



ESS2222

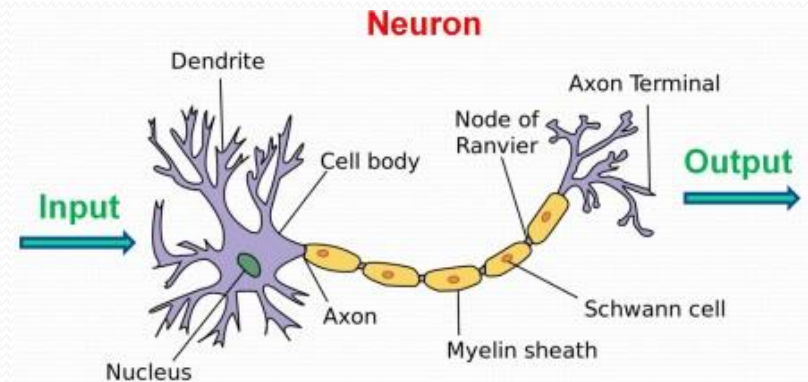
Lecture 9 – Decision Trees and Random Forest

Hosein Shahnas

University of Toronto, Department of Earth Sciences,

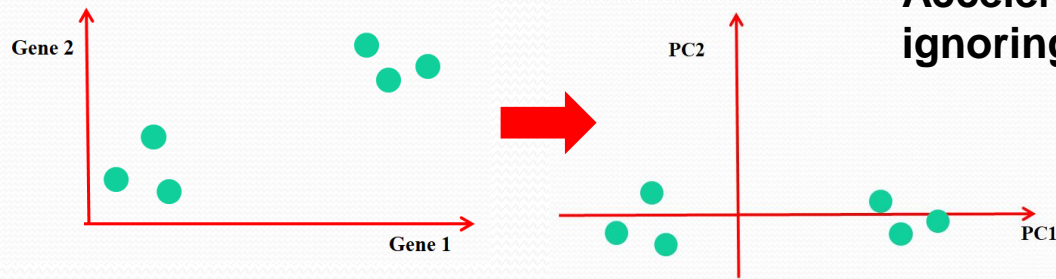
Outline

- ❑ Decision Trees
- ❑ Gini index
- ❑ Decision Tree – Iris Classification Problem
- ❑ Random Forest Algorithm
- ❑ Bagging – Bootstrap Sampling
- ❑ Random Forest – Image recognition
- ❑ Feature Importance



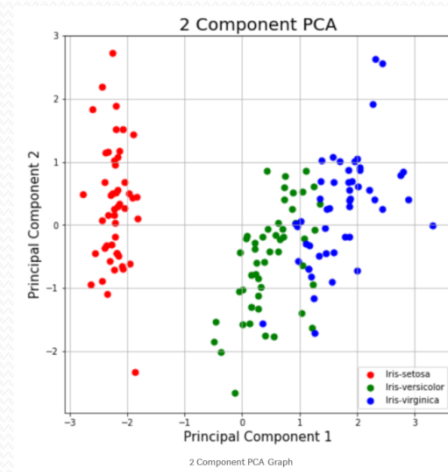
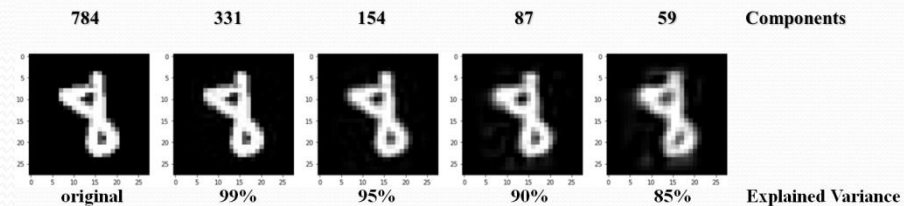
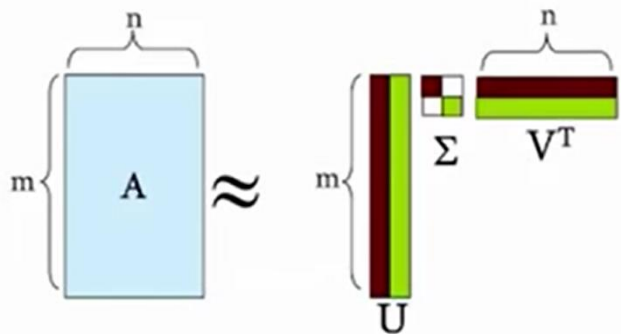
Review of Lecture 8

Accelerating the computational process by ignoring less important features



Singular Value Decomposition (SVD)

$$A_{[mn]} = U_{[mr]} \Sigma_{[rr]} V_{[nr]}^T$$

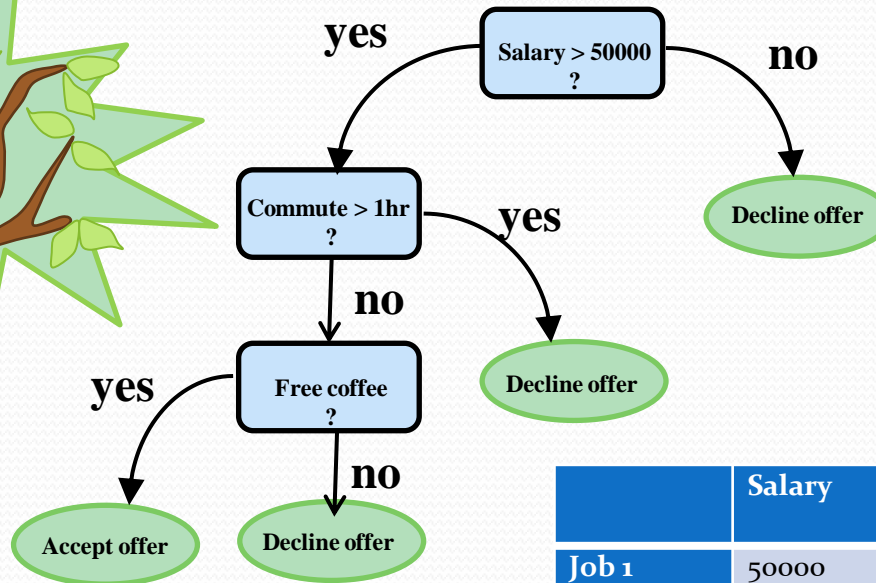
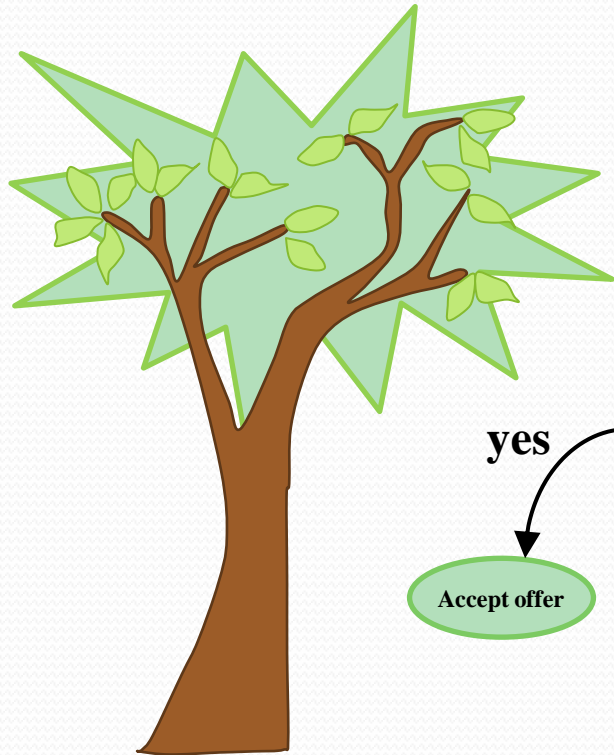


Visualization

Classification and Regression Tree (CART)

Decision Tree:

Should I accept the job?



Root node: The top most decision node

Decision nodes: Nodes that have two or more branches

Leaf nodes: Nodes of the tree that have no additional nodes coming off them

	Salary	Commute	Free coffee	Class
Job 1	50000	2 hr	Yes	1
Job 2	45000	3 hr	No	0
...
...
Job n	53000	1 hr	No	1

Tree-based methods are one of the commonly used supervised learning algorithms in machine learning, both for regression and classification problems.

Decision Trees

Decision trees:

A decision tree is a **simple** but **powerful** supervised learning method that uses tree-like model of decisions and their possible consequences. They are used in both **classification** and **regression** problems.

Unlike the highly dimensional SVM, where it is impossible for a human brain to even imagine how the hyperplane built looks like, the decision trees provide very **good visualization** on the steps of decisions and **the relative importance of the features**.

<code>class sklearn.tree.DecisionTreeRegressor</code>	<i>for classification</i>
<code>class sklearn.tree.DecisionTreeClassifier</code>	<i>for regression</i>

Gini Index

Gini index: The gini index is a number describing the quality of the split of a node on a variable (feature).

If a data set D contains samples from C classes, gini index is defined as:

$$gini(D) = 1 - \sum_{c=1}^C P_c^2$$

where P_c is the relative frequency of class c in D

If a data set D splits on S into two subsets D_1 and D_2 , the gini index is defined as:

$$gini_S(D) = \frac{D_1}{D} gini(D_1) + \frac{D_2}{D} gini(D_2)$$

where $gini(D_1) < gini(D)$, $gini(D_2) < gini(D)$

Reduction in impurity:

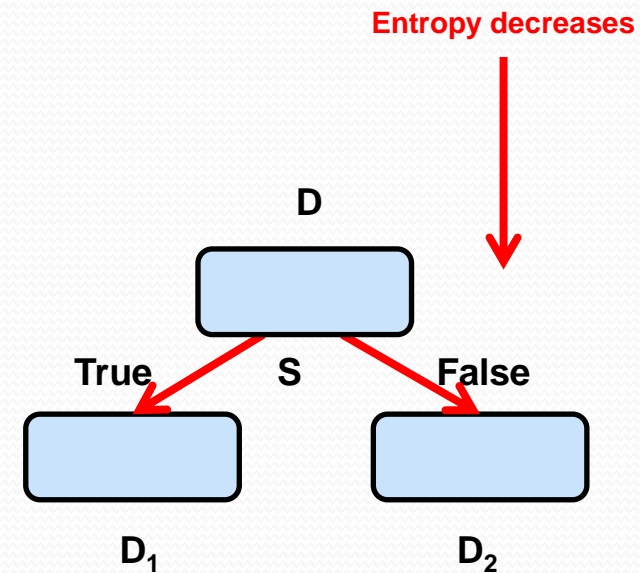
$$\Delta gini(S) = gini(D) - gini_S(D)$$

Information Gain Entropy:

Information gain is a measure of decrease in entropy after the data is split.

Information gain entropy is another criterion for choosing the features in splitting. Entropy is a measure of the degree of disorder or randomness in the system.

$$H(D) = 1 - \sum_{c=1}^C P_c \log(P_c)$$




Gini Index

Example: Playing Tennis

Features

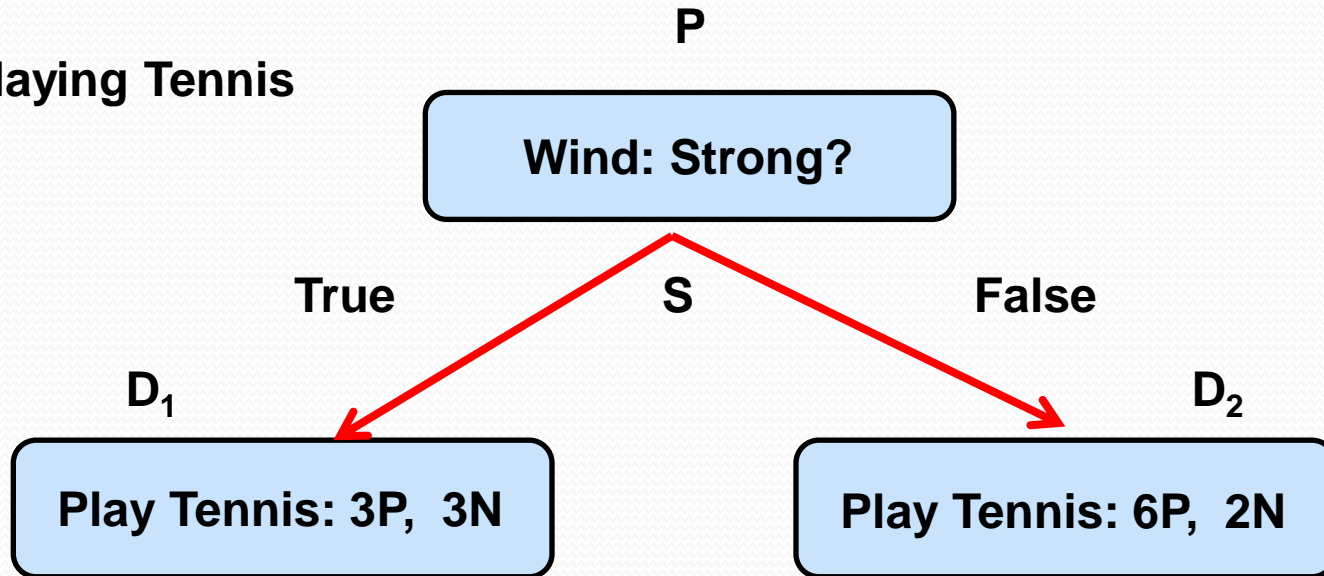
Class labels



Day	Outlook	Temp.	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Gini Index (Gini Impurity)

Example: Playing Tennis



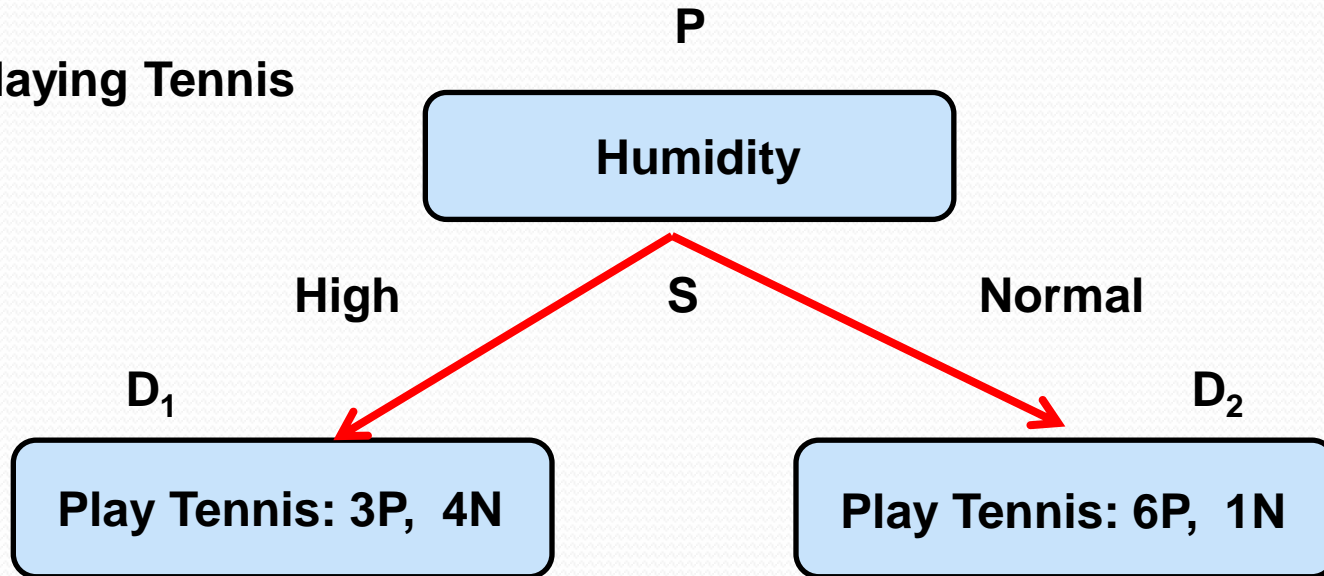
$$\text{gini}(\text{Play Tennis} \mid \text{Wind} = \text{True}) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.5$$

$$\text{gini}(\text{Play Tennis} \mid \text{Wind} = \text{False}) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.375$$

$$\begin{aligned} \text{gini}_S(D) &= \frac{D_1}{D} \text{gini}(D_1) + \frac{D_2}{D} \text{gini}(D_2) \\ &= \frac{6}{14} \times 0.5 + \frac{8}{14} \times 0.375 = 0.4286 \end{aligned}$$

Gini Index (Gini Impurity)

Example: Playing Tennis



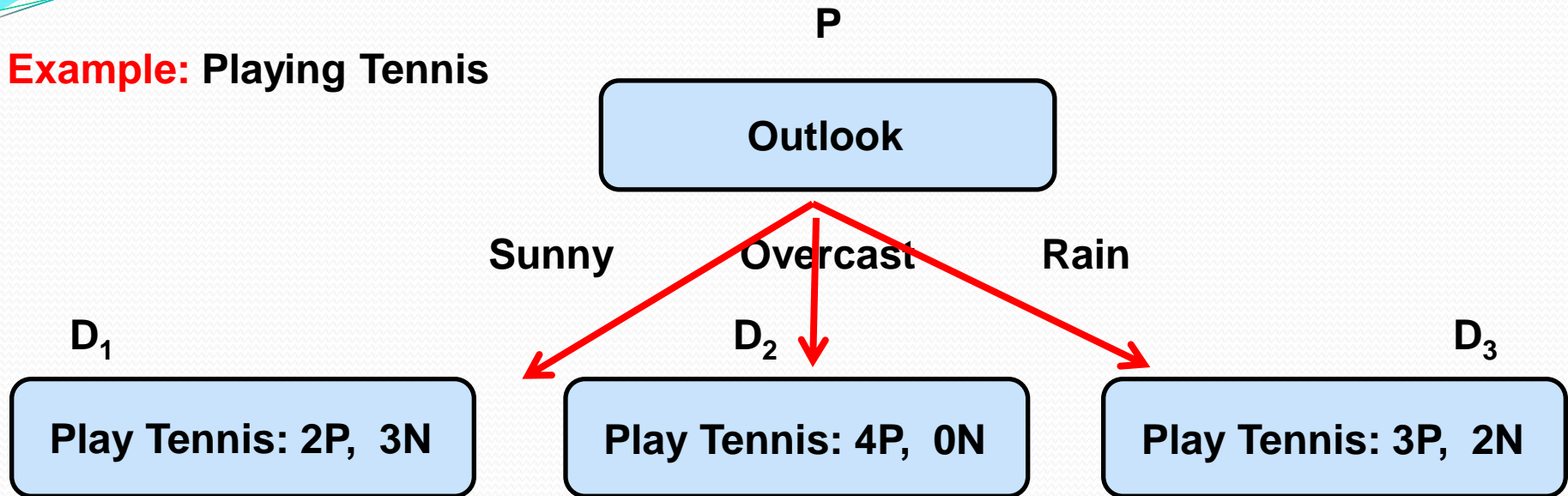
$$\text{gini (Play Tennis | Humidity = High)} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.4898$$

$$\text{gini (Play Tennis | Humidity = Normal)} = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.2449$$

$$\begin{aligned} \text{gini}_S(D) &= \frac{D_1}{D} \text{gini}(D_1) + \frac{D_2}{D} \text{gini}(D_2) \\ &= \frac{6}{14} \times 0.4898 + \frac{8}{14} \times 0.2449 = 0.3674 \end{aligned}$$

Gini Index (Gini Impurity)

Example: Playing Tennis



$$\text{gini}(\text{Play Tennis} \mid \text{Outlook} = \text{Sunny}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

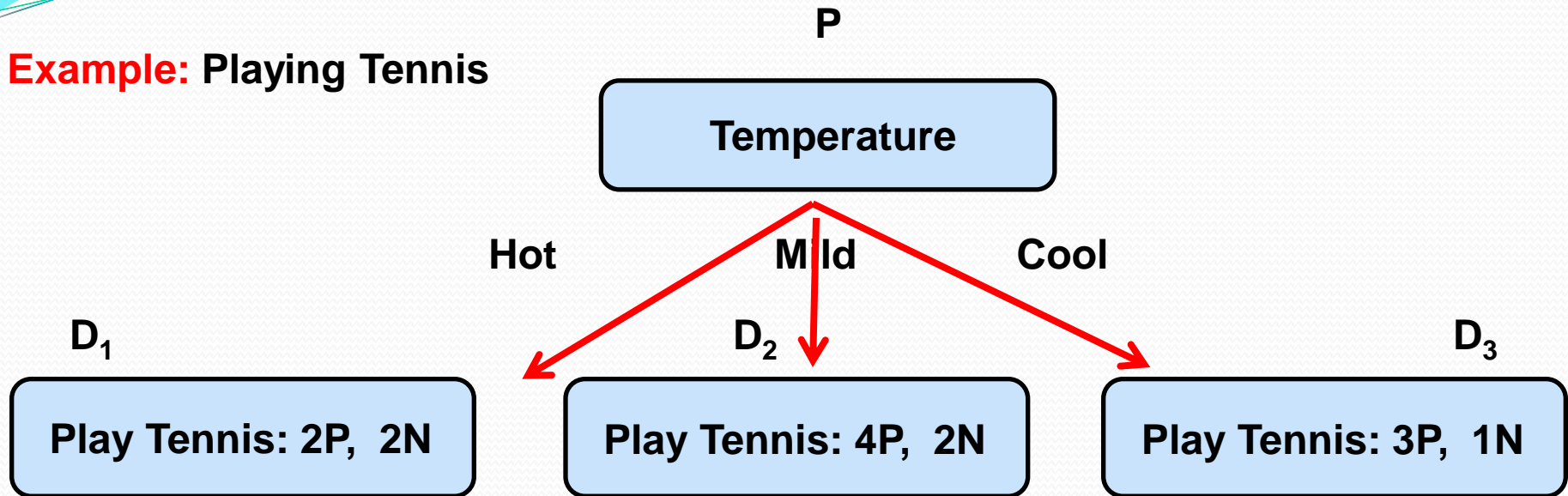
$$\text{gini}(\text{Play Tennis} \mid \text{Outlook} = \text{Overcast}) = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0$$

$$\text{gini}(\text{Play Tennis} \mid \text{Outlook} = \text{Rain}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

$$\begin{aligned} \text{gini}_S(D) &= \frac{D_1}{D} \text{gini}(D_1) + \frac{D_2}{D} \text{gini}(D_2) + \frac{D_3}{D} \text{gini}(D_3) \\ &= \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48 = 0.3429 \end{aligned}$$

Gini Index (Gini Impurity)

Example: Playing Tennis



$$\text{gini (Play Tennis | Temperature = Hot)} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.5$$

$$\text{gini (Play Tennis | Temperature = Mild)} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.4444$$

$$\text{gini (Play Tennis | Temperature = Cool)} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

$$\begin{aligned} \text{gini}_S(D) &= \frac{D_1}{D} \text{gini}(D_1) + \frac{D_2}{D} \text{gini}(D_2) + \frac{D_3}{D} \text{gini}(D_3) \\ &= \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.4444 + \frac{4}{14} \times 0.375 = 0.4405 \end{aligned}$$

Gini Index (Gini Impurity)

Example: Playing Tennis

$$gini_S(Windy) = 0.4286$$

$$gini_S(Humidity) = 0.3674$$

$$gini_S(Outlook) = 0.3429$$

$$gini_S(Temperature) = 0.4405$$

The splitting gain for Temperature feature is high, therefore we chose temperature as the feature for splitting.

$$gini(D) = 1 - \sum_{c=1}^C P_c^2$$

$$gini(D) = 1 - \left(\frac{5}{14}\right)^2 - \left(\frac{9}{14}\right)^2 = 0.541$$

Parameters for DecisionTreeClassifier

criterion : string, optional (default="gini")

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : string, optional (default="best")

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int or None, optional (default=None)

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

min_samples_split : int, float, optional (default=2)

The minimum number of samples required to split an internal node.

min_samples_leaf : int, float, optional (default=1)

The minimum number of samples required to be at a leaf node.

max_features : int, float, string or None, optional (default=None)

The number of features to consider when looking for the best split.

max_leaf_nodes : int or None, optional (default=None)

Grow a tree with max_leaf_nodes in best-first fashion. Best nodes are defined as relative reduction in impurity. If None then unlimited number of leaf nodes.

DecisionTreeRegressor

criterion : strstring, optional (default="mse")

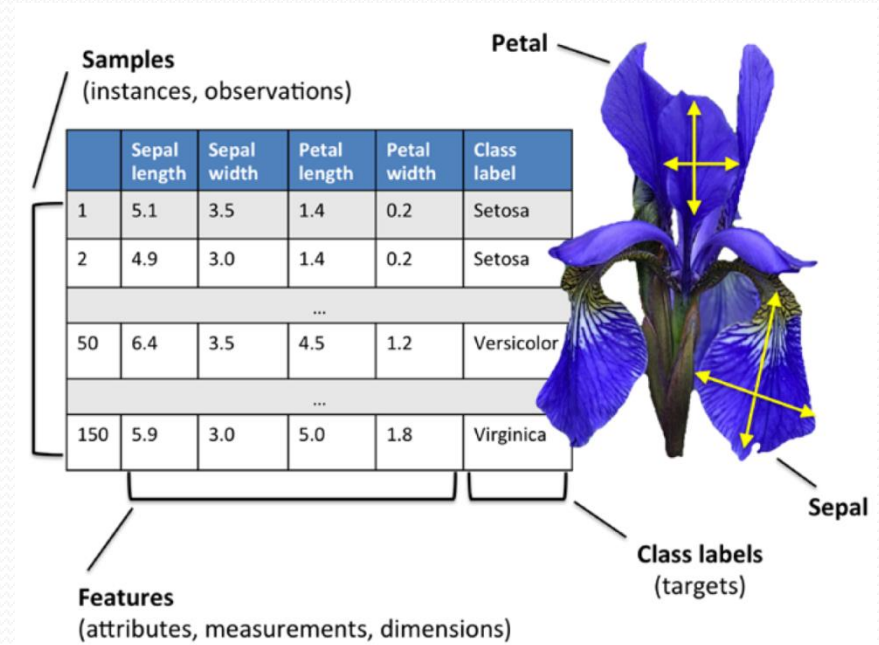
The function to measure the quality of a split. Supported criteria are "mse" for the mean squared error, which is equal to variance reduction as feature selection criterion and minimizes the L2 loss

Decision Tree - Iris Classification Problem

Iris problem:

Three class problem with four features.

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV for grid search
```



Decision Tree - Iris Problem

```
1#=====
2import pandas as pd
3import graphviz
4from sklearn.preprocessing import LabelEncoder
5from sklearn.tree import DecisionTreeClassifier
6from sklearn.model_selection import train_test_split
7from sklearn.metrics import accuracy_score
8import numpy as np
9import sys
10import os.path
11# Setting random seed.
12seed = 10
13#=====
14save_path = os.path.dirname(os.path.abspath(__file__))
15print(save_path)
16#=====
17#===== read data
18# Loading Iris dataset.
19data = pd.read_csv('data/iris.csv')
20print('data.head(3) = ', data.head(3))
21#===== read data
22
23#===== construct arrays from raw data
24n_s = 150 #number of samples
25y = data.iloc[0:n_s, 4].values # get data at the fifth (4) column for the first 100 lines
26X = data.iloc[0:n_s, [0,1 , 2, 3]].values
27print('X[0] = ', X[0])
28print('y[0] = ', y[0])
29
30Feature_size = int(X.size/y.size)
31print('Feature_size = ', Feature_size)
32print('X.shape = ', X.shape, 'X.shape[0] = ', X.shape[0], 'X.shape[1] = ', X.shape[1])
33
34y = np.reshape(y, (n_s,1))
35
36print('X[0]2 = ', X[0])
37print('y[0]2 = ', y[0])
38#===== construct arrays from raw data
39
```

Decision Tree - Iris Problem

```
39
40#===== train-test split
41X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y, random_state=seed)
42print('X_train.shape = ', X_train.shape)
43print('y_train.shape = ', y_train.shape)
44print('X_test.shape = ', X_test.shape)
45print('y_test.shape = ', y_test.shape)
46print('X_test.size = ', X_train.size)
47print('y_train.size = ', y_train.size)
48print('X_test.size = ', X_test.size)
49print('y_test.size = ', y_test.size)
50n_test = y_test.size
51n_train = y_train.size
52#===== train-test split
53
54#===== grid search
55'''
56from sklearn.model_selection import GridSearchCV
57parm_grid = dict(min_samples_leaf = [1,3],
58                min_samples_split = [2,5],
59                max_depth = [1,20],
60                max_features = [1,4])
61clf_tree=DecisionTreeClassifier(criterion='gini')
62grid=GridSearchCV(clf_tree,parm_grid, cv=7)
63grid.fit(X_train,y_train)
64print('=====')
65print("grid.cv_results_ {}".format(grid.cv_results_))
66print('=====')
67sys.exit()
68'''
69#===== grid search
70
71#===== decision tree
72'''
73tree = DecisionTreeClassifier(criterion='gini',
74                              splitter='best',
75                              max_depth=None,
76                              min_samples_split=2,
77                              min_samples_leaf=1,
78                              min_weight_fraction_leaf=0.0,
```


Decision Tree - Iris Problem

```
79         max_features=None,
80         random_state=None,
81         max_leaf_nodes=None,
82         min_impurity_decrease=0.0,
83         min_impurity_split=None,
84         class_weight=None,
85         presort=False)
86 '''
87 tree = DecisionTreeClassifier(criterion='gini',
88                               min_samples_leaf=5,
89                               min_samples_split=5,
90                               max_depth=None,
91                               random_state=seed)
92 #===== decision tree
93
94 #===== train
95 tree.fit(X_train, y_train)
96 #===== train
97
98 #===== feature importance
99 print()
100 print('tree.feature_importances_ = ', tree.feature_importances_)
101 print()
102 #===== feature importance
103
104 #===== predict test samples
105 y_pred = tree.predict(X_test)
106 accuracy = accuracy_score(y_test, y_pred)
107 print('DecisionTreeClassifier accuracy score: {}'.format(accuracy))
108 print('y_test.shape = ', y_test.shape)
109 print('y_pred.shape = ', y_pred.shape)
110 #print(y_test,y_pred)
111 #===== predict test samples
112
113 #===== predict train samples
114 y_pred_train = tree.predict(X_train)
115 #===== predict train samples
116
117 #===== Acc. Test data
```

Decision Tree - Iris Problem

```
118 file_name = 'Mis_class_deg_Test'
119 completeName = os.path.join(save_path, file_name+".dat")
120 file1 = open(completeName, "w")
121 file_name = 'Score_Test'
122 completeName = os.path.join(save_path, file_name+".dat")
123 file2 = open(completeName, "w")
124 file_name = 'Prediction_Test'
125 completeName = os.path.join(save_path, file_name+".dat")
126 file3 = open(completeName, "w")
127
128 rs = np.zeros(n_test)
129 rs2 = np.ones((n_test), dtype=bool)
130
131 for i in range(0,n_test):
132     if (y_pred[i] == y_test [i] ):
133         rs[i] = 0.0
134     else:
135         rs[i] = 1.0
136
137     if rs[i]==0.0:
138         rs2[i] = True
139         file3.write("%s %s \n" % (y_pred[i], y_test[i]))
140     else:
141         rs2[i] = False
142         file3.write("%s %s %s\n" % (y_pred[i], y_test[i], str(rs2[i])))
143
144     file1.write('{:4.2f}\n'.format(rs[i]))
145     file2.write(str(rs2[i]))
146     file2.write('\n')
147
148 file1.close()
149 file2.close()
150 file3.close()
151
152 err_test = 0.0
153 for i in range(0,n_test):
154     err_test = err_test + rs[i]
155 print('Number of errors (Test samples) = ', err_test)
156 print('Error (Test samples) = ', err_test/float(n_test))
```

Decision Tree - Iris Problem

```
157 print('Accuracy (Test samples) = ', 1.0 - err_test/float(n_test))
158 #===== Acc. Test data
159
160 #===== Acc. Train data
161 file_name = 'Mis_class_deg_Train'
162 completeName = os.path.join(save_path, file_name+".dat")
163 file1 = open(completeName, "w")
164 file_name = 'Score_Train'
165 completeName = os.path.join(save_path, file_name+".dat")
166 file2 = open(completeName, "w")
167 file_name = 'Prediction_Train'
168 completeName = os.path.join(save_path, file_name+".dat")
169 file3 = open(completeName, "w")
170
171 rs = np.zeros(n_train)
172 rs2 = np.ones((n_train), dtype=bool)
173
174 for i in range(0,n_train):
175     if (y_pred_train[i] == y_train [i] ):
176         rs[i] = 0.0
177     else:
178         rs[i] = 1.0
179
180     if rs[i]==0.0:
181         rs2[i] = True
182         file3.write("%s %s \n" % (y_pred_train[i], y_train[i]))
183     else:
184         rs2[i] = False
185         file3.write("%s %s %s\n" % (y_pred_train[i], y_train[i], str(rs2[i])))
186
187     file1.write('{:4.2f}\n'.format(rs[i]))
188     file2.write(str(rs2[i]))
189     file2.write('\n')
190
191 file1.close()
192 file2.close()
193 file3.close()
194
195 err_train = 0.0
```

Decision Tree - Iris Problem

```
196 for i in range(0,n_train):
197     err_train = err_train + rs[i]
198
199 print('Number of errors (Train samples) = ', err_train)
200 print('Error (Train samples) = ', err_train/float(n_train))
201 print('Accuracy (Train samples) = ', 1.0 - err_train/float(n_train))
202 #===== Acc. Train data
203
204 #===== write results
205 file_name = 'Results'
206 completeName = os.path.join(save_path, file_name+".dat")
207 file1 = open(completeName, "w")
208
209 file1.write("%s %4.2f \n" % ('Number of errors (Test samples) = ', err_test))
210 file1.write('\n')
211 file1.write("%s %4.2f \n" % ('Error (Test samples) = ', err_test/float(n_test)))
212 file1.write('\n')
213 file1.write("%s %4.2f \n" % ('Accuracy (Test samples) = ', 1.0 - err_test/float(n_test)))
214 file1.write('\n')
215
216 file1.write("%s %4.2f \n" % ('Number of errors (Train samples) = ', err_train))
217 file1.write('\n')
218 file1.write("%s %4.2f \n" % ('Error (Train samples) = ', err_train/float(n_train)))
219 file1.write('\n')
220 file1.write("%s %4.2f \n" % ('Accuracy (Train samples) = ', 1.0 - err_train/float(n_train)))
221 file1.write('\n')
222
223 file1.write("%s %s \n" % ('tree.feature_importances_ = ', tree.feature_importances_))
224 print('tree.feature_importances_ = ', tree.feature_importances_)
225 file1.write('\n')
226 file1.close()
227 #===== write results
228
229 #===== predict using test samples
230
```

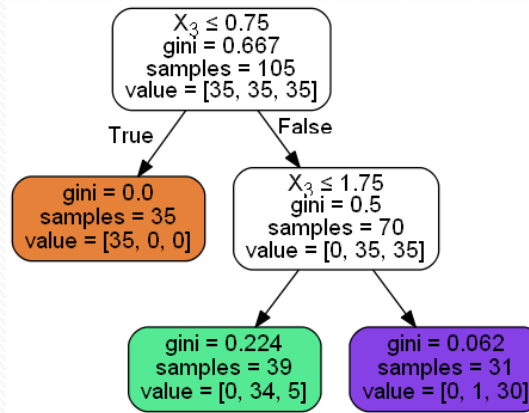
Decision Tree - Iris Problem

```
231#===== visualization of tree
232from sklearn.externals.six import StringIO
233from IPython.display import Image
234from sklearn.tree import export_graphviz
235import pydotplus
236
237'''
238export_graphviz(decision_tree,
239                out_file=None,
240                max_depth=None,
241                feature_names=None,
242                class_names=None,
243                label='all',
244                filled=False,
245                leaves_parallel=False,
246                impurity=True,
247                node_ids=False,
248                proportion=False,
249                rotate=False,
250                rounded=False,
251                special_characters=False,
252                precision=3
253'''
254dot_data = StringIO()
255
256export_graphviz(tree, out_file=dot_data,
257                filled=True, rounded=True,
258                special_characters=True)
259
260graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
261Image(graph.create_png())
262
263# Create PDF
264graph.write_pdf("iris.pdf")
265
266# Create PNG
267graph.write_png("iris.png")
268#===== visualization of tree
269sys.exit()
270
```

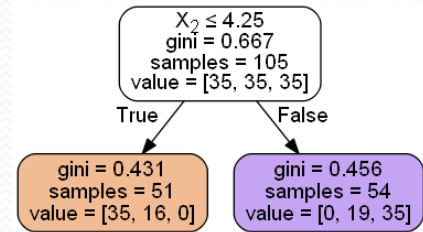
Decision Tree

Example 1: Iris problem

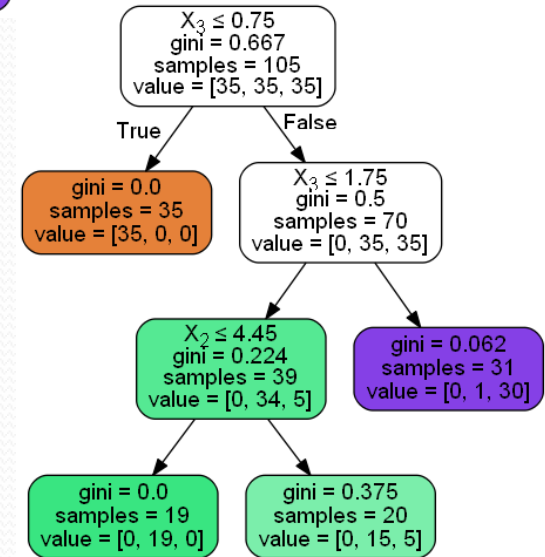
criterion='gini',
splitter='best',
min_samples_split=5,
min_samples_leaf = ?,
max_depth=20,
max_features=None,



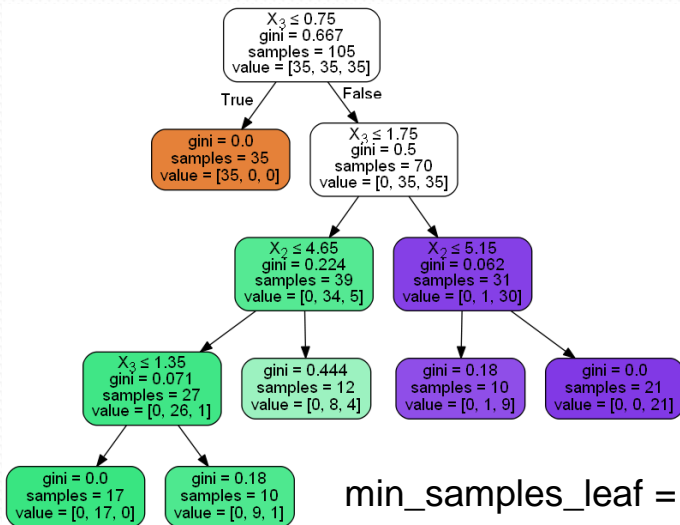
$\text{min_samples_leaf} = 25$



$\text{min_samples_leaf} = 50$



$\text{min_samples_leaf} = 17$

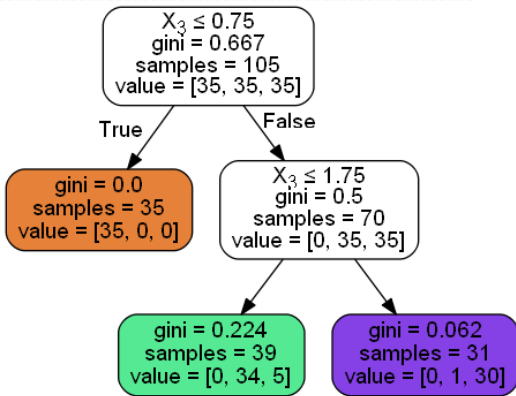


$\text{min_samples_leaf} = 10$

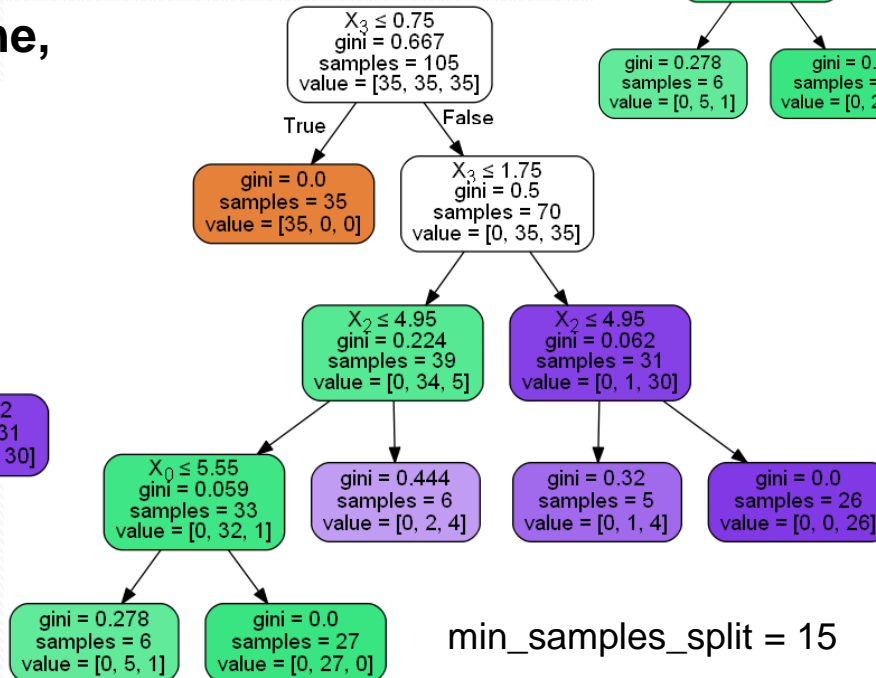
Decision Tree

Example 1: Iris problem

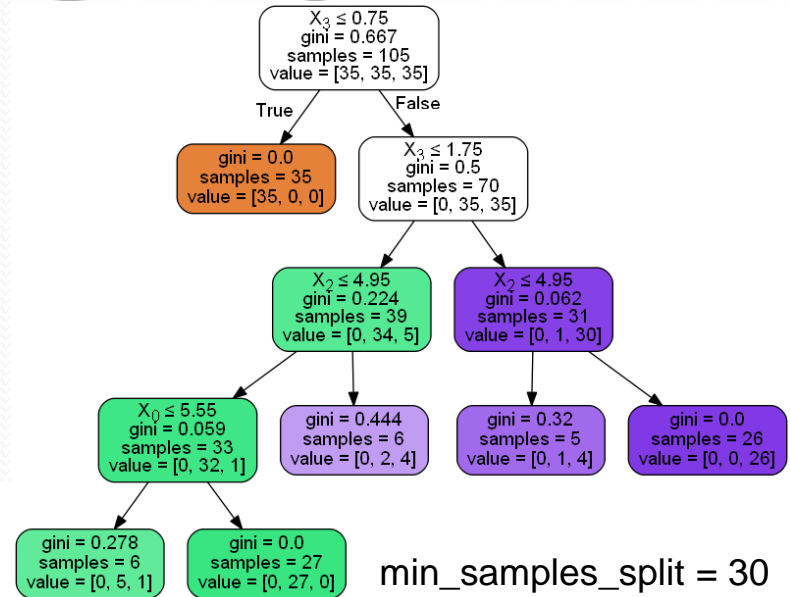
criterion='gini',
splitter='best',
min_samples_split=?,
min_samples_leaf=5,,
max_depth=20,
max_features=None,



min_samples_split = 65



min_samples_split = 15

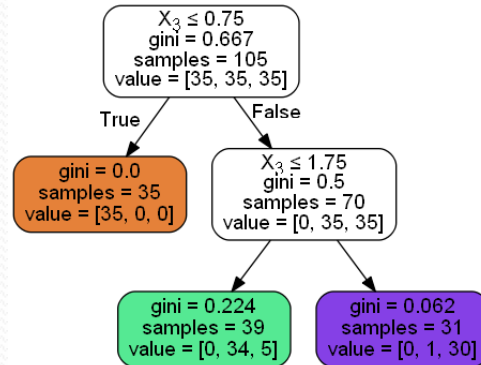


min_samples_split = 30

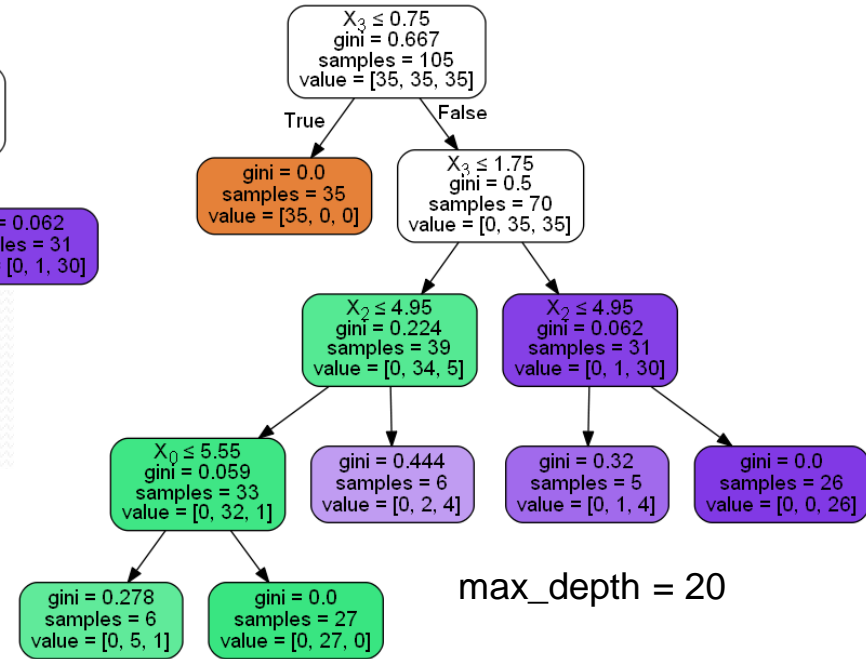
Decision Tree

Example 1: Iris problem

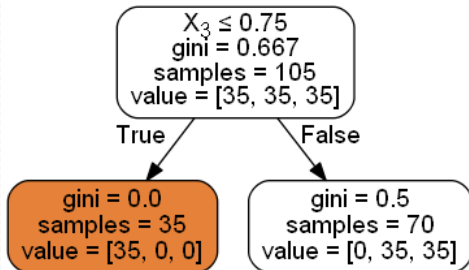
criterion='gini',
splitter='best',
min_samples_split=5,
min_samples_leaf=5,
max_depth = ?,
max_features=None,



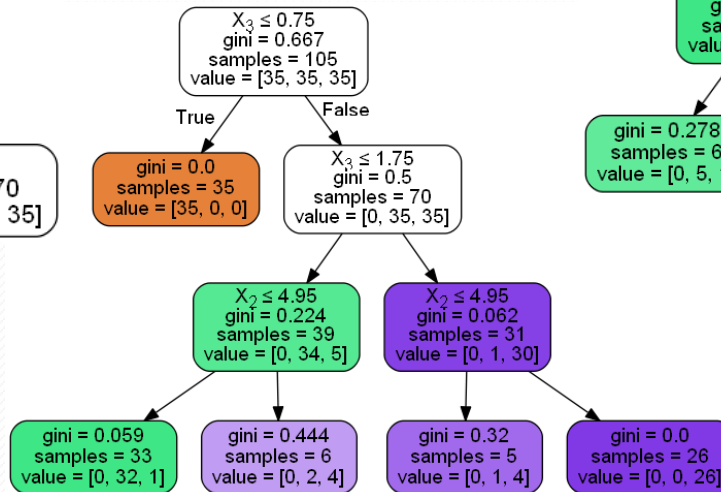
max_depth = 2



max_depth = 20



max_depth = 1

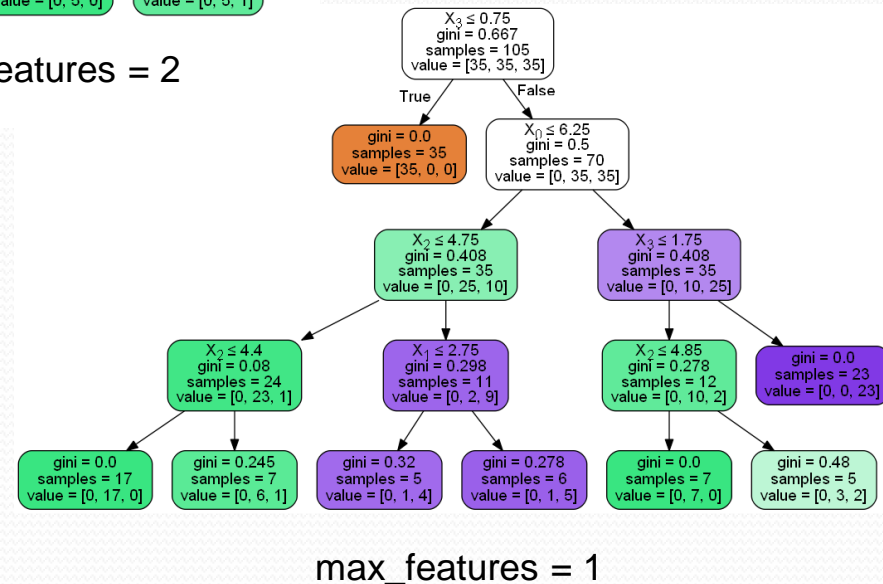
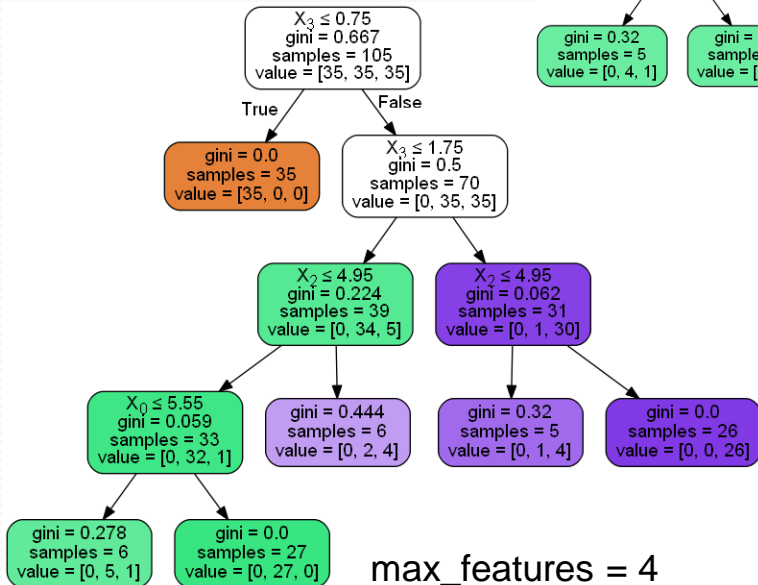
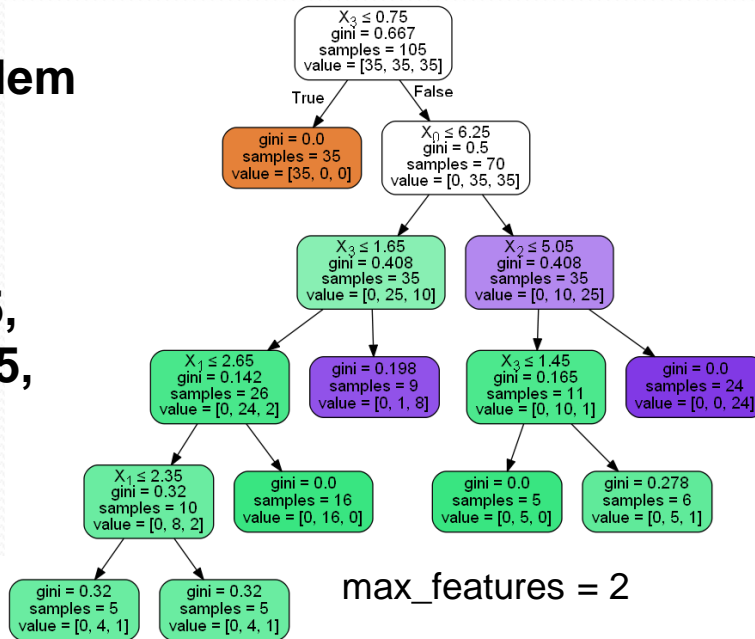


max_depth = 3

Decision Tree

Example 1: Iris problem

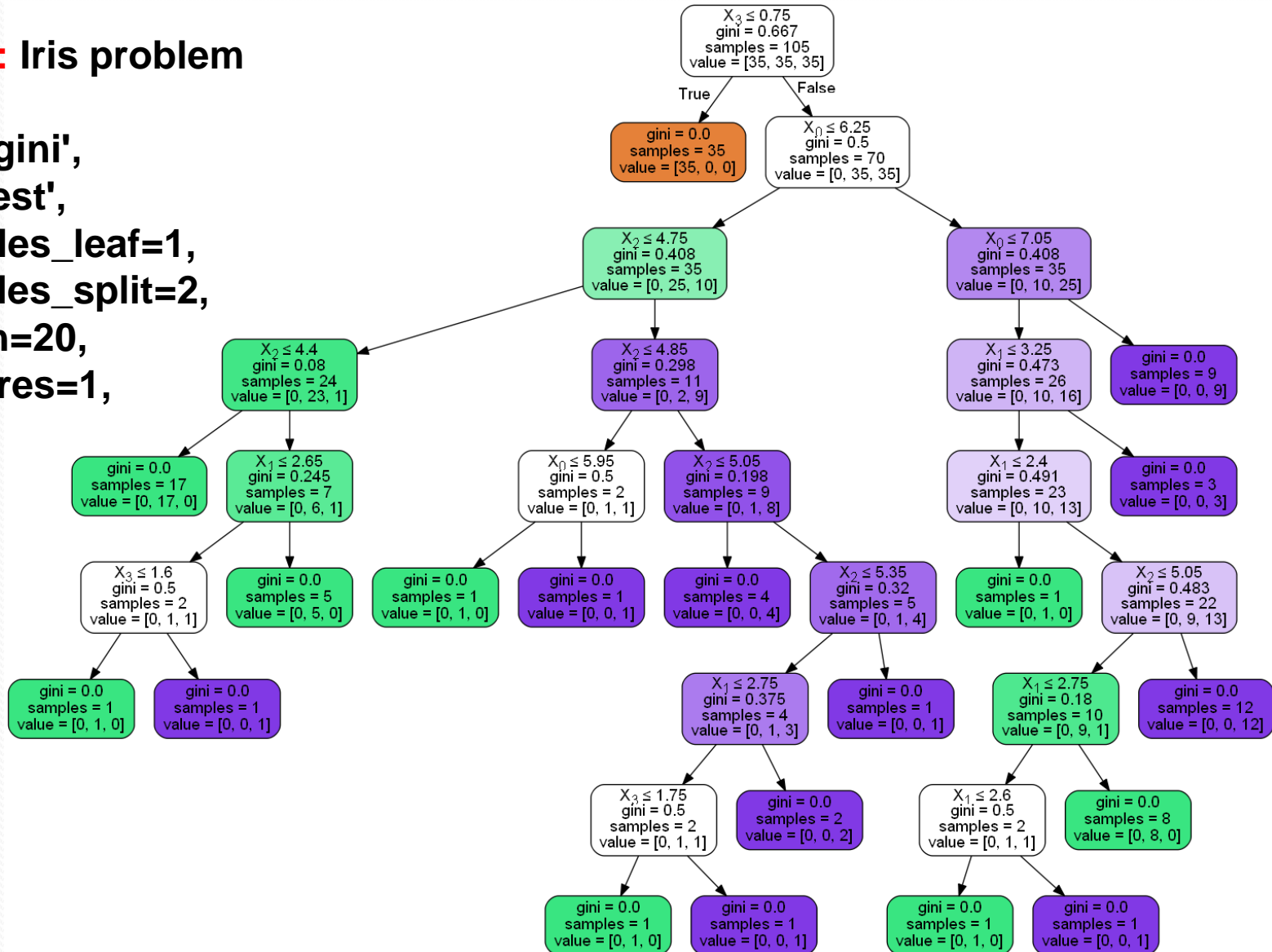
critterion='gini',
splitter='best',
min_samples_leaf=5,
min_samples_split=5,
max_depth=20,
max_features = ?,



Decision Tree

Example 1: Iris problem

**criterion='gini',
splitter='best',
min_samples_leaf=1,
min_samples_split=2,
max_depth=20,
max_features=1,**



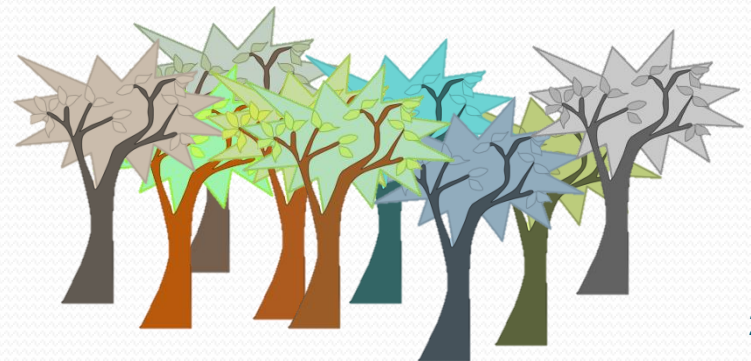
Random Forest Algorithm

The random forest method is another way to elaborate nonlinear problems. Classification and Regression Trees (CART) were first introduced by [Leo Breiman \(2001\)](#) for classification or regression predictive modeling problems.

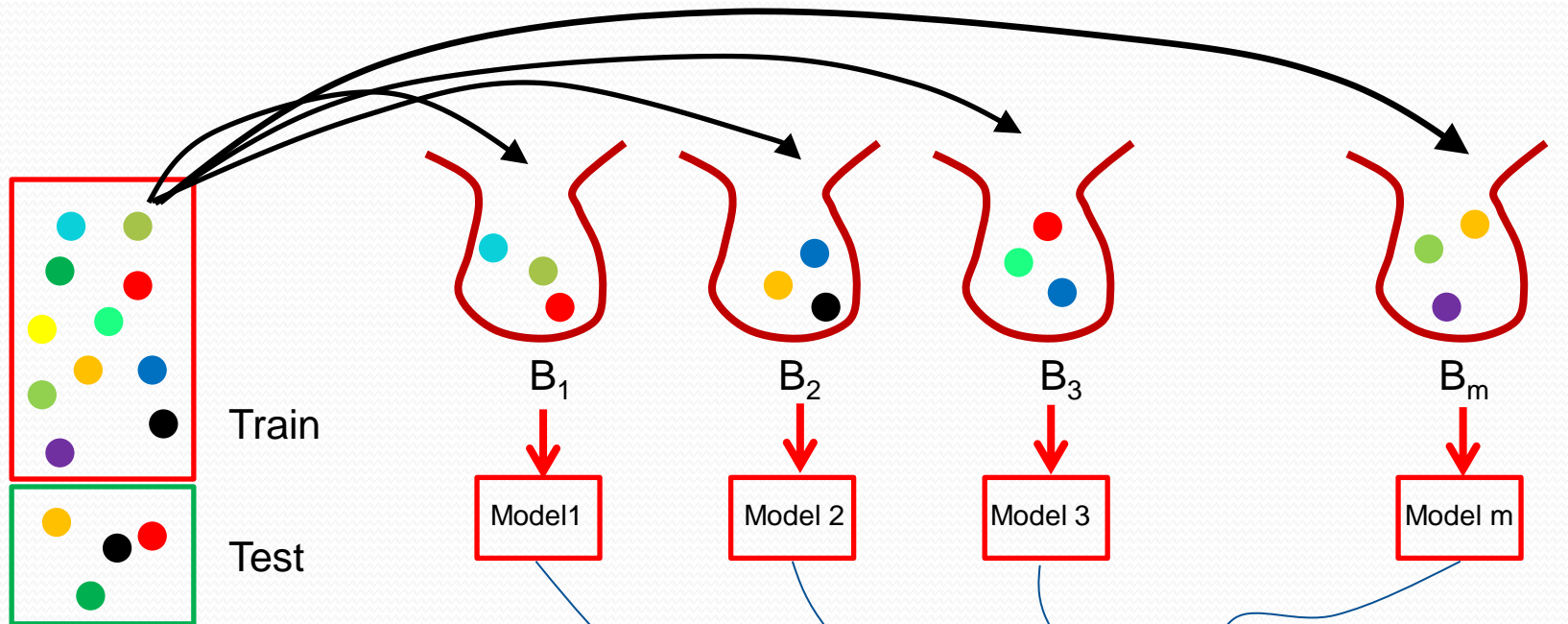
Trees: The trees are grown from different subsets of the training data by a [bagging procedure \(Breiman, 1996\)](#) which ensures the diversity of the trees and minimizes the similarities between them.

The bagging procedure is based on a [bootstrap technique](#). In this method the mean value of a variable is calculated from the average of the mean values of the random subsets of that variable with replacement (i.e., the same value of the variable can be selected few times randomly in a subset). This [reduces](#) the correlation between the trees that may cause [overfitting](#). The prediction for unseen data is then based on the average estimations from all regression or classification trees.

It has been shown that the [generalization error](#) for the forest [converges](#) as the [number of the trees](#) in the forest [increases](#) which prevents an [overfitting](#) problem. On the other hand, a reduction in the number of features in splitting reduces the correlation among trees and increases the model prediction accuracy ([Breiman, 1996](#)).



Bagging – Bootstrap Sampling

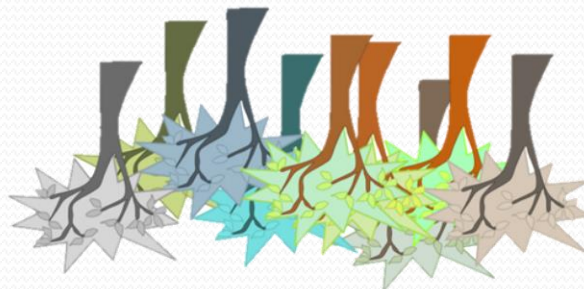


n : Number training samples

n' : Number of samples in a bag

$n' < n$ (e.g., $n' = 60\% n$)

m : Number of bags



Parameters

The random forest method has a number of tuning parameters (e.g., Raschka, 2015) which can be adjusted by a grid search or manually.

n_estimators : integer, optional (default=10)

The number of trees in the forest. Increasing this parameter to a certain level reduces the possibility of overfitting to the cost of computational time. However, a very large number of random trees may not result in a significant gain in the prediction accuracy and the score could even drop.

criterion : string, optional (default="gini")

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

max_depth : integer or None, optional (default=None)

If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples. The depth of the tree is proportional to the number of splits of the nodes into the child nodes. Deeper trees are more complex and are more likely to overfit the training data.

Parameters

min_samples_split: int, float, optional (default=2)

The minimum number of samples required to split at each internal node of a tree (from one to all sample). However, splitting a **larger** number of samples at each node will cause **underfitting** and may abruptly cause a decrease in learning. This parameter may be declared as integer (the number of samples required to split at each node) or be set as the fraction of the samples in the range of $0 < \text{min_samples_split} \leq 1$ (floating).

min_samples_leaf : int, float, optional (default=1)

The minimum number of samples at a leaf node to avoid further splitting (integer or floating). **Increasing** this parameter may cause **underfitting**.

The parameters “**max_depth**” and “**min_samples_leaf**” control the level of regularization; that is, **decreasing** “**max_depth**” and **increasing** “**min_samples_leaf**” will result in a better regularization.

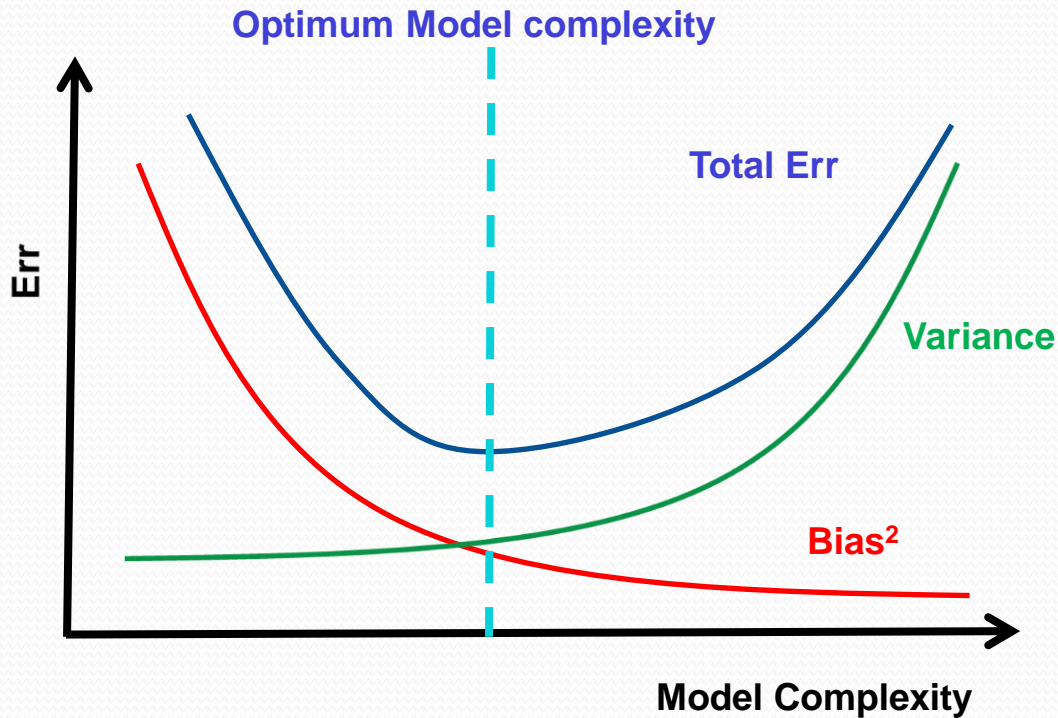
max_features : int, float, string or None, optional (default=”auto”)

The number of features to consider when looking for the best split.

The **smaller values** prevents **overfitting**, but **too small** values may introduce **underfitting**.

max_leaf_nodes: The maximum number of leaves in the tree.

Bias-Variance Trade-off

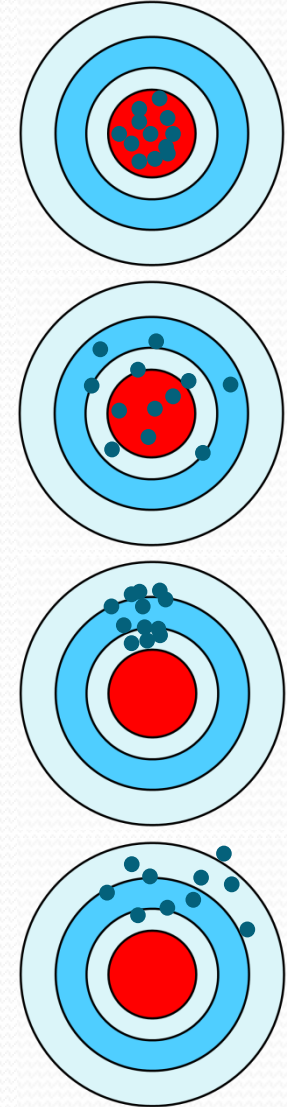


$$Err(x) = (\bar{g}(x) - f(x))^2 + (g(x) - \bar{g}(x))^2 + \sigma^2$$

Bias

Variance

Irreducible error



Random Forest Model

Image Recognition – Fashion Data

Fashion problem:

$$precision = t_p / (t_p + f_p)$$

$$recall = t_p / (t_p + f_n)$$

t_p : true positive

f_p : false positive

f_n : true positive

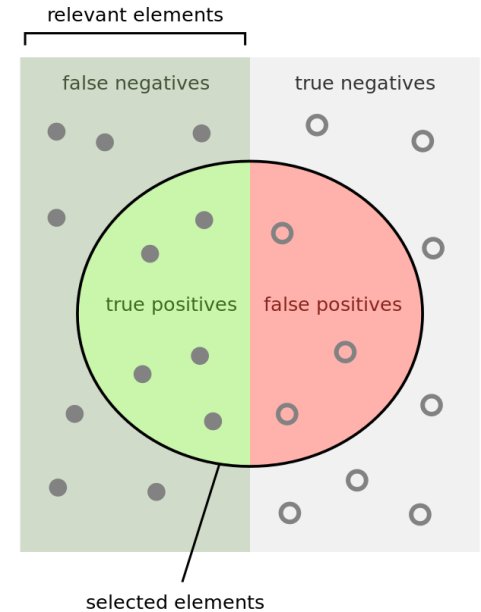
$$f_\beta = (1 + \beta^2) \frac{precision \times recall}{\beta^2 precision + recall}$$

$$f_1 = 2 \times \frac{precision \times recall}{precision + recall}$$

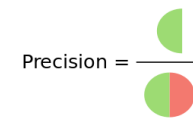
The f_β score can be interpreted as a weighted harmonic mean of the precision and recall, where an f_β score reaches its best value at 1 and worst score at 0.

The support is the number of occurrences of each class in y_true .

```
===== final test
Test data metrics:
  precision  recall  f1-score  support
class 0     0.64   0.81   0.72    91
class 1     0.93   0.91   0.92    92
class 2     0.52   0.75   0.61    91
class 3     0.79   0.79   0.79   105
class 4     0.75   0.54   0.62    99
class 5     0.74   0.74   0.74   105
class 6     0.37   0.28   0.32    99
class 7     0.73   0.76   0.74    94
class 8     0.91   0.83   0.87   115
class 9     0.87   0.84   0.86   109
avg / total 0.73   0.73   0.72  1000
===== final test
```

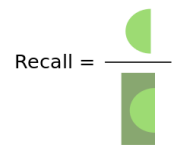


How many selected items are relevant?



Precision =

How many relevant items are selected?



Recall =

Random Forest Model

Image Recognition – Fashion Data

```
1#=====
2#https://github.com/PraveenDubba/Image-Classification-using-Random-Forest/blob/master/Random_Forest_Latest.py
3import sys
4import os
5import pandas as pd
6import numpy as np
7import sklearn
8from sklearn.model_selection import StratifiedKFold
9from sklearn.ensemble import RandomForestClassifier
10from sklearn.grid_search import GridSearchCV
11from sklearn.cluster import KMeans
12from sklearn.cross_validation import train_test_split
13#=====
14
15#===== import data
16train_data = pd.read_csv("fashion_train.csv")
17final_test_data = pd.read_csv("fashion_test.csv")
18#===== import data
19
20#===== split data
21print('===== split data')
22X_train = train_data.iloc[:,0:784]
23y_train = train_data.iloc[:,784:]
24print('train_data.shape = ', train_data.shape)
25print('X_train.shape = ', X_train.shape)
26#print('X_train = ', X_train.tail(5))
27print('y_train = ', y_train.tail(5))
28
29X_final_test = final_test_data.iloc[:,0:784]
30y_final_test = final_test_data.iloc[:,784:]
31print('final_test_data.shape = ', final_test_data.shape)
32print('X_final_test.shape = ', X_final_test.shape)
33#print('X_final_test = ', X_final_test.tail(5))
34print('y_final_test = ', y_final_test.tail(5))
35
36# Splitting train data into training and validation datasets
37x_train_v, x_test_v, y_train_v, y_test_v = train_test_split(X_train,y_train, test_size = 0.3, random_state = 2)
38
```

Random Forest Model

Image Recognition – Fashion Data

```
38
39 print('x_train_v.shape = ', x_train_v.shape)
40 print('x_test_v.shape = ', x_test_v.shape)
41 print('y_train_v.shape = ', y_train_v.shape)
42 print('y_test_v.shape = ', y_test_v.shape)
43 y_train_v = np.ravel(y_train_v)
44 print('y_train_v.shape = ', y_train_v.shape)
45 y_test_v = np.ravel(y_test_v)
46 print('y_test_v.shape = ', y_test_v.shape)
47 print('===== split data')
48 #===== split data
49
50 #===== a simple random forest without grid search
51 # ===== Using Random Forest without hyper parameter tuning and clustering =====
52 print('===== simple model')
53 rf = RandomForestClassifier(n_estimators=6)
54 rf.fit(x_train_v,y_train_v)
55
56 # Predictions on training and validation
57 y_pred_train = rf.predict(x_train_v)
58 # predictions for test
59 y_pred_test = rf.predict(x_test_v)
60
61 my_target_names=['class 0', 'class 1', 'class 2','class 3', 'class 4', 'class 5', 'class 6', 'class 7', 'class 8', 'class 9']
62
63 print('===== metrics 1')
64 print('===== train')
65 # training metrics
66 print("Training metrics:")
67 print(sklearn.metrics.classification_report(y_true= y_train_v, y_pred= y_pred_train, target_names=my_target_names))
68 print('===== train')
69
70 print('===== test')
71 # test data metrics
72 print("Test data metrics:")
73 print(sklearn.metrics.classification_report(y_true= y_test_v, y_pred= y_pred_test, target_names=my_target_names))
74 print('===== test')
75
```

Random Forest Model

Image Recognition – Fashion Data

```
76 print('===== final test')
77 # Predictions on testset
78 y_pred_test = rf.predict(X_final_test)
79 # test data metrics
80 print("Test data metrics:")
81 print(sklearn.metrics.classification_report(y_true= y_final_test, y_pred= y_pred_test, target_names=my_target_names))
82 print('===== final test')
83 print('y_pred_test[113] = ', y_pred_test[113]) #check one of them
84 print('===== metrics 1')
85 print('===== simple model')
86 #===== a simple random forest without grid search
87 #sys.exit()
88 #===== random forest with grid search
89 print('===== model with grid search')
90 print("=====random forest with grid search started=====")
91 #Using Grid Search for hyper parameter tuning
92 #===== grid search
93 print("=====grid search=====")
94 rf = RandomForestClassifier()
95 param_grid = dict(n_estimators=[10,20],
96                  min_samples_leaf=[2,3])
97 grid = GridSearchCV(rf, param_grid, cv=5, scoring=None)
98 model = grid.fit(x_train_v, y_train_v)
99 print('===== grid best estimator')
100 print('grid.best_estimator_ = ', grid.best_estimator_)
101 print('===== grid best estimator')
102 #===== grid results
103 print("=====grid scores=====")
104 print("grid.cv_results_ {}".format(grid.grid_scores_))
105 print("=====grid scores=====")
106 print("=====grid best score=====")
107 print("-grid.best_score_ (r2 / variance) = {} ".format(-grid.best_score_))
108 print("=====grid best score=====")
109 print("=====grid best parameters=====")
110 print("grid.best_params_ = {} ".format(grid.best_params_))
111 print("=====grid best parameters=====")
112 print("=====grid best estimator=====")
113 print("grid.best_estimator_ = {} ".format(grid.best_estimator_))
114 print("=====grid best estimator=====")
115 #===== grid results
```

Random Forest Model

Image Recognition – Fashion Data

```
116
117#===== predict
118y_pred_train = model.predict(x_train_v)
119    # predictions for test
120y_pred_test = model.predict(x_test_v)
121print('===== metrics 2')
122    # training metrics
123print("Training metrics:")
124print(sklearn.metrics.classification_report(y_true= y_train_v, y_pred= y_pred_train))
125
126    # test data metrics
127print("Test data metrics:")
128print(sklearn.metrics.classification_report(y_true= y_test_v, y_pred= y_pred_test))
129
130
131# Predictions on testset
132y_pred_test = model.predict(X_final_test)
133    # test data metrics
134print("Test data metrics:")
135print(sklearn.metrics.classification_report(y_true= y_final_test, y_pred= y_pred_test))
136print('===== metrics 2')
137#===== predict
138print("=====grid search=====")
139#===== grid search
140print('===== model with grid search')
141#===== a simple random forest without grid searc
142sys.exit('Program finished')
143
144
```

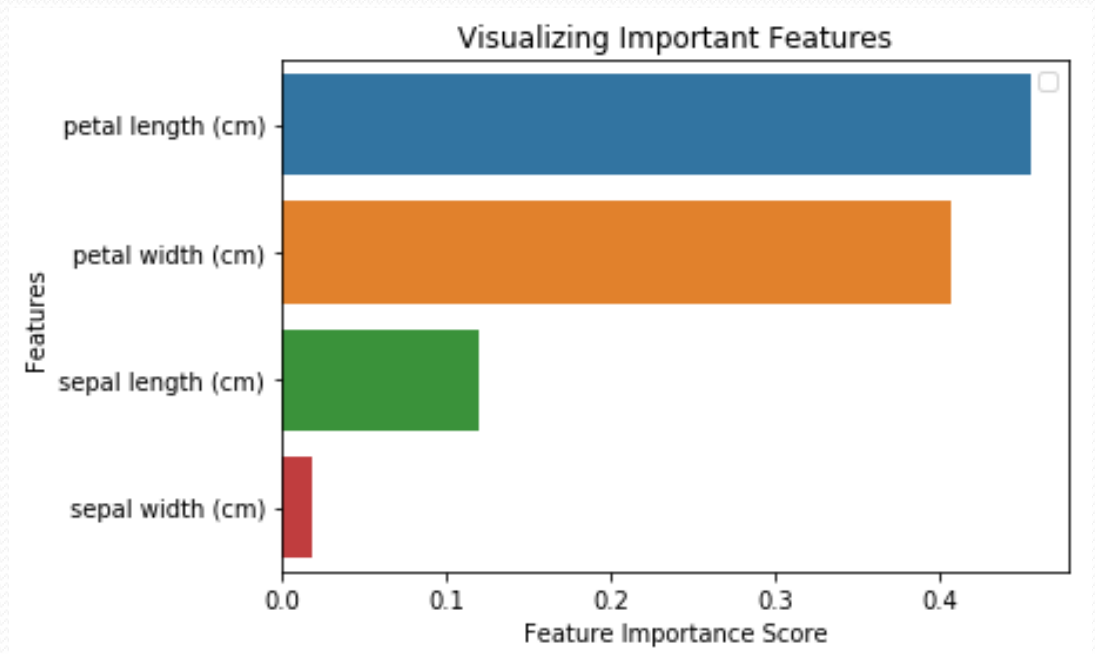
Feature Importance

```
tree.fit(X_train, y_train)  
print('tree.feature_importances_ = ', tree.feature_importances_)
```

```
tree.feature_importances_ = [0.00425693 0.          0.06941788 0.92632519]
```

Or

```
clf.fit(X_train,y_train)  
import pandas as pd  
feature_imp = pd.Series(clf.feature_importances_,index=iris.feature_names).sort_values(ascending=False)
```



Feature Importance

```
1#=====
2from sklearn import datasets # scikit-Learn dataset library
3import sys
4#===== Load iris data
5iris = datasets.load_iris()
6print('iris.target_names = ', iris.target_names) # Labels (setosa, versicolor, virginica)
7print('iris.feature_names = ', iris.feature_names) # feature names
8print('iris.data[0:5] = ', iris.data[0:5]) # print data (features) from row 0 to row 4
9print('iris.target[0:5] = ', iris.target[0:5]) # print target (class labels) from row 0 to row 4
10#print('iris.target = ', iris.target)
11#print('iris.data = ', iris.data)
12print('iris.target.shape = ', iris.target.shape)
13print('iris.data.shape = ', iris.data.shape)
14#===== Load iris data
15
16#===== create a DataFrame of iris dataset
17import pandas as pd
18data=pd.DataFrame({
19    'sepal length':iris.data[:,0],
20    'sepal width':iris.data[:,1],
21    'petal length':iris.data[:,2],
22    'petal width':iris.data[:,3],
23    'species':iris.target
24})
25data.head()
26print('data.head(3) = ', data.head(3))
27print('data.tail(3) = ', data.tail(3))
28#===== create a DataFrame of iris dataset
29
30#===== split data into train and test
31from sklearn.model_selection import train_test_split
32X=data[['sepal length', 'sepal width', 'petal length', 'petal width']] # Features
33y=data['species'] # Labels
34X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state = 19) # 70% training and 30% test
35
36print('=====')
37#print('y = ', y)
38#print('X = ', X)
39print('y_train.shape = ', y_train.shape)
40print('X_train.shape = ', X_train.shape)
```


Feature Importance

```
41 print('y_test.shape = ', y_test.shape)
42 print('X_test.shape = ', X_test.shape)
43 print('=====')
44 #===== split data into train and test
45
46 #===== random forest model (training)
47 from sklearn.ensemble import RandomForestClassifier
48 #clf=RandomForestClassifier(n_estimators=100) # n_estimators: the number of trees in the forest
49
50 clf=RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
51                             max_depth=None, max_features='auto', max_leaf_nodes=None,
52                             min_impurity_decrease=0.0, min_impurity_split=None,
53                             min_samples_leaf=1, min_samples_split=2,
54                             min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=1,
55                             oob_score=False, random_state=None, verbose=0, # set random_state = i and get the ave.
56                             warm_start=False)
57 clf.fit(X_train,y_train) # train
58 #===== random forest model (training)
59
60 #===== testing
61 y_pred=clf.predict(X_test)
62 y_pred2=clf.predict(X_train)
63 print('y_pred = ', y_pred)
64 from sklearn import metrics #scikit-Learn metrics module for accuracy calculation
65 print("Accuracy (test):",metrics.accuracy_score(y_test, y_pred))
66 print("Accuracy (train):",metrics.accuracy_score(y_train, y_pred2))
67 #===== testing
68
69 #===== predict a single item (can be used for plotting decision boundaries)
70 single_pred = clf.predict([[3, 5, 4, 2]])
71 print('single_pred = ', single_pred)
72 #===== predict a single item (can be used for plotting decision boundaries)
73 #sys.exit()
74 #=====feature importance
75 import pandas as pd
76 feature_imp = pd.Series(clf.feature_importances_,index=iris.feature_names).sort_values(ascending=False)
77 print('feature_imp = ', feature_imp)
78 #=====feature importance
79
```

Feature Importance

```
79
80#=====feature importance (visualization)
81import matplotlib.pyplot as plt
82import seaborn as sns # statistical data visualization
83#%matplotlib inline
84# Creating a bar plot
85sns.barplot(x=feature_imp, y=feature_imp.index)
86# Add Labels to your graph
87plt.xlabel('Feature Importance Score')
88plt.ylabel('Features')
89plt.title("Visualizing Important Features")
90plt.legend()
91plt.show()
92#=====feature importance (visualization)
93
94#===== Generating the Model on Selected Features
95#from sklearn.cross_validation import train_test_split
96#from sklearn.ensemble import RandomForestClassifier
97#Create a Gaussian Classifierclf=RandomForestClassifier(n_estimators=100)
98#from sklearn import metrics #scikit-learn metrics module for accuracy calculation
99
100# Split dataset into features and Labels
101X=data[['petal length', 'petal width','sepal length']] # Removed feature "sepal Length"
102y=data['species']
103# Split dataset into training set and test set
104X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.70, random_state=5) # 70% training and 30% test
105
106clf.fit(X_train,y_train) # train
107y_pred=clf.predict(X_test) # test
108y_pred2=clf.predict(X_train)
109print("Accuracy (test with reduced features):",metrics.accuracy_score(y_test, y_pred))
110print("Accuracy (train with reduced features):",metrics.accuracy_score(y_train, y_pred2))
111
112#===== Generating the Model on Selected Features
113
114
```