': 0.1}

Best cross-validation score: 0.97

Best estimator:

```
SVC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma=0.1, kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

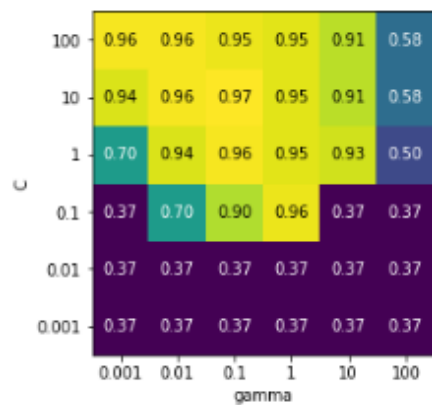
| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_C | param_gamma | params | split0_test_score | split1_test_score | s |
|---|---------------|--------------|-----------------|----------------|---------|-------------|------------------------------|-------------------|-------------------|---|
| 0 | 0.001134 | 0.000283 | 0.000450 | 0.000129 | 0.001 | 0.001 | {'C': 0.001, 'gamma': 0.001} | 0.347826 | 0.347826 | 0 |
| 1 | 0.000839 | 0.000094 | 0.000343 | 0.000028 | 0.001 | 0.01 | {'C': 0.001, 'gamma': 0.01} | 0.347826 | 0.347826 | 0 |
| 2 | 0.000913 | 0.000022 | 0.000378 | 0.000010 | 0.001 | 0.1 | {'C': 0.001, 'gamma': 0.1} | 0.347826 | 0.347826 | 0 |
| 3 | 0.001001 | 0.000170 | 0.000397 | 0.000059 | 0.001 | 1 | {'C': 0.001, 'gamma': 1} | 0.347826 | 0.347826 | 0 |
| 4 | 0.000906 | 0.000022 | 0.000374 | 0.000015 | 0.001 | 10 | {'C': 0.001, 'gamma': 10} | 0.347826 | 0.347826 | 0 |

```

1 scores = np.array(results.mean_test_score).reshape(6, 6)
2
3 # plot the mean cross-validation scores
4 mglearn.tools.heatmap(scores, xlabel='gamma', xticklabels=param_grid['gamma'],
5                          ylabel='C', yticklabels=param_grid['C'], cmap="viridis")

```

<matplotlib.collections.PolyCollection at 0x7f7374238a10>

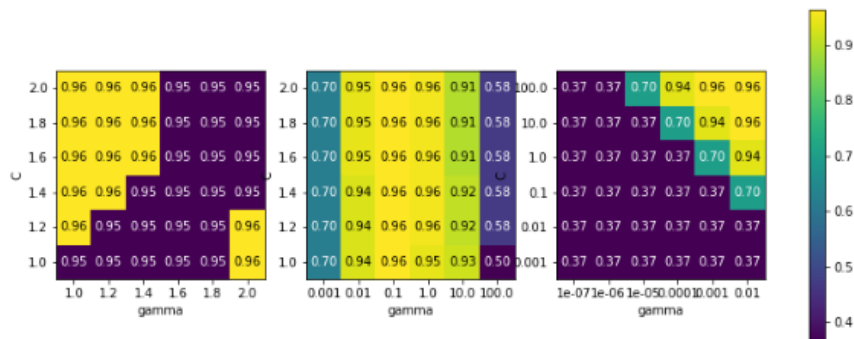


```

1 import matplotlib.pyplot as plt
2 from sklearn.datasets import load_iris
3 fig, axes = plt.subplots(1, 3, figsize=(13, 5))
4
5 param_grid_linear = {'C': np.linspace(1, 2, 6),
6                       'gamma': np.linspace(1, 2, 6)}
7
8 param_grid_one_log = {'C': np.linspace(1, 2, 6),
9                       'gamma': np.logspace(-3, 2, 6)}
10
11 param_grid_range = {'C': np.logspace(-3, 2, 6),
12                     'gamma': np.logspace(-7, -2, 6)}
13
14 for param_grid, ax in zip([param_grid_linear, param_grid_one_log,
15                           param_grid_range], axes):
16     grid_search = GridSearchCV(SVC(), param_grid, cv=5)
17     grid_search.fit(X_train, y_train)
18     scores = grid_search.cv_results_['mean_test_score'].reshape(6, 6)
19
20     # plot the mean cross-validation scores
21     scores_image = mglearn.tools.heatmap(
22         scores, xlabel='gamma', ylabel='C', xticklabels=param_grid['gamma'],
23         yticklabels=param_grid['C'], cmap="viridis", ax=ax)
24
25 plt.colorbar(scores_image, ax=axes.tolist())

```

<matplotlib.colorbar.Colorbar at 0x7f7261926dd0>



Use at least two metrics other than accuracy to evaluate your models.

1- Confusion Matrix

Let's first make sure we know the basic terminologies used in classification problems before going through the detail of each metrics.

One of the key concepts in classification performance is confusion matrix (AKA error matrix), which is a tabular visualization of the model predictions versus the ground-truth labels. Each row of confusion matrix represents the instances in a predicted class and each column represents the instances in an actual class.

Regression Related Metrics

Regression models are another family of machine learning and statistical models, which are used to predict continuous target values. They have a wide range of applications, from house price prediction, E-commerce pricing systems, weather forecasting, stock market prediction, to image super-resolution, feature learning via auto-encoders, and image compression.

⌕

Confusion matrices

```
1 import warnings
2 import pandas as pd
3 from sklearn import model_selection
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import confusion_matrix
6 import matplotlib.pyplot as plt
7 %matplotlib inline
8
9 #ignore warnings
10 warnings.filterwarnings('ignore')
11 # Load digits dataset
12 url = "http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
13 df = pd.read_csv(url)
14 # df = df.values
15 X = df.iloc[:,0:4]
16 y = df.iloc[:,4]
17 #test size
18 test_size = 0.33
19 #generate the same set of random numbers
20 seed = 7
21 #Split data into train and test set.
22 X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, test_size=test_size, random_state=seed)
23 #Train Model
24 model = LogisticRegression()
25 model.fit(X_train, y_train)
26 pred = model.predict(X_test)
```

^

```

27
28 #Construct the Confusion Matrix
29 labels = ['Iris-setosa', 'Iris-versicolor', 'Iris-virginica']
30 cm = confusion_matrix(y_test, pred, labels)
31 print(cm)
32 fig = plt.figure()
33 ax = fig.add_subplot(111)
34 cax = ax.matshow(cm)
35 plt.title('Confusion matrix')
36 fig.colorbar(cax)
37 ax.set_xticklabels([''] + labels)
38 ax.set_yticklabels([''] + labels)
39 plt.xlabel('Predicted Values')
40 plt.ylabel('Actual Values')
41 plt.show()

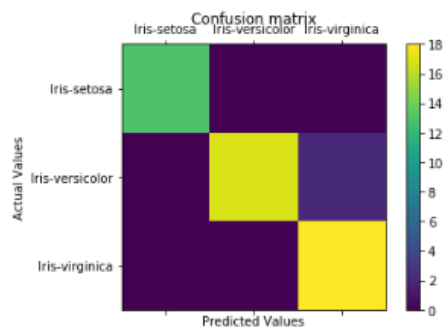
```

Execution time: 1seconds639milliseconds Finished at: 2022-05-25 19:58:12

```

[[13  0  0]
 [ 0 17  2]
 [ 0  0 18]]

```



Accuracy

```
> 1 #import modules
2 import warnings
3 import pandas as pd
4 import numpy as np
5 from sklearn import model_selection
6 from sklearn.linear_model import LogisticRegression
7 from sklearn import datasets
8 from sklearn.metrics import accuracy_score
9 #ignore warnings
10 warnings.filterwarnings('ignore')
11 # Load digits dataset
12 iris = datasets.load_iris()
13 # Create feature matrix
14 X = iris.data
15 # Create target vector
16 y = iris.target
17 #test size
18 test_size = 0.33
19 #generate the same set of random numbers
20 seed = 7
21 #cross-validation settings
22 kfold = model_selection.KFold(n_splits=10, random_state=seed)
23 #Model instance
24 model = LogisticRegression()
25 #Evaluate model performance
26 scoring = 'accuracy'
27 results = model_selection.cross_val_score(model, X, y, cv=kfold, scoring=scoring)
28 print('Accuracy -val set: %.2f%% (%.2f)' % (results.mean()*100, results.std()))
29 #split data
30 X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, test_size=test_size, random_state=seed)
31 #fit model
32 model.fit(X_train, y_train)
33 #accuracy on test set
34 result = model.score(X_test, y_test)
35 print("Accuracy - test set: %.2f%%" % (result*100.0))

Accuracy -val set: 94.67% (0.06)
Accuracy - test set: 92.00%
```

```

> 1 # confusion matrix in sklearn
2 from sklearn.metrics import confusion_matrix
3 from sklearn.metrics import classification_report
4
5 # actual values
6 actual = [1,0,0,1,0,0,1,0,0,1]
7 # predicted values
8 predicted = [1,0,0,1,0,0,0,1,0,0]
9
10 # confusion matrix
11 matrix = confusion_matrix(actual,predicted, labels=[1,0])
12 print('Confusion matrix : \n',matrix)
13
14 # outcome values order in sklearn
15 tp, fn, fp, tn = confusion_matrix(actual,predicted,labels=[1,0]).reshape(-1)
16 print('Outcome values : \n', tp, fn, fp, tn)
17
18 # classification report for precision, recall f1-score and accuracy
19 matrix = classification_report(actual,predicted,labels=[1,0])
20 print('Classification report : \n',matrix)

```

⇒ Confusion matrix :

```

[[2 2]
 [1 5]]

```

Outcome values :

```

2 2 1 5

```

Classification report :

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 1 | 0.67 | 0.50 | 0.57 | 4 |
| 0 | 0.71 | 0.83 | 0.77 | 6 |
| accuracy | | | 0.70 | 10 |
| macro avg | 0.69 | 0.67 | 0.67 | 10 |
| weighted avg | 0.70 | 0.70 | 0.69 | 10 |

Iris (Plots correlation matrices)

```
[26] 1 import pandas as pd
    2 pd.options.mode.chained_assignment = None
    3 dataframe = pd.read_csv("data/Iris.csv")
    4
    5 from sklearn.naive_bayes import GaussianNB
    6 from sklearn import svm
    7 import seaborn as sns
    8 import matplotlib.pyplot as plt
    9 from sklearn.metrics import mean_squared_error
   10 from math import sqrt
```

Execution time: 9milliseconds Finished at: 2022-05-25 19:15:35

```
▷ 1 # Lets se how our dataframe looks like
   2 dataframe.head()
```

Execution time: 11milliseconds Finished at: 2022-05-25 19:14:08

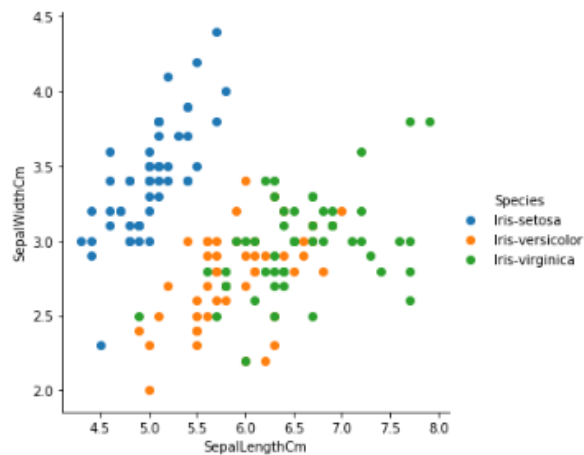


| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|----|---------------|--------------|---------------|--------------|-------------|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

```
[24] 1 sns.FacetGrid(dataframe, hue="Species", size=5) \
    2 .map(plt.scatter, "SepalLengthCm", "SepalWidthCm") \
    3 .add_legend()
```

Execution time: 417milliseconds Finished at: 2022-05-25 19:14:21

```
[> /opt/conda/envs/python35-paddle120-env/lib/python3.7/site-packages/seaborn/axisgrid.py:243: UserWarning: The 'size'
renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)
<seaborn.axisgrid.FacetGrid at 0x7fdd74a2e950>
```



```
> 1 #In order to use Naive_Bayes classifier we have to replace the "Species" values
    2 dataframe['Species'].replace("Iris-setosa",1,inplace=True)
    3 dataframe['Species'].replace("Iris-virginica",2,inplace=True)
    4 dataframe['Species'].replace("Iris-versicolor",3,inplace=True)
```

Execution time: 7milliseconds Finished at: 2022-05-25 19:19:21

```
[33] 1 #Now check if everything was changed properly
    2 dataframe['Species'].unique()
```

Execution time: 7milliseconds Finished at: 2022-05-25 19:19:28

```
[> array([1, 3, 2])
```

```
[32] 1 #In order to use Naive_Bayes classifier we have to replace the "Species" values
2 dataframe['Species'].replace("Iris-setosa",1,inplace=True)
3 dataframe['Species'].replace("Iris-virginica",2,inplace=True)
4 dataframe['Species'].replace("Iris-versicolor",3,inplace=True)
```

Execution time: 7milliseconds Finished at: 2022-05-25 19:19:21

```
[33] 1 #Now check if everything was changed properly
2 dataframe['Species'].unique()
```

Execution time: 7milliseconds Finished at: 2022-05-25 19:19:28

```
[> array([1, 2, 3])
```

```
[34] 1 X = dataframe.iloc[:, 0:4]
2 Y = dataframe['Species']
```

Execution time: 5milliseconds Finished at: 2022-05-25 19:19:37

```
[> 1 # I prefer to use train_test_split for cross-validation
2 # This piece will prove us if we have overfitting
3 X_train, X_test, y_train, y_test = train_test_split(
4     X, y, test_size=0.4, random_state=0)
5 print("X_train",X_train)
6 print("X_test",X_test)
7 print("y_train",y_train)
8 print("y_test",y_test)
```

Execution time: 21milliseconds Finished at: 2022-05-25 19:19:47

```
[> X_train      Id  SepalLengthCm  SepalWidthCm  PetalLengthCm
88      86           6.0           2.4           4.5
90      91           4.8           1.1           1.6
101     102           5.8           2.7           5.1
94      95           5.6           2.7           4.2
64      65           5.6           2.9           3.6
..      ..           ..           ..           ..
9       10           4.9           1.1           1.5
103     104           6.3           2.9           5.6
67      68           5.8           2.7           4.1
117     118           7.7           3.8           6.7
47      48           4.6           1.2           1.4
```

[90 rows x 5 columns]

```
[> 1 #Train and test model
2 clf = GaussianNB()
3 clf = clf.fit(X_train, y_train)
4 clf.score(X_test, y_test)
```

Execution time: 14milliseconds Finished at: 2022-05-25 19:20:09

```
[> 0.9833333333333333
```

Regression metrics

```
[43] 1 import pandas as pd
      2 # used to read the data set
      3 import numpy as np
      4 # used to do some operations with the arrays
      5 import os
      6 # used handle some files
      7 import matplotlib.pyplot as plt
      8 # used to visualize the data using graphs
      9 import seaborn as sns
     10 # plotting the chart in a single line
```

Execution time: 5milliseconds Finished at: 2022-05-25 19:28:56

```
> 1 df = pd.read_csv("data/Iris.csv")
```

Execution time: 8milliseconds Finished at: 2022-05-25 19:29:26

```
[45] 1 df.head(5)
```

Execution time: 13milliseconds Finished at: 2022-05-25 19:29:37



| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|----|---------------|--------------|---------------|--------------|-------------|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

```
> 1 df = df.drop(columns = ['Id'])
   2 df.head(5)
```

Execution time: 13milliseconds Finished at: 2022-05-25 19:29:46



| | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---------------|--------------|---------------|--------------|-------------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

```
[47] 1 df.info()
```

Execution time: 12milliseconds Finished at: 2022-05-25 19:29:55

```
> <class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 SepalLengthCm 150 non-null float64
1 SepalWidthCm 150 non-null float64
2 PetalLengthCm 150 non-null float64
3 PetalWidthCm 150 non-null float64
4 Species 150 non-null object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

```
> 1 df['Species'].value_counts()
```

Execution time: 9milliseconds Finished at: 2022-05-25 19:30:03

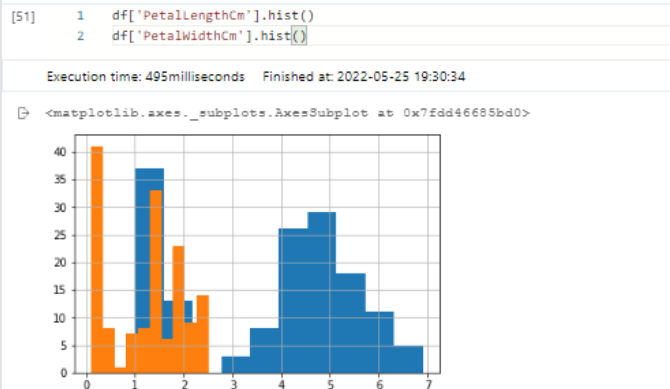
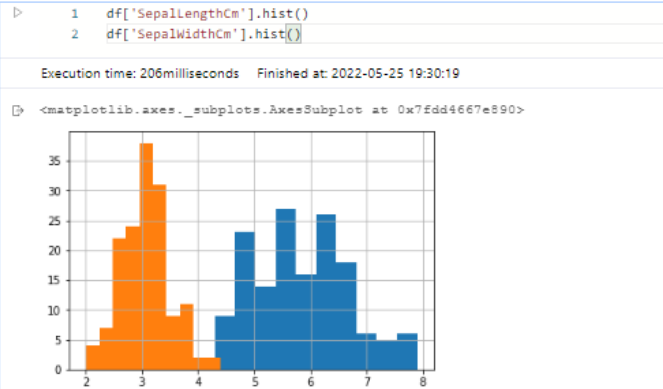
```
> Iris-virginica 50
Iris-versicolor 50
Iris-setosa 50
Name: Species, dtype: int64
```

```
[49] 1 df.isnull().sum()
```

Execution time: 9milliseconds Finished at: 2022-05-25 19:30:10

```
> SepalLengthCm 0
SepalWidthCm 0
PetalLengthCm 0
PetalWidthCm 0
Species 0
dtype: int64
```

```
[50] 1 df['SepalLengthCm'].hist()
2 df['SepalWidthCm'].hist()
```



```
1 df.corr()
```

Execution time: 12milliseconds Finished at: 2022-05-25 19:30:46

```
1 from sklearn.preprocessing import LabelEncoder
2 le = LabelEncoder()
3 df['Species'] = le.fit_transform(df['Species'])
4 df.head(100)
```

Execution time: 20milliseconds Finished at: 2022-05-25 19:30:54

| | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|-----|---------------|--------------|---------------|--------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | 0 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | 0 |
| ... | ... | ... | ... | ... | ... |
| 95 | 5.7 | 3.0 | 4.2 | 1.2 | 1 |
| 96 | 5.7 | 2.9 | 4.2 | 1.3 | 1 |
| 97 | 6.2 | 2.9 | 4.3 | 1.3 | 1 |
| 98 | 5.1 | 2.5 | 3.0 | 1.1 | 1 |
| 99 | 5.7 | 2.8 | 4.1 | 1.3 | 1 |

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

Execution time: 4milliseconds Finished at: 2022-05-25 19:31:04

```
1 X = df.drop(columns = ['Species'])
2 Y = df['Species']
3 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.25)
4 from sklearn.linear_model import LogisticRegression
5 model = LogisticRegression()
6 model.fit(X_train, Y_train)
```

Execution time: 49milliseconds Finished at: 2022-05-25 19:31:31

```
1 LogisticRegression()
```

```
1 print("Accuracy: ", model.score(X_test, Y_test) * 100)
```

Execution time: 8milliseconds Finished at: 2022-05-25 19:31:41

```
Accuracy: 97.36842105263158
```

Conclusion: I did this project using AI Baidu studio on my computer I also did take a lot of help from recommended books it was a great project for starting I did learn a lot of things. I did every class attentively and also did every homework that's why I did learn so much from this course I believe what I learn from this course I will be able to use this in my future.

Thank you, teacher,