

# Supervised Learning – Part 4

---

ESM3081 Programming for Data Science

Seokho Kang



# Learning algorithms covered in this course

---

- **Supervised Learning** (Classification/Regression)

- K-Nearest Neighbors
- Linear Models (Logistic/Linear Regression)
- **Decision Trees**
- **Random Forests**
- Support Vector Machines
- Neural Networks

*\* Many algorithms have a classification and a regression variant, and we will describe both.*

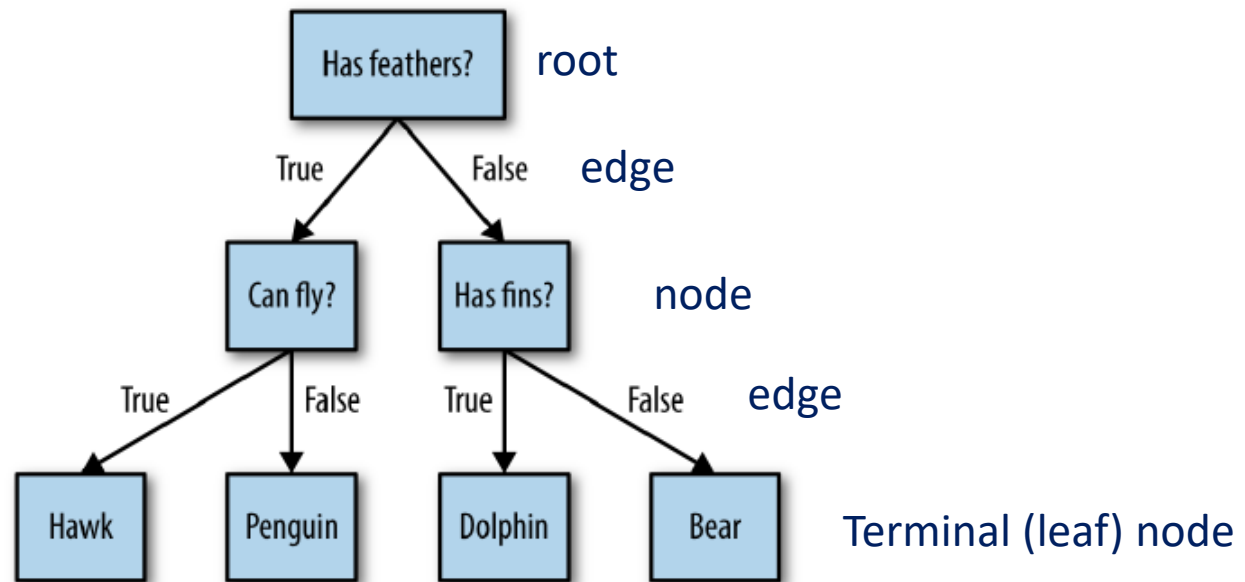
*\* We will review the most popular machine learning algorithms, explain how they learn from data and how they make predictions, and examine the strengths and weaknesses of each algorithm.*

# Decision Trees

---

# Decision Trees

- **A decision tree is a hierarchy of if/else questions leading to a decision.**
  - Each node either represents a question or a terminal node (also called a leaf) that contains the answer.
  - The edges connect the answers to a question with the next question you would ask.



# Decision Trees

- **Example – Skiing Classification**

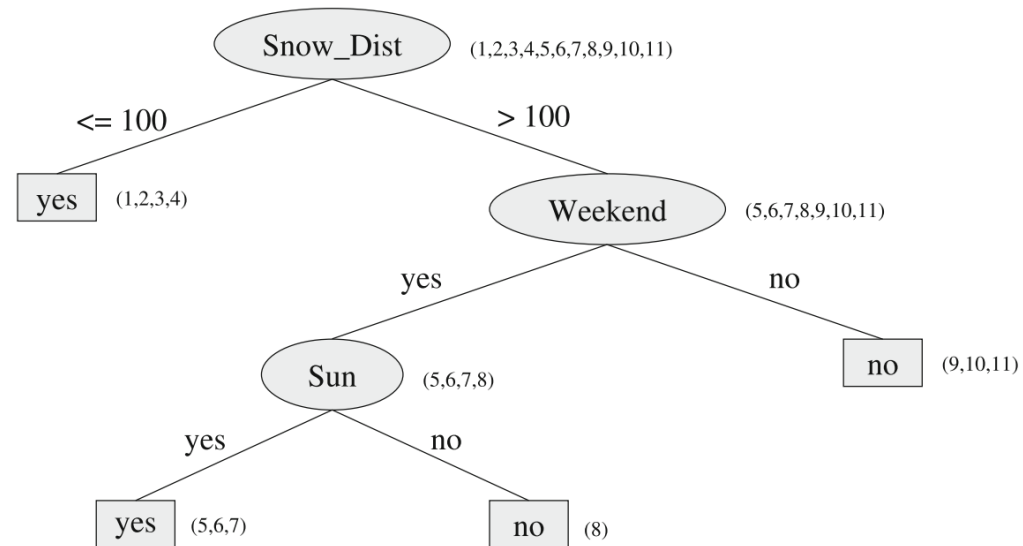
- **Features:** “*Snow\_Dist*”, “*Weekend*”, and “*Sun*”.
- **Target:** “*Skiing*” = *yes* or *no*

**Dataset**

Day	<i>Snow_Dist</i>	<i>Weekend</i>	<i>Sun</i>	<i>Skiing</i>
1	$\leq 100$	yes	yes	yes
2	$\leq 100$	yes	yes	yes
3	$\leq 100$	yes	no	yes
4	$\leq 100$	no	yes	yes
5	$> 100$	yes	yes	yes
6	$> 100$	yes	yes	yes
7	$> 100$	yes	yes	no
8	$> 100$	yes	no	no
9	$> 100$	no	yes	no
10	$> 100$	no	yes	no
11	$> 100$	no	no	no

<i>Ski</i> (goal variable)	yes, no	Should I drive to the nearest ski resort with enough snow?
<i>Sun</i> (feature)	yes, no	Is there sunshine today?
<i>Snow_Dist</i> (feature)	$\leq 100$ , $> 100$	Distance to the nearest ski resort with good snow conditions (over/under 100 km)
<i>Weekend</i> (feature)	yes, no	Is it the weekend today?

**Decision Tree**



# Decision Trees

---

- **Instead of building the tree by hand, we can learn them from data**
  - the output is a set of if/else rules represented by a tree diagram
    - Each inner node represents a feature.
    - Each edge stands for the condition of a feature value.
    - At each leaf node, a prediction is given.
- **The extracted knowledge can be easily understood, interpreted, and controlled by humans in the form of a readable decision tree**
- **Decision trees are widely used for classification and regression tasks**

# Decision Trees

---

- **There are finitely many different decision trees.**
  - **Optimal algorithm** for the construction of a tree is to simply generate all possible trees and choose the best one.
  - The obvious disadvantage of this algorithm is its unacceptably high computation time, as soon as the number of features becomes somewhat larger.
- **Decision Tree Learning**
  - We use **heuristic algorithms** with greedy strategy.
  - Because greedy strategy is used for construction of the tree, the trees are in general sub-optimal.

# Decision Trees

---

- Given a **(training)** dataset  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$  such that  $\mathbf{x}_i = (x_{i1}, \dots, x_{id}) \in \mathbb{R}^d$  is the  $i$ -th input vector of  $d$  features and  $y_i$  is the corresponding target label.
- **General Procedure** – Repeatedly split a node into two parts so as to minimize the impurity of outcome within the new parts.
  - The training dataset  $D$  constitutes the root node
  - Repeat the following process
    - Try all possible splits in all nodes and features to find the “**best split**”
    - Split the node



# Decision Trees

- **How to determine the best split (for classification)**
  - Nodes with low degree of impurity are preferred

C0: 5
C1: 5

High degree of impurity

C0: 9
C1: 1

Low degree of impurity

- **Measures of node impurity**
  - \*  $p(j|t)$  is the fraction of class  $j$  at node  $t$

$$\text{GINI}(t) = 1 - \sum_j [p(j|t)]^2$$

$$\text{Entropy}(t) = - \sum_j p(j|t) \log p(j|t)$$

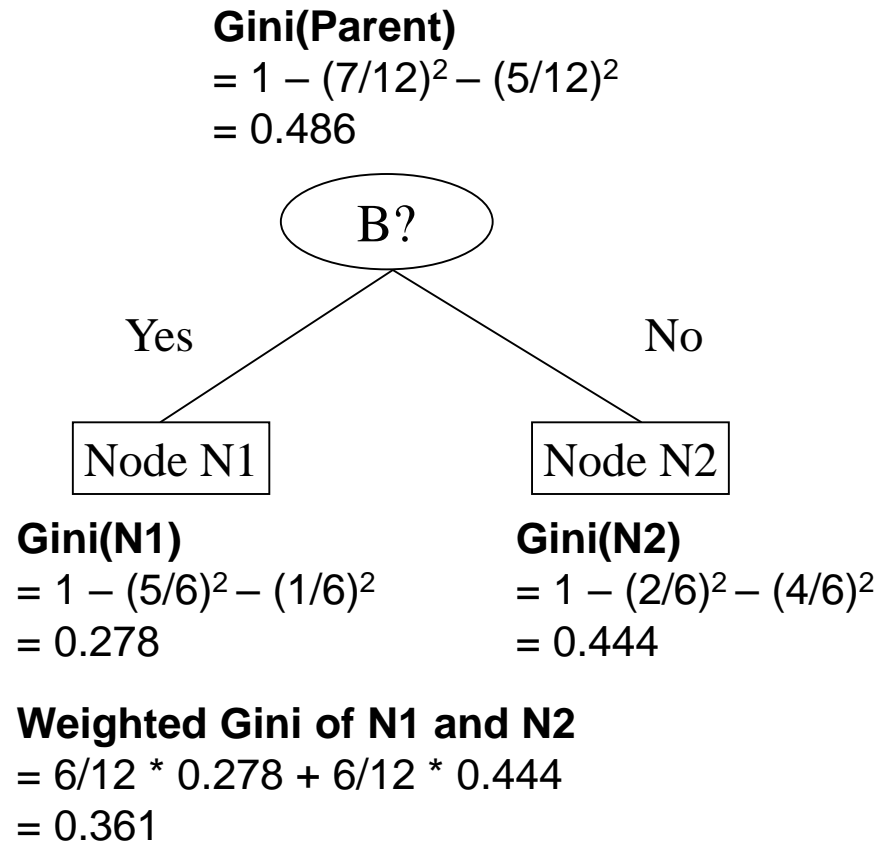
# Decision Trees

---

- **How to determine the best split (for classification)**
  1. **Compute impurity measure (P) before splitting**
  2. **Compute impurity measure (M) after splitting**
    - Compute impurity measure of each child node
    - M is the weighted impurity of children
  3. **Choose the feature that produces the highest gain**  
or equivalently, lowest impurity measure after splitting (M)
    - **Gain = P - M**

# Decision Trees

- How to determine the best split (for classification)



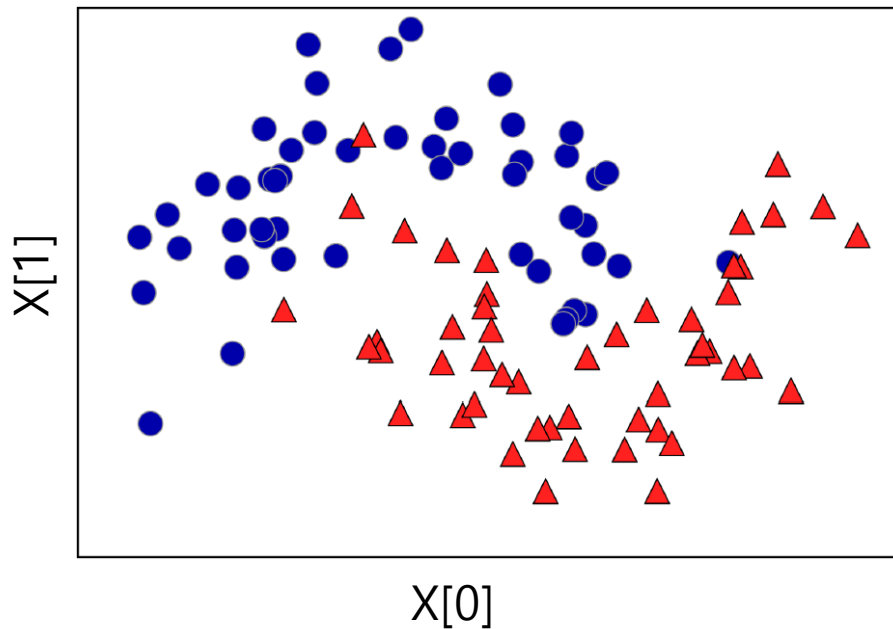
	Parent
C1	7
C2	5
<b>Gini = 0.486</b>	

	N1	N2
C1	5	2
C2	1	4
<b>Gini=0.361</b>		

**Gain = 0.486 - 0.361 = 0.125**

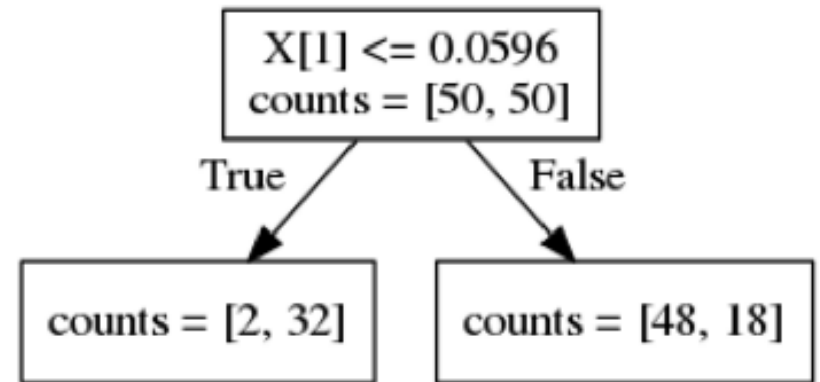
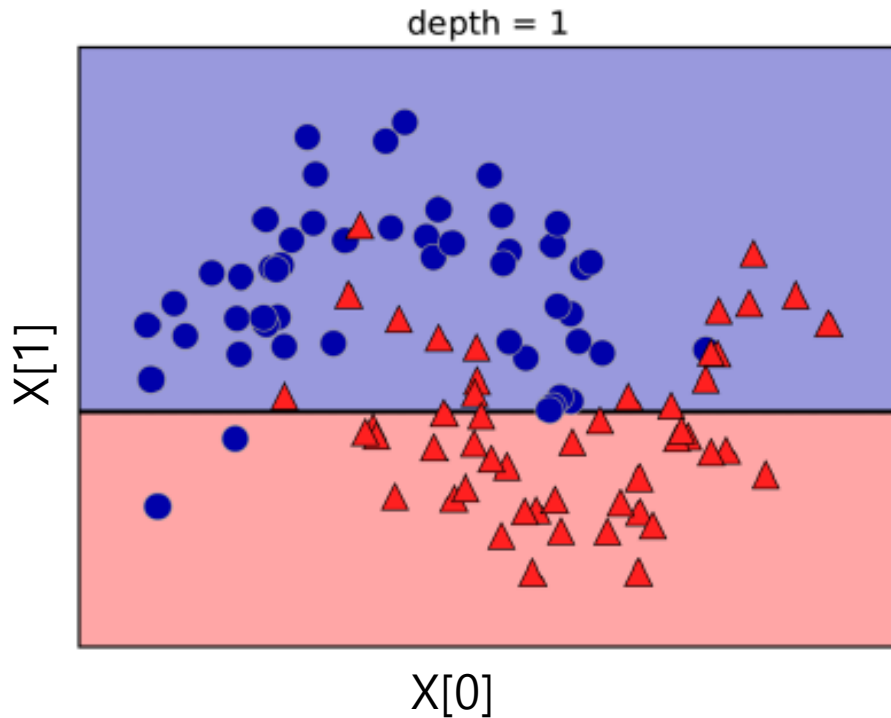
# Decision Trees

- **Example of building a decision tree (for classification)**
  - The *two-moons* dataset consists of two half-moon shapes, with each class consisting of 50 data points.



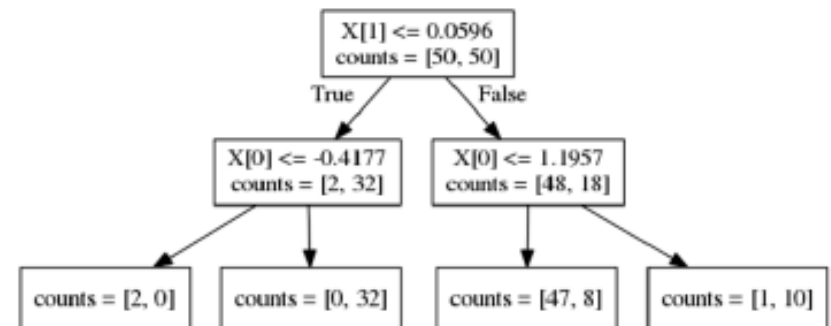
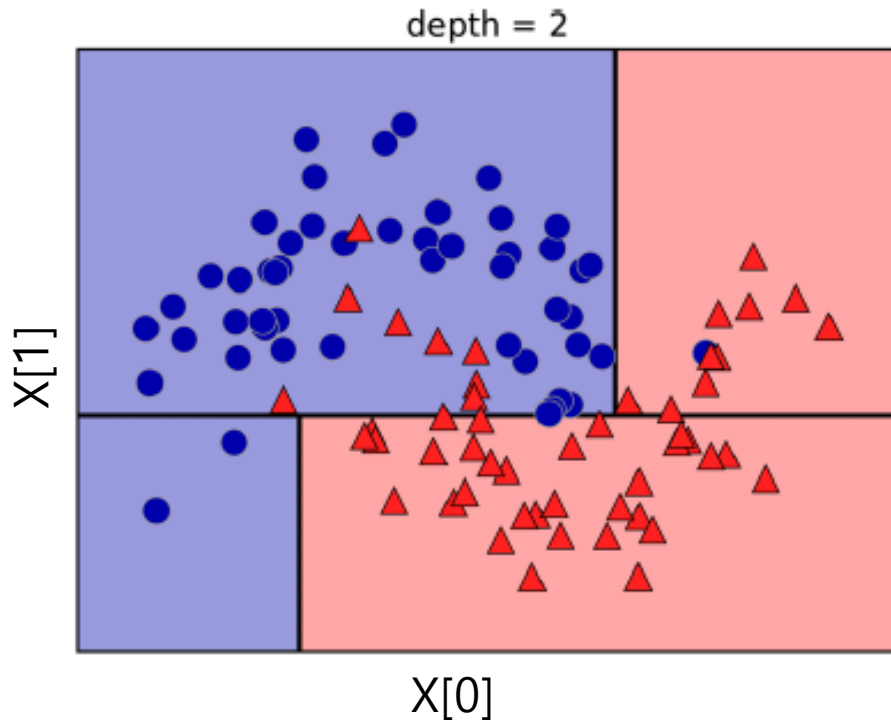
# Decision Trees

- **Example of building a decision tree (for classification)**
  - The algorithm searches over all possible tests and finds the one that is most informative about the target.



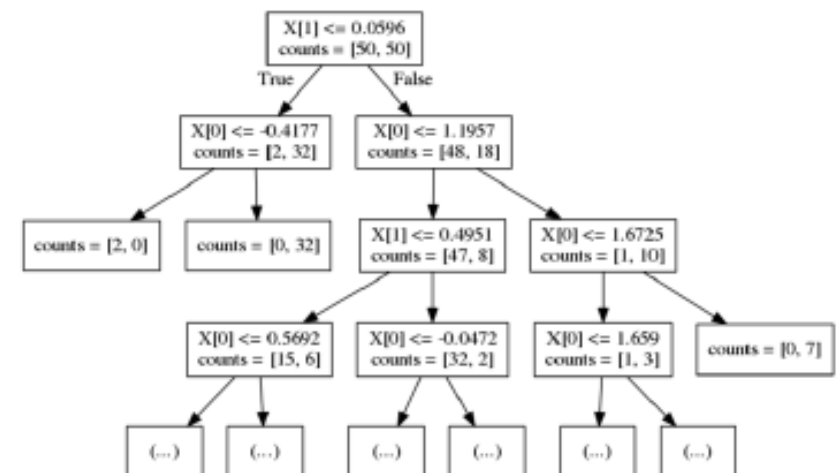
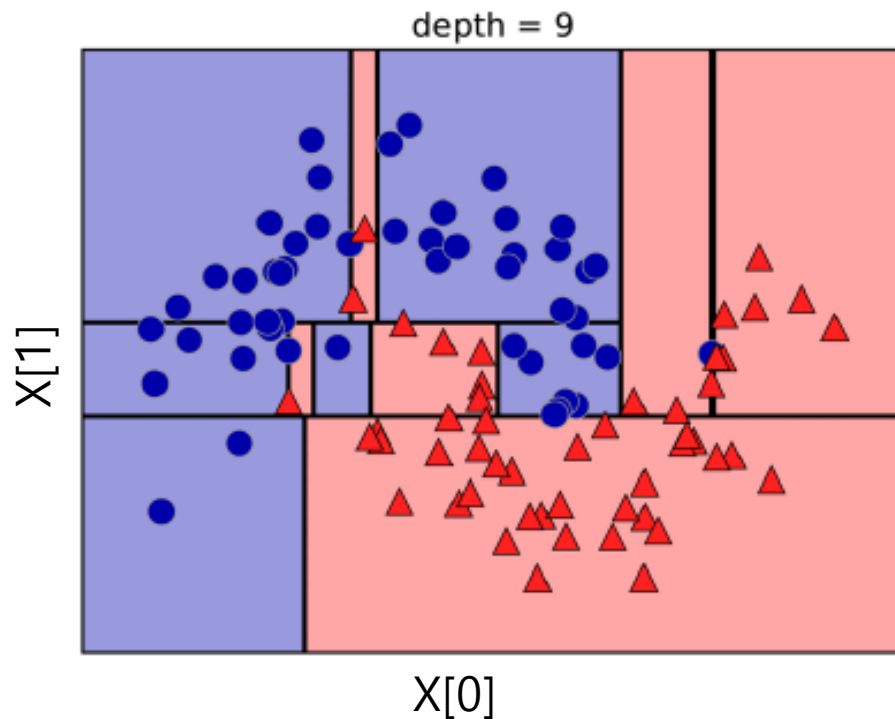
# Decision Trees

- **Example of building a decision tree (for classification)**
  - We can build a more accurate model by repeating the process of looking for the most informative next split for the left and the right region.



# Decision Trees

- **Example of building a decision tree (for classification)**
  - The recursive partitioning of the data is repeated until each region in the partition (each leaf in the tree) become homogeneous.



# Decision Trees

---

- To make a prediction for a new data point, we traverse the tree to find the leaf the data point falls into.
- **For Classification,**
  - The output for the data point is the majority class of the training points in this leaf.
- **For Regression**
  - The output for the data point is the mean target of the training points in this leaf.



# Decision Trees

---

- **Building a tree typically leads to a model that is very complex and highly overfit to the training data.**
  - Continuing until all leaves are pure means that a tree is 100% accurate on the training set, but fails to generalize on new data
- **Two common strategies to prevent overfitting**
  - **Pre-pruning:** stopping the creation of the tree early (*e.g.*, below)
    - *max\_depth*: limiting the maximum depth of the tree
    - *min\_sample\_leaf*: requiring a minimum number of points in a node to keep splitting it
    - *max\_leaf\_nodes*: limiting the maximum number of leaves
  - **Post-pruning:** building the tree but then removing or collapsing nodes that contain little information

# scikit-learn Practice: *DecisionTreeClassifier*

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

## `sklearn.tree.DecisionTreeClassifier`

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2,  
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,  
min_impurity_decrease=0.0, min_impurity_split=None, class_weight=None, ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

*\* It only implements pre-pruning, not post-pruning.*

Read more in the [User Guide](#).

### Parameters:

**criterion** : {"gini", "entropy"}, 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** : {"best", "random"}, 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, 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.

# scikit-learn Practice: *DecisionTreeClassifier*

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

## **`min_samples_leaf` : *int or float, default=1***

The minimum number of samples required to be at a leaf node. A split point at any depth will only be considered if it leaves at least `min_samples_leaf` training samples in each of the left and right branches. This may have the effect of smoothing the model, especially in regression.

- If int, then consider `min_samples_leaf` as the minimum number.
- If float, then `min_samples_leaf` is a fraction and `ceil(min_samples_leaf * n_samples)` are the minimum number of samples for each node.

*Changed in version 0.18:* Added float values for fractions.

## **`max_leaf_nodes` : *int, 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.

## scikit-learn Practice: *DecisionTreeClassifier*

- **Example (*breast\_cancer* dataset)**

```
In [2]: from sklearn.datasets import load_breast_cancer  
cancer = load_breast_cancer()
```

```
In [3]: from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(  
    cancer.data, cancer.target, stratify=cancer.target, random_state=42)
```

```
In [4]: from sklearn.tree import DecisionTreeClassifier  
clf = DecisionTreeClassifier(random_state=0)  
clf.fit(X_train, y_train)
```

```
Out[4]: ▾ DecisionTreeClassifier  
DecisionTreeClassifier(random_state=0)
```

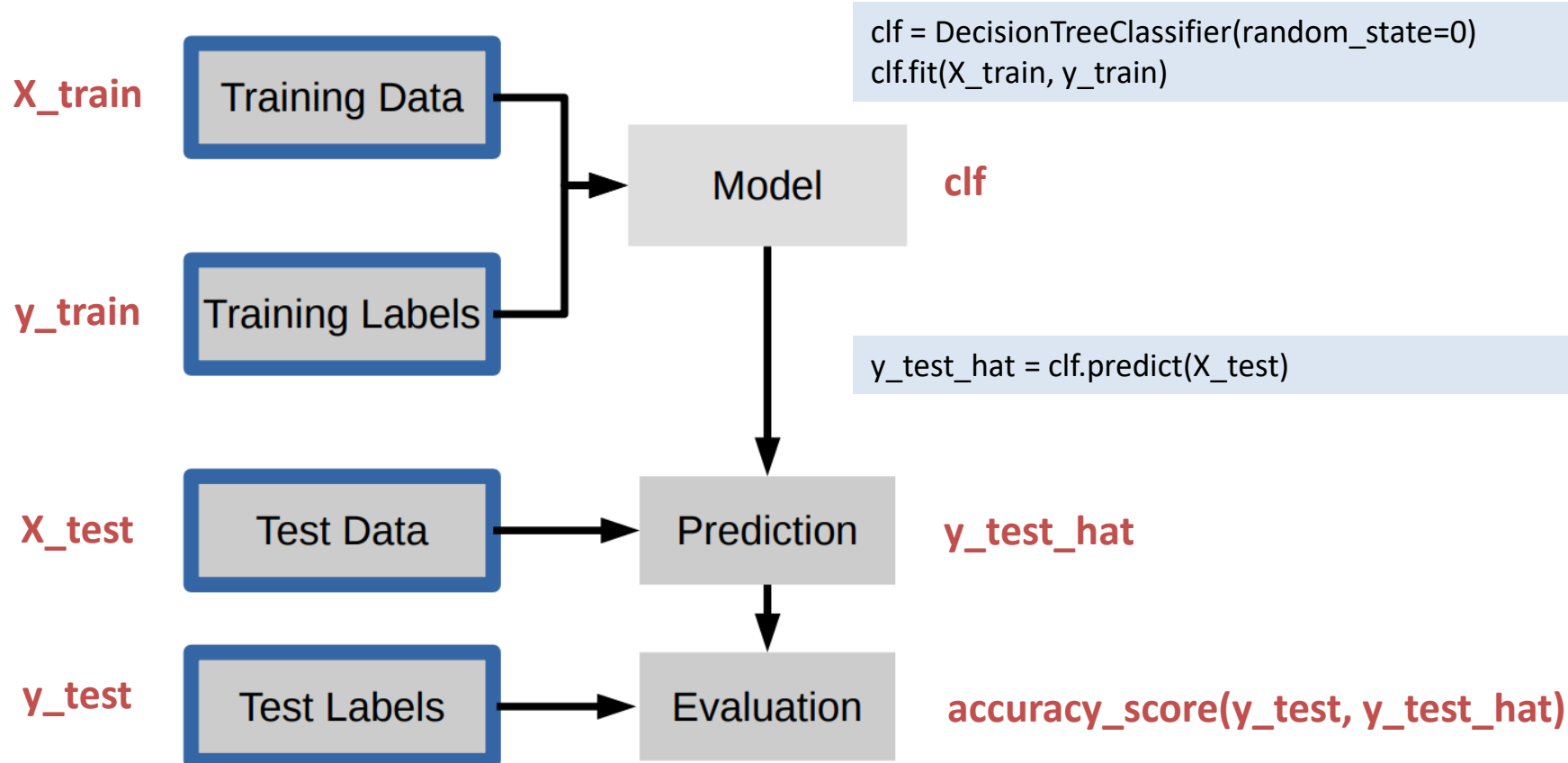
```
In [5]: from sklearn.metrics import accuracy_score  
y_train_hat = clf.predict(X_train)  
print('train accuracy: %.5f'%accuracy_score(y_train, y_train_hat))  
y_test_hat = clf.predict(X_test)  
print('test accuracy: %.5f'%accuracy_score(y_test, y_test_hat))
```

```
train accuracy: 1.00000  
test accuracy: 0.93706
```

# scikit-learn Practice: *DecisionTreeClassifier*

- Example (*breast\_cancer* dataset)

```
cancer = load_breast_cancer()  
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target, stratify=cancer.target, random_state=42)
```



## scikit-learn Practice: *DecisionTreeClassifier*

- Example (*breast\_cancer* dataset): varying the hyperparameter *min\_samples\_leaf*

```
In [6]: from sklearn.datasets import load_breast_cancer
        from sklearn.model_selection import train_test_split
        from sklearn.tree import DecisionTreeClassifier
        from sklearn.metrics import accuracy_score

        cancer = load_breast_cancer()
        X_train, X_test, y_train, y_test = train_test_split(
            cancer.data, cancer.target, stratify=cancer.target, random_state=42)
```

```
In [7]: training_accuracy = []
        test_accuracy = []

        m_settings = [1, 2, 5, 7, 10, 20]
        for m in m_settings:
            # build the model
            clf = DecisionTreeClassifier(min_samples_leaf=m, random_state=0)
            clf.fit(X_train, y_train)

            # accuracy on the training set
            y_train_hat = clf.predict(X_train)
            training_accuracy.append(accuracy_score(y_train, y_train_hat))

            # accuracy on the test set (generalization)
            y_test_hat = clf.predict(X_test)
            test_accuracy.append(accuracy_score(y_test, y_test_hat))
```

	min_samples_leaf	training accuracy	test accuracy
0	1	1.00000	0.93706
1	2	0.99061	0.93706
2	5	0.97653	0.93706
3	7	0.96244	0.95105
4	10	0.96244	0.95105
5	20	0.94601	0.91608

## Visualizing and Analyzing Decision Trees

- The visualization of the tree provides a great in-depth view of how the algorithm makes predictions

```
In [9]: from sklearn import tree
         tree.plot_tree(clf, precision=5)
```

```
Out[9]: [Text(0.5833333333333334, 0.9, 'X[20] <= 16.795#ngini = 0.46786#nsamples = 426#nvalue = [159, 267]'),
Text(0.4166666666666667, 0.7, 'X[27] <= 0.13595#ngini = 0.16056#nsamples = 284#nvalue = [25, 259]'),
Text(0.3333333333333333, 0.5, 'X[1] <= 21.56#ngini = 0.03124#nsamples = 252#nvalue = [4, 248]'),
Text(0.16666666666666666, 0.3, 'X[10] <= 0.40995#ngini = 0.00948#nsamples = 210#nvalue = [1, 209]'),
Text(0.08333333333333333, 0.1, 'gini = 0.0#nsamples = 190#nvalue = [0, 190]'),
Text(0.25, 0.1, 'gini = 0.095#nsamples = 20#nvalue = [1, 19]'),
Text(0.5, 0.3, 'X[18] <= 0.01888#ngini = 0.13265#nsamples = 42#nvalue = [3, 39]'),
Text(0.4166666666666667, 0.1, 'gini = 0.255#nsamples = 20#nvalue = [3, 17]'),
Text(0.5833333333333334, 0.1, 'gini = 0.0#nsamples = 22#nvalue = [0, 22]'),
Text(0.5, 0.5, 'gini = 0.45117#nsamples = 32#nvalue = [21, 11]'),
Text(0.75, 0.7, 'X[21] <= 23.375#ngini = 0.10633#nsamples = 142#nvalue = [134, 8]'),
Text(0.6666666666666666, 0.5, 'gini = 0.455#nsamples = 20#nvalue = [13, 7]'),
Text(0.8333333333333334, 0.5, 'X[5] <= 0.09352#ngini = 0.01626#nsamples = 122#nvalue = [121, 1]'),
Text(0.75, 0.3, 'gini = 0.095#nsamples = 20#nvalue = [19, 1]'),
Text(0.9166666666666666, 0.3, 'gini = 0.0#nsamples = 102#nvalue = [102, 0]')]
```



# Feature Importance in Decision Trees

- ***feature importance*** summarizes the workings of a tree by rating how important each feature is for the decision the tree makes.
  - The importance of a feature is computed as the (normalized) total reduction of the criterion brought by that feature.
  - It is a number between 0 and 1 for each feature, where 0 means “not used at all” and 1 means “perfectly predicts the target.”

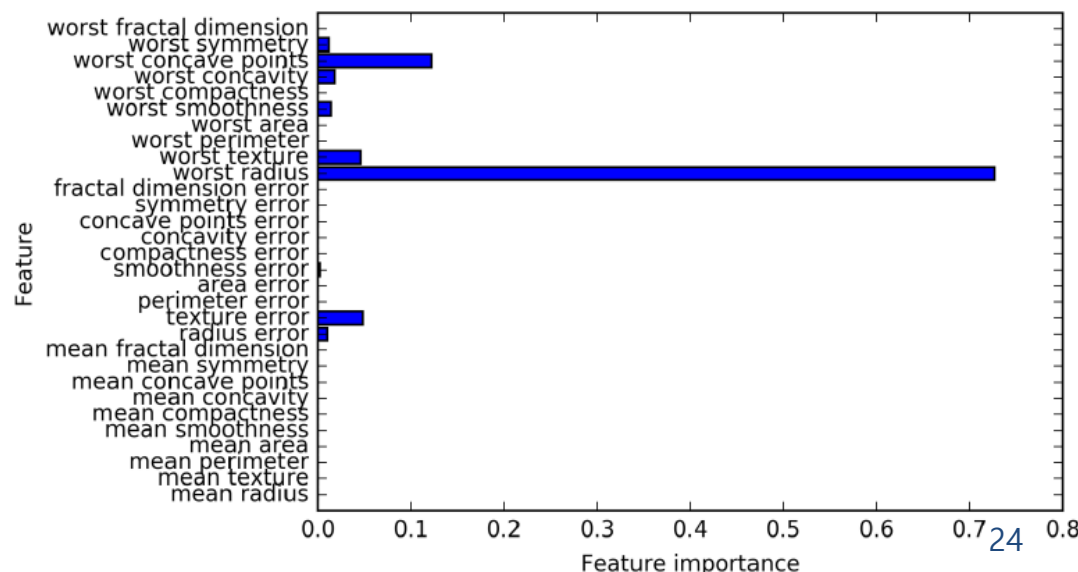
Attributes:

`feature_importances_` : ndarray of shape (n\_features,)

Return the feature importances.

```
In [10]: clf.feature_importances_
```

```
Out[10]: array([0.00000e+00, 1.86437e-03, 0.00000e+00, 0.00000e+00, 0.00000e+00,  
 5.01021e-04, 0.00000e+00, 0.00000e+00, 0.00000e+00, 0.00000e+00,  
 5.42188e-04, 0.00000e+00, 0.00000e+00, 0.00000e+00, 0.00000e+00,  
 0.00000e+00, 0.00000e+00, 0.00000e+00, 2.82508e-03, 0.00000e+00,  
 8.30651e-01, 2.40602e-02, 0.00000e+00, 0.00000e+00, 0.00000e+00,  
 0.00000e+00, 0.00000e+00, 1.39556e-01, 0.00000e+00, 0.00000e+00])
```





# scikit-learn Practice: *DecisionTreeRegressor*

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html>

## sklearn.tree.DecisionTreeRegressor

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='mse', splitter='best', max_depth=None, min_samples_split=2,  
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,  
min_impurity_decrease=0.0, min_impurity_split=None, ccp_alpha=0.0)
```

[source]

A decision tree regressor.

Read more in the [User Guide](#).

### Parameters:

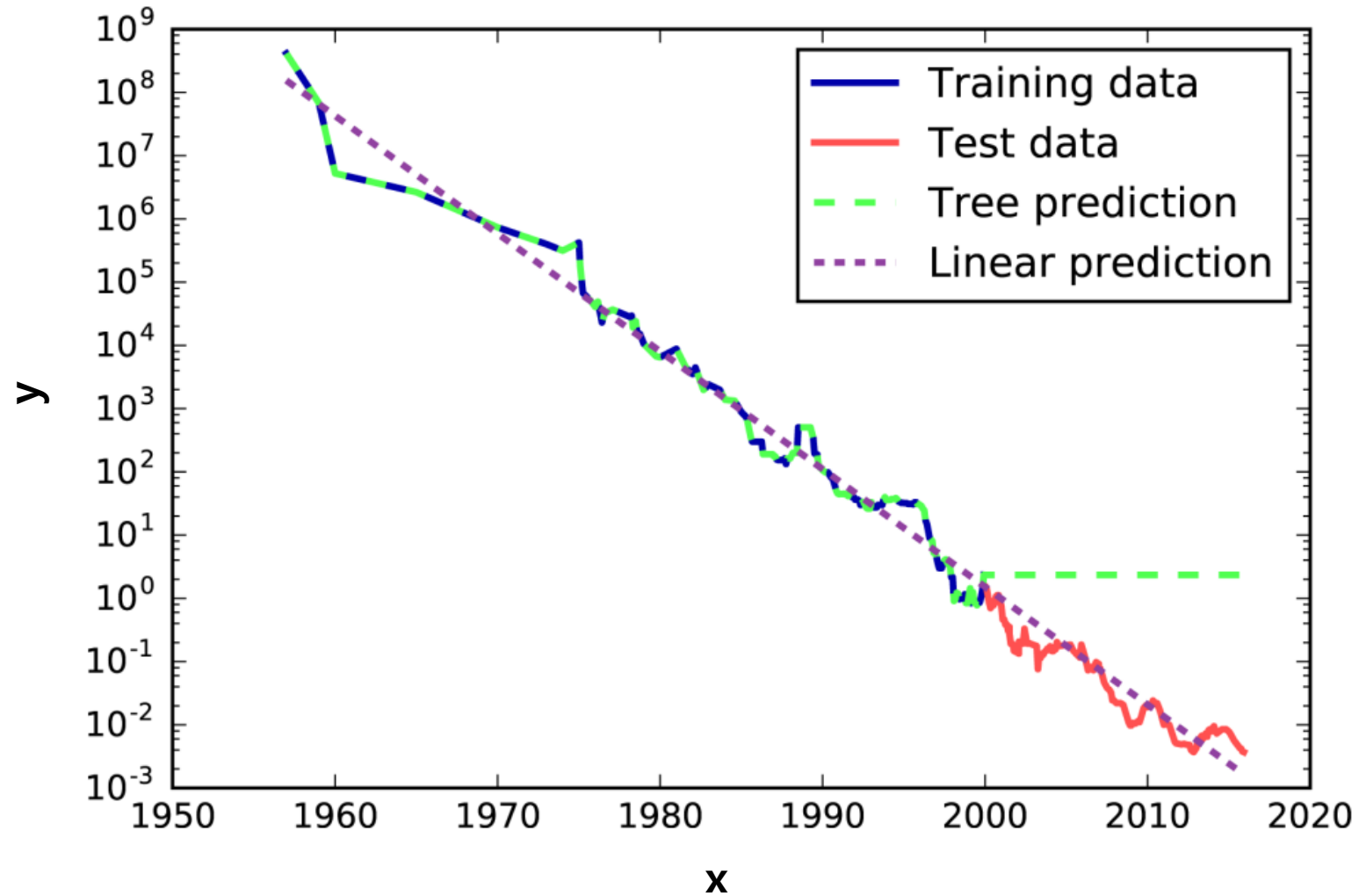
**criterion** : {"mse", "friedman\_mse", "mae", "poisson"}, default="mse"

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 using the mean of each terminal node, "friedman\_mse", which uses mean squared error with Friedman's improvement score for potential splits, "mae" for the mean absolute error, which minimizes the L1 loss using the median of each terminal node, and "poisson" which uses reduction in Poisson deviance to find splits.

*New in version 0.18:* Mean Absolute Error (MAE) criterion.

*New in version 0.24:* Poisson deviance criterion.

## scikit-learn Practice: *DecisionTreeRegressor*



# Discussion

---

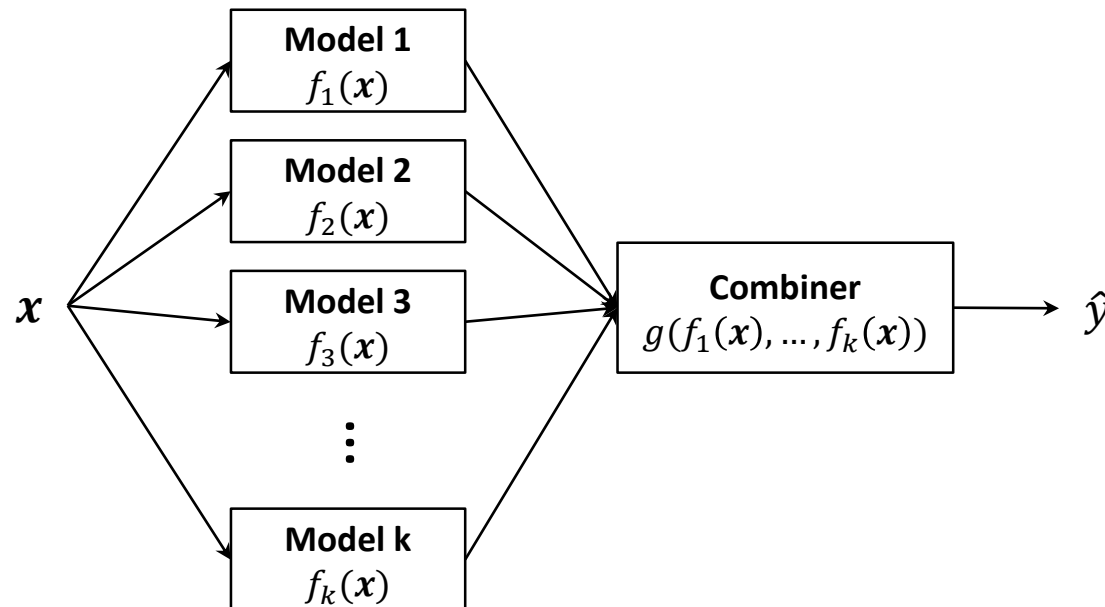
- **The main hyperparameters of decision trees**
  - *criterion*: the function to measure the quality of a split (impurity measure)
  - Picking one of the pre-pruning strategies—setting either *max\_depth*, *max\_leaf\_nodes*, or *min\_samples\_leaf*—is sufficient to prevent overfitting.
  - \* Typically chosen to have the highest performance in *validation data*
- **Strengths**
  - Decision trees tend to work well when you have a mix of continuous and categorical features.
  - The algorithm is completely invariant to the scales of features (*no data scaling is needed*)
  - Feature selection & reduction is automatic.
  - It is robust to outliers.
  - The resulting model can easily be visualized and understood.
- **Weaknesses**
  - Even with the use of pre-pruning, they tend to overfit and provide poor generalization performance.  
→ the ensemble methods are usually used in place of a single decision tree.
  - It does not take into account interactions between features.
  - Space of possible decision trees is exponentially large. Greedy approaches are often unable to find the optimal tree.

# Random Forests

---

# Ensembles of Decision Trees

- Decision trees tend to overfit the training data  
→ Ensembles are one way to address this problem.
- **Ensemble methods** combine multiple machine learning models to create more powerful models.



# Random Forests

---

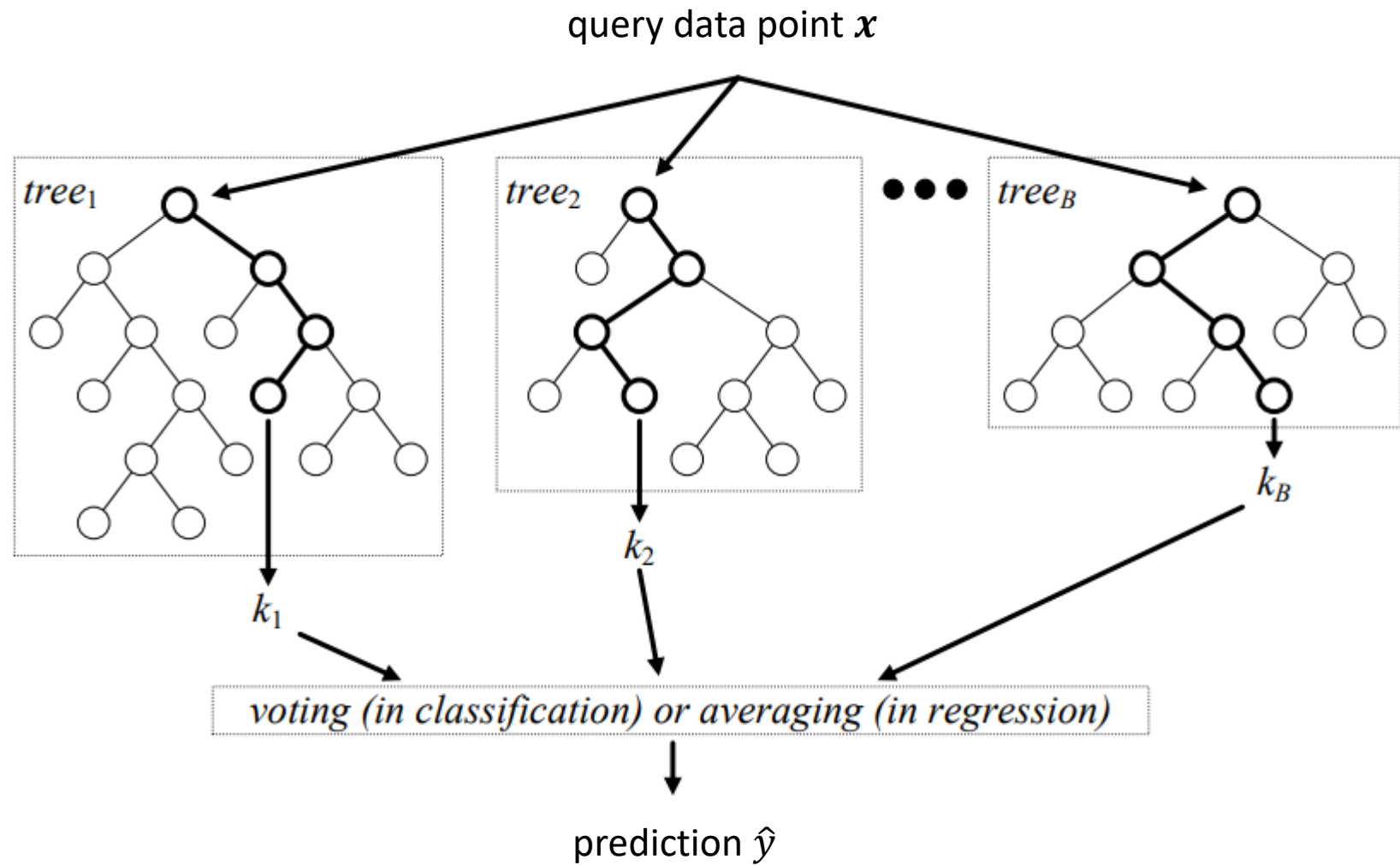
- **A random forest is an ensemble of decision trees, where each tree is slightly different from the others.**
  - Each tree might do a relatively good job of predicting, but will likely overfit on part of the data in different ways.
  - If we build many trees, we can reduce the amount of overfitting by averaging their results while retaining the predictive power of the trees.
  - Random forests get their name from injecting randomness into the tree building to ensure each tree is different.

# Random Forests

---

- Two ways in which the trees in a random forest are randomized
    - **by selecting the data points used to build a tree**
      - *bootstrap*: It leads to each decision tree in the random forest being built on a slightly different dataset.
        - From a list ['a', 'b', 'c', 'd'], possible examples of bootstrap samples are ['b', 'd', 'd', 'c'] and ['d', 'a', 'd', 'a'].
    - **by selecting the features in each split test.**
      - *max\_features*: in each node, the algorithm randomly selects a subset of the features, and it looks for the best possible test involving one of these features
        - each node in a tree can make a decision using a different subset of the features.
        - A high *max\_features* means that the trees in the random forest will be quite similar, and they will be able to fit the data easily, using the most distinctive features.
        - A low *max\_features* means that the trees in the random forest will be quite different, and that each tree might need to be very deep in order to fit the data well.
- \* Typically, for a classification problem with  $p$  features,  $\sqrt{p}$  (rounded down) features are used in each split. For regression problems the inventors recommend  $p/3$  (rounded down) as the default. In practice the best value will depend on the problem.*

# Random Forests





# Random Forests

---

- To make a prediction for a new data point, we first make a prediction for every tree in the forest.
- **For Classification,**
  - Each tree makes a “soft” prediction, providing a probability for each possible output label.
  - The probabilities predicted by all the trees are averaged.
  - The output for the data point is the class with the highest average probability.
- **For Regression**
  - The output for the data point is the mean prediction of the trees in the forest.

# scikit-learn Practice: *RandomForestClassifier*

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

## `sklearn.ensemble.RandomForestClassifier`

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

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

### Parameters:

**`n_estimators` : int, default=100**

The number of trees in the forest.

*Changed in version 0.22:* The default value of `n_estimators` changed from 10 to 100 in 0.22.

# scikit-learn Practice: *RandomForestClassifier*

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

**`max_features : {"auto", "sqrt", "log2"}, int or float, default="auto"`**

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

- If int, then consider `max_features` features at each split.
- If float, then `max_features` is a fraction and `round(max_features * n_features)` features are considered at each split.
- If "auto", then `max_features=sqrt(n_features)`.
- If "sqrt", then `max_features=sqrt(n_features)` (same as "auto").
- If "log2", then `max_features=log2(n_features)`.
- If None, then `max_features=n_features`.

Note: the search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than `max_features` features.

**`bootstrap : bool, default=True`**

Whether bootstrap samples are used when building trees. If False, the whole dataset is used to build each tree.

# scikit-learn Practice: *RandomForestClassifier*

- Example (*breast\_cancer* dataset): varying the hyperparameter *n\_estimators*

```
In [11]: from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

cancer = load_breast_cancer()
X_train, X_test, y_train, y_test = train_test_split(
    cancer.data, cancer.target, stratify=cancer.target, random_state=42)
```

```
In [12]: training_accuracy = []
test_accuracy = []

n_settings = [1, 2, 5, 10, 20, 50, 100]
for n in n_settings:
    # build the model
    clf = RandomForestClassifier(n_estimators=n)
    clf.fit(X_train, y_train)

    # accuracy on the training set
    y_train_hat = clf.predict(X_train)
    training_accuracy.append(accuracy_score(y_train, y_train_hat))

    # accuracy on the test set (generalization)
    y_test_hat = clf.predict(X_test)
    test_accuracy.append(accuracy_score(y_test, y_test_hat))
```

	n_estimators	training accuracy	test accuracy
0	1	0.98357	0.93706
1	2	0.97887	0.91608
2	5	0.98826	0.93706
3	10	0.99765	0.95105
4	20	1.00000	0.95105
5	50	1.00000	0.95804
6	100	1.00000	0.95804

# How Many Trees in a Random Forest?



[International Workshop on Machine Learning and Data Mining in Pattern Recognition](#)

MLDM 2012: [Machine Learning and Data Mining in Pattern Recognition](#) pp 154-168 | [Cite as](#)

## How Many Trees in a Random Forest?

Authors

[Authors and affiliations](#)

Thais Mayumi Oshiro, Pedro Santoro Perez, José Augusto Baranauskas

Conference paper

80

Citations

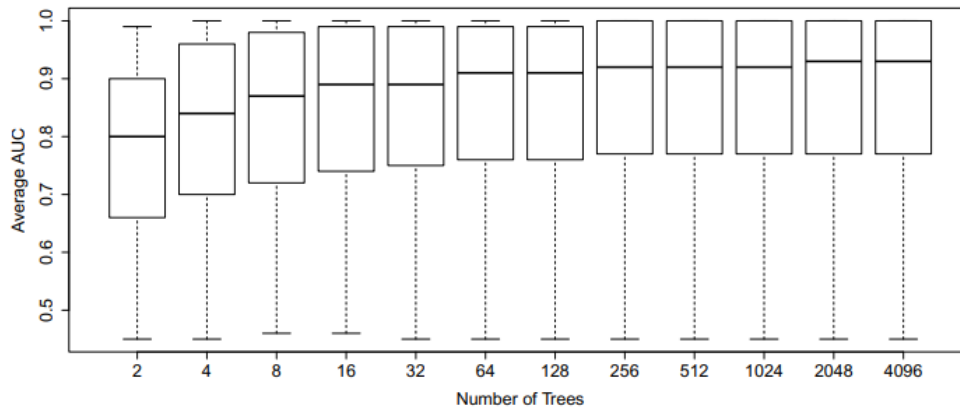
2

Readers

5.5k

Downloads

Part of the [Lecture Notes in Computer Science](#) book series (LNCS, volume 7376)



**Fig. 1.** AUC in all datasets

... .. The analysis of 29 datasets shows that from 128 trees there is no more significant difference between the forests using 256, 512, 1024, 2048 and 4096 trees. The mean and the median AUC values do not present major changes from 64 trees. Therefore, it is possible to suggest, based on the experiments, a range between 64 and 128 trees in a forest. With these numbers of trees it is possible to obtain a good balance between AUC, processing time, and memory usage. ... ..

# Feature Importance in Random Forest

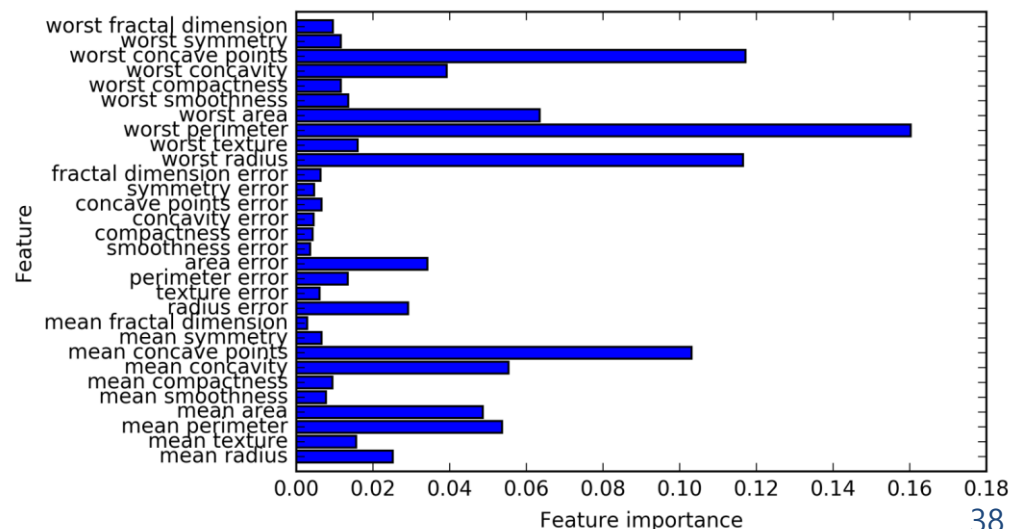
- Similarly to the decision tree, the random forest provides feature importances
  - Computed by aggregating the feature importances over the trees in the forest.
  - Typically, the feature importances provided by the random forest are more reliable than the ones provided by a single tree.

```
In [14]: clf = RandomForestClassifier()  
         clf.fit(X_train, y_train)
```

```
Out[14]: ▾ RandomForestClassifier  
         RandomForestClassifier()
```

```
In [15]: clf.feature_importances_
```

```
Out[15]: array([0.03884, 0.01009, 0.04455, 0.02885, 0.00516, 0.01385, 0.06065,  
                0.07779, 0.004 , 0.00228, 0.00507, 0.00663, 0.0177 , 0.02541,  
                0.00482, 0.00413, 0.00298, 0.00738, 0.00495, 0.00633, 0.12762,  
                0.01444, 0.11571, 0.14592, 0.01136, 0.01029, 0.03425, 0.15056,  
                0.01104, 0.00735])
```



# scikit-learn Practice: *RandomForestRegressor*

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

## sklearn.ensemble.RandomForestRegressor

```
class sklearn.ensemble.RandomForestRegressor(n_estimators=100, *, criterion='mse', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest regressor.

A random forest is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

### Parameters:

**n\_estimators : int, default=100**

The number of trees in the forest.

Changed in version 0.22: The default value of `n_estimators` changed from 10 to 100 in 0.22.

**max\_features : {"auto", "sqrt", "log2"}, int or float, default="auto"**

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

- If "auto", then `max_features=n_features`.

# Discussion

---

- **The main hyperparameters of random forests**
  - *n\_estimators*: larger is always better
  - *bootstrap*: whether bootstrap samples are used
  - *max\_features*: how random each tree is, the default values in scikit-learn are  $\sqrt{n\_features}$  for classification and  $n\_features$  for regression
  - Plus, the main hyperparameters of decision trees
  - \* Typically chosen to have the highest performance in validation data
- **Strengths**
  - They are very powerful, often work well without heavy tuning of the hyperparameters
- **Weaknesses**
  - It is basically impossible to interpret all the trees in detail
  - Building random forests on large datasets might be somewhat time consuming
    - but, it can be parallelized across multiple CPU (use *n\_jobs*)



## Do we Need Hundreds of Classifiers to Solve Real World Classification Problems?

Manuel Fernández-Delgado

Eva Cernadas

Senén Barro

*CITIUS: Centro de Investigación en Tecnologías da Información da USC*

*University of Santiago de Compostela*

*Campus Vida, 15872, Santiago de Compostela, Spain*

MANUEL.FERNANDEZ.DELGADO@USC.ES

EVA.CERNADAS@USC.ES

SENEN.BARRO@USC.ES

Dinani Amorim

*Departamento de Tecnologia e Ciências Sociais- DTCS*

*Universidade do Estado da Bahia*

*Av. Edgard Chastinet S/N - São Geraldo - Juazeiro-BA, CEP: 48.305-680, Brasil*

DINANIAMORIM@GMAIL.COM

### Abstract

**Editor:** Russ Greiner

We evaluate **179 classifiers** arising from **17 families** (discriminant analysis, Bayesian, neural networks, support vector machines, decision trees, rule-based classifiers, boosting, bagging, stacking, random forests and other ensembles, generalized linear models, nearest-neighbors, partial least squares and principal component regression, logistic and multinomial regression, multiple adaptive regression splines and other methods), implemented in Weka, R (with and without the caret package), C and Matlab, including all the relevant classifiers available today. We use **121 data sets**, which represent **the whole UCI** data base (excluding the large-scale problems) and other own real problems, in order to achieve significant conclusions about the classifier behavior, not dependent on the data set collection. **The classifiers most likely to be the bests are the random forest (RF)** versions, the best of which (implemented in R and accessed via caret) achieves 94.1% of the maximum accuracy overcoming 90% in the 84.3% of the data sets. However, the difference is not statistically significant with the second best, the SVM with Gaussian kernel implemented in C using LibSVM, which achieves 92.3% of the maximum accuracy. A few models are clearly better than the remaining ones: random forest, SVM with Gaussian and polynomial kernels, extreme learning machine with Gaussian kernel, C5.0 and avNNet (a committee of multi-layer perceptrons implemented in R with the caret package). The random forest is clearly the best family of classifiers (3 out of 5 bests classifiers are RF), followed by SVM (4 classifiers in the top-10), neural networks and boosting ensembles (5 and 3 members in the top-20, respectively).

## Are Random Forests Truly the Best Classifiers?

**Michael Wainberg**

*Department of Electrical and Computer Engineering  
University of Toronto, Toronto, ON M5S 3G4, Canada;  
Deep Genomics, Toronto, ON M5G 1L7, Canada*

M.WAINBERG@UTORONTO.CA

**Babak Alipanahi**

*Department of Electrical and Computer Engineering  
University of Toronto, Toronto, ON M5S 3G4, Canada*

BABAK@PSI.TORONTO.EDU

**Brendan J. Frey**

*Department of Electrical and Computer Engineering  
University of Toronto, Toronto, ON M5S 3G4, Canada;  
Deep Genomics, Toronto, ON M5G 1L7, Canada*

FREY@PSI.TORONTO.EDU

**Editor:** Nando de Freitas

## Abstract

The JMLR study *Do we need hundreds of classifiers to solve real world classification problems?* benchmarks 179 classifiers in 17 families on 121 data sets from the UCI repository and claims that “the random forest is clearly the best family of classifier”. In this response, we show that the study’s results are biased by the lack of a held-out test set and the exclusion of trials with errors. Further, the study’s own statistical tests indicate that random forests do not have significantly higher percent accuracy than support vector machines and neural networks, calling into question the conclusion that random forests are the best classifiers.

