

▼ 예제 데이터

IRIS 데이터

```
1 import pandas as pd
2 from sklearn.datasets import load_iris
3 from sklearn.tree import DecisionTreeClassifier
4
5 # Load data and store it into pandas DataFrame objects
6 iris=load_iris()
7 X = pd.DataFrame(iris.data[:, :], columns = iris.feature_names[:])
8 y = pd.DataFrame(iris.target, columns =["Species"])

1 import pandas as pd
2 df=pd.read_csv('http://203.247.53.31/Stat_Notes/example_data/iris.csv')
3 df.count()
```

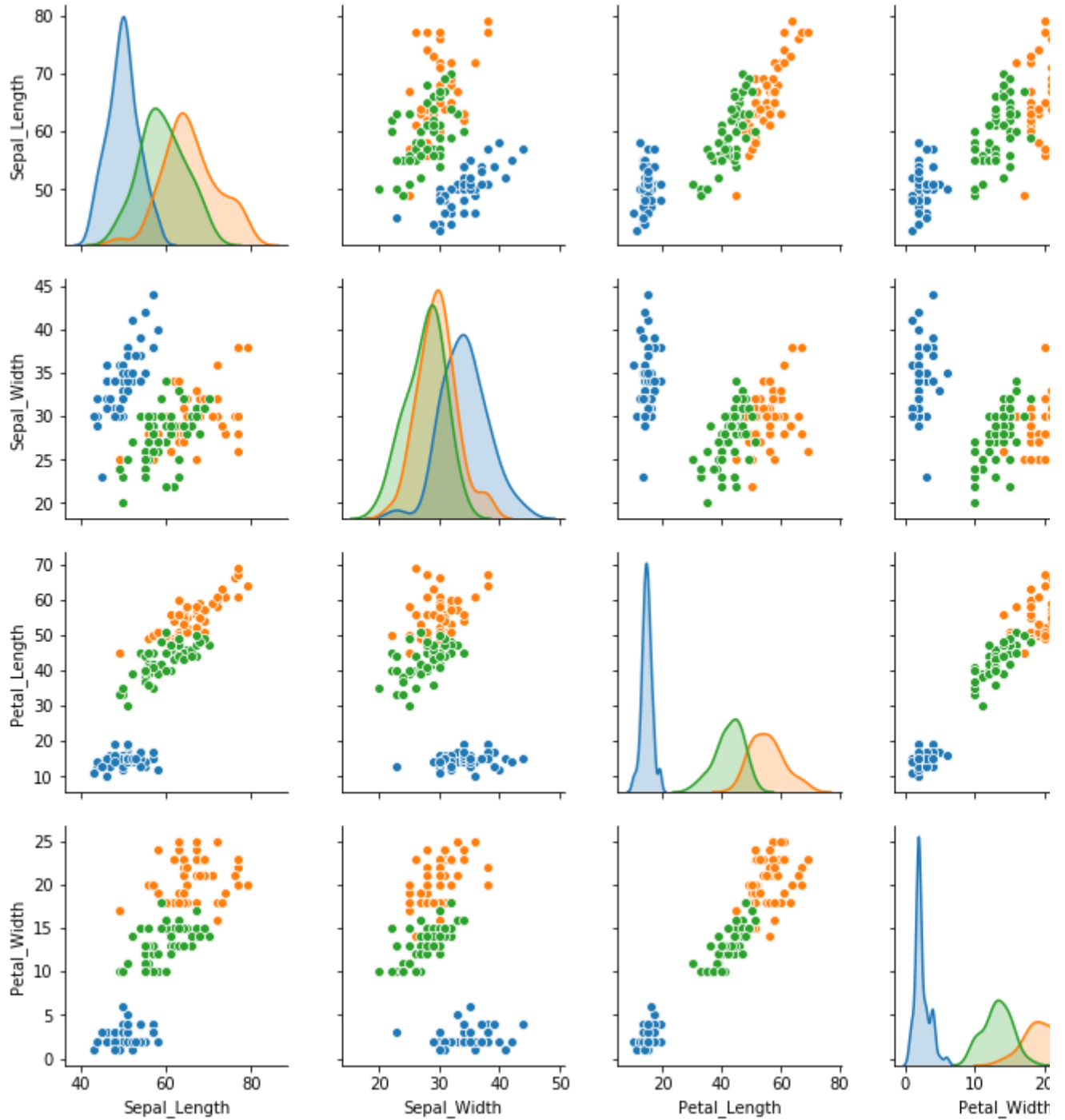
```
↳ Sepal_Length      150
   Sepal_Width       150
   Petal_Length      150
   Petal_Width       150
   group             150
   dtype: int64
```

▼ 집단별 판별변수 산점도

```
1 import seaborn as sns
2 sns.pairplot(df, hue='group')
```

```
↳
```

<seaborn.axisgrid.PairGrid at 0x7f680b3f2080>



▼ 데이터 검증 : 결측치여부, 변수정보

```
1 df.isnull().any()
```

```

1  Sepal_Length  False
2  Sepal_Width   False
3  Petal_Length  False
4  Petal_Width   False
5  group         False
6  dtype: bool

```

```
1 df.dtypes
```

```

↳ Sepal_Length      int64
   Sepal_Width      int64
   Petal_Length     int64
   Petal_Width      int64
   group            object
   dtype: object

```

▼ 판별변수, 목표집단변수 설정

```

1 X=df.iloc[:,0:4]
2 y=df.iloc[:,4]

```

▼ 방법1

```

1 from sklearn import tree
2 clf = tree.DecisionTreeClassifier(criterion='entropy', max_depth=2,min_samples_leaf=5)
3 clf = clf.fit(X,y)

1 from sklearn import metrics
2 def measure_performance(X,y,clf, show_accuracy=True, show_classification_report=True, show_cor
3     y_pred=clf.predict(X)
4     if show_accuracy:
5         print ("Accuracy:{0:.3f}".format(metrics.accuracy_score(y,y_pred)),"\n")
6
7     if show_classification_report:
8         print ("Classification report")
9         print (metrics.classification_report(y,y_pred),"\n")
10
11    if show_confusion_matrix:
12        print ("Confusion matrix")
13        print (metrics.confusion_matrix(y,y_pred),"\n")
14
15 measure_performance(X,y,clf, show_classification_report=False, show_confusion_matrix=False)

```

```

↳ Accuracy:0.960

```

test data evaluating

```

1 clf_dt=tree.DecisionTreeClassifier(criterion='entropy', max_depth=3,min_samples_leaf=5)
2 clf_dt.fit(X_train,y_train)
3 measure_performance(X_test,y_test,clf_dt)

```

▼ 방법2

```

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import confusion_matrix
3 from sklearn.tree import export_graphviz
4 from sklearn.externals.six import StringIO
5 from IPython.display import Image

```

```
6 from pydot import graph_from_dot_data
7 dt = DecisionTreeClassifier()
8 dt.fit(X, y)
```

```
↳ /usr/local/lib/python3.6/dist-packages/sklearn/externals/six.py:31: DeprecationWarning:
  "(https://pypi.org/project/six/).", DeprecationWarning)
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=None, splitter='best')
```

```
1 dot_data = StringIO()
2 export_graphviz(dt, out_file=dot_data, feature_names=iris.feature_names)
3 (graph, ) = graph_from_dot_data(dot_data.getvalue())
4 Image(graph.create_png())
```

```
↳
```

