

# Module-4.5-Decision Trees and Ensemble Models

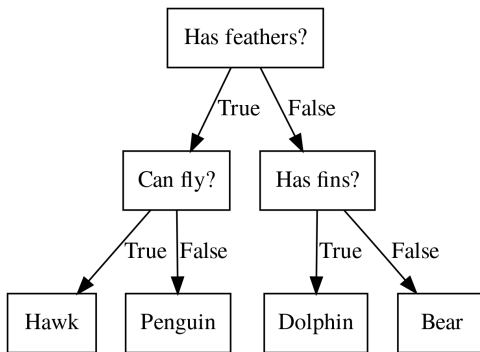
Presented by Yasin Ceran

# Table of Contents

## 1 Decision Trees

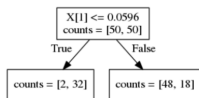
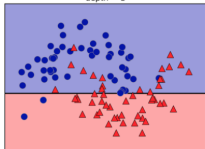
## 2 Ensemble Models

# Idea: Series of Binary of Binary Questions

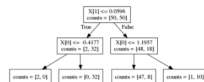
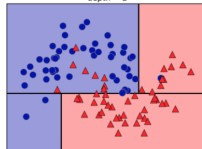


# Building Trees

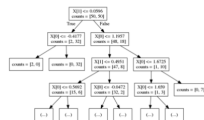
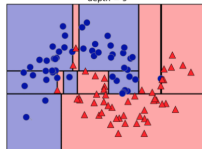
depth = 1



depth = 2



depth = 9



Continuous features:

- “questions” are thresholds on single features.
- Minimize impurity

# Criteria for Classification

- Gini Index:

$$H_{\text{gini}}(X_m) = \sum_{k \in \mathcal{Y}} p_{mk}(1 - p_{mk})$$

- Cross-Entropy:

$$H_{\text{CE}}(X_m) = - \sum_{k \in \mathcal{Y}} p_{mk} \log(p_{mk})$$

- $X_m$  observations in node  $m$
- $\mathcal{Y}$  classes
- $p_m$ . distribution over classes in node  $m$

# An Example

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

# Entropy

Entropy  $H(\mathbf{S})$  is a measure of the amount of uncertainty in the data set  $\mathbf{S}$  (i.e., entropy characterizes the data set  $\mathbf{S}$ ).

$$H(S) = \sum_{c \in C} -p(c) \log_2 p(c)$$

where

- $S$  - The current data set for which the entropy is being calculated
- $C$  - Set of classes in  $S$ , for example,  $C = \{yes, no\}$
- $p(c)$  - The proportion of the number of elements in class  $c$  to the number of elements in set  $S$

When  $H(S) = 0$ , the set  $S$  is perfectly classified.

# Information Gain

Information gain **IG(A)** is the measure of the difference in entropy from before to after the set **S** is split on an attribute **A**. In other words, how much uncertainty in **S** was reduced after splitting **S** on attribute **A**.

$$IG(A, S) = \mathbf{H}(S) - \sum_{t \in T} p(t)H(t)$$

where

- $H(S)$  - Entropy of set **S**
- $T$  - The subsets created from splitting set **S** by attribute **A** such that  $\mathbf{S} = \bigcup_{t \in T} t$
- $p(t)$  - The proportion of the number of elements in  $t$  to the number of elements in set  $S$
- $H(t)$  - Entropy of subset  $t$



# Compute the Entropy for the Weather Data Set

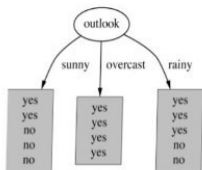
$$H(S) = \sum_{c \in C} -p(c) \log_2 p(c)$$

$$C = \{yes, no\}$$

Out of 14 instances, 9 are classified as 'yes' and 5 as 'no'

- $p(yes) = -(9/14) * \log_2(9/14)$
- $p(no) = -(5/14) * \log_2(5/14)$
- $H(S) = p(yes) + p(no) = 0.94$

# Calculate Entropy and IG for All Features



$$E(\text{Outlook}=\text{sunny}) = -\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) = 0.971$$

$$E(\text{Outlook}=\text{overcast}) = -1 \log(1) - 0 \log(0) = 0$$

$$E(\text{Outlook}=\text{rainy}) = -\frac{3}{5} \log\left(\frac{3}{5}\right) - \frac{2}{5} \log\left(\frac{2}{5}\right) = 0.971$$

$$\left. \begin{array}{l} E(\text{Outlook}=\text{sunny}) \\ E(\text{Outlook}=\text{overcast}) \\ E(\text{Outlook}=\text{rainy}) \end{array} \right\} H(S, \text{Outlook})$$

Average Entropy information for Outlook

$$I(\text{Outlook}) = \frac{5}{14} * 0.971 + \frac{4}{14} * 0 + \frac{5}{14} * 0.971 = 0.693$$

$$\text{Gain}(\text{Outlook}) = E(S) - I(\text{outlook}) = 0.94 - 0.693 = 0.247$$

$$\left. \begin{array}{l} I(\text{Outlook}) \end{array} \right\} \sum_{t \in T} p(t) H(t)$$



$$IG(A, S) = H(S) - \sum_{t \in T} p(t) H(t)$$

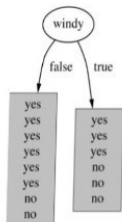
$$E(\text{Windy}=\text{false}) = -\frac{6}{8} \log\left(\frac{6}{8}\right) - \frac{2}{8} \log\left(\frac{2}{8}\right) = 0.811$$

$$E(\text{Windy}=\text{true}) = -\frac{3}{6} \log\left(\frac{3}{6}\right) - \frac{3}{6} \log\left(\frac{3}{6}\right) = 1$$

Average entropy information for Windy

$$I(\text{Windy}) = \frac{8}{14} * 0.811 + \frac{6}{14} * 1 = 0.892$$

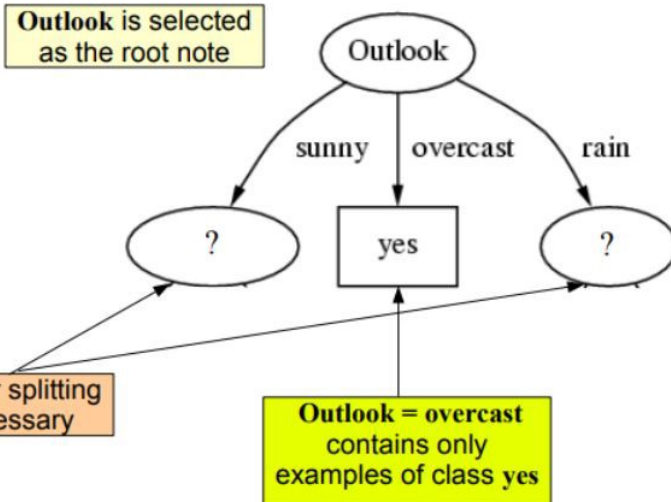
$$\text{Gain}(\text{Windy}) = E(S) - I(\text{Windy}) = 0.94 - 0.892 = 0.048$$



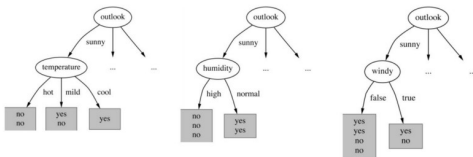
# Pick the Highest Gain Attribute

Outlook		Temperature	
Info:	0.693	Info:	0.911
Gain: $0.940 - 0.693$	0.247	Gain: $0.940 - 0.911$	0.029
Humidity		Windy	
Info:	0.788	Info:	0.892
Gain: $0.940 - 0.788$	0.152	Gain: $0.940 - 0.892$	0.048

# Which One is the Root?



# Create the Tree

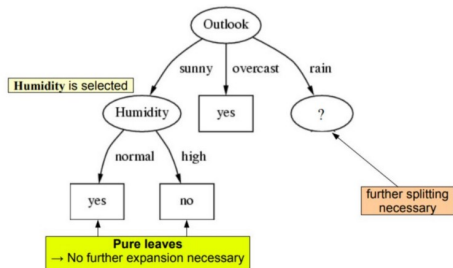


$\text{Gain}(\text{Temperature}) = 0.571 \text{ bits}$

$\text{Gain}(\text{Humidity}) = 0.971 \text{ bits}$

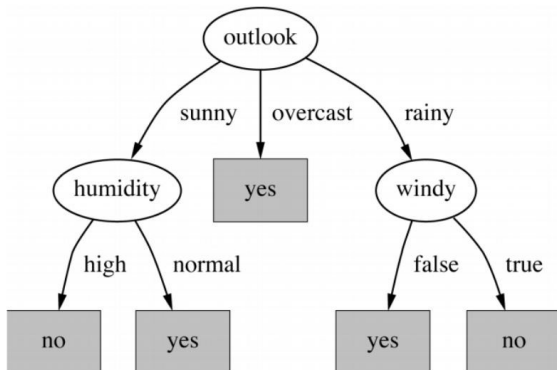
$\text{Gain}(\text{Windy}) = 0.020 \text{ bits}$

**Humidity is selected**

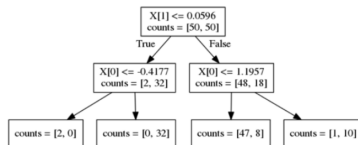
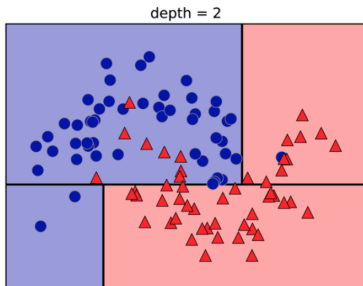


# Final Decision

## Final decision tree



# Predictions with Decision Trees



# Regression Trees

- Prediction:

$$\bar{y}_m = \frac{1}{N_m} \sum_{i \in N_m} y_i$$

- Mean Squared Error:

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} (y_i - \bar{y}_m)^2$$

- Mean Absolute Error:

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} |y_i - \bar{y}_m|$$

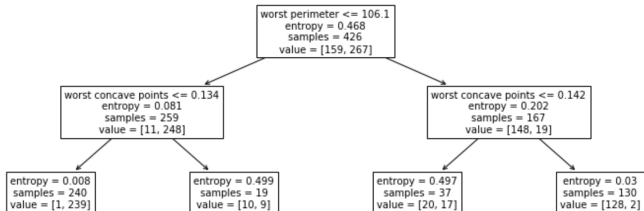


# A Simple Example with Tree

```
from sklearn.datasets import load_breast_cancer
cancer = load_breast_cancer()
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
                                                    stratify=cancer.target,
                                                    random_state=0)
```

```
from sklearn.tree import DecisionTreeClassifier, export_graphviz
tree = DecisionTreeClassifier(max_depth=2)
tree.fit(X_train, y_train)
```

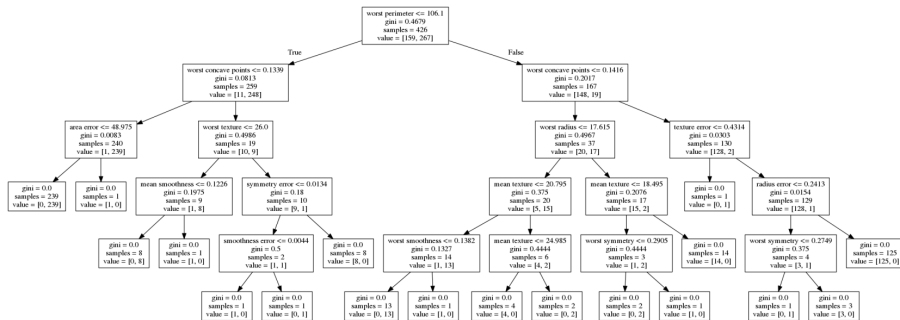
```
from sklearn.tree import plot_tree
tree_dot = plot_tree(tree, feature_names=cancer.feature_names)
```



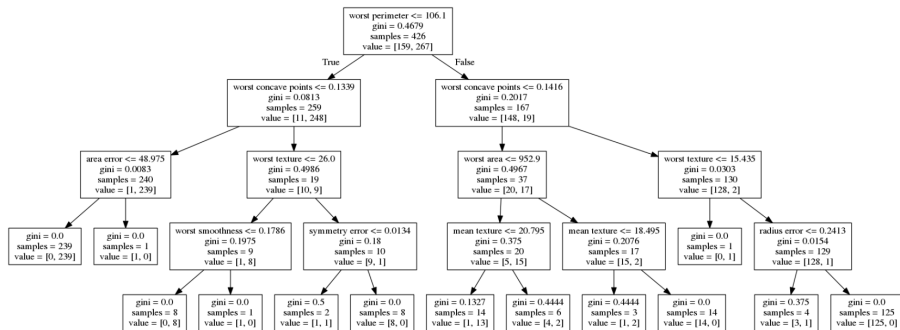
# Parameter Tuning

- Limit tree size (pick one, maybe two):
  - max\_depth
  - max\_leaf\_nodes
  - min\_samples\_split
  - min\_impurity\_decrease

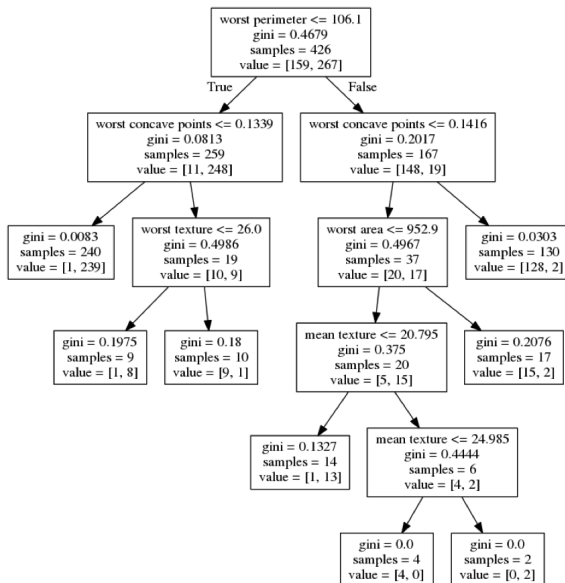
# No Pruning



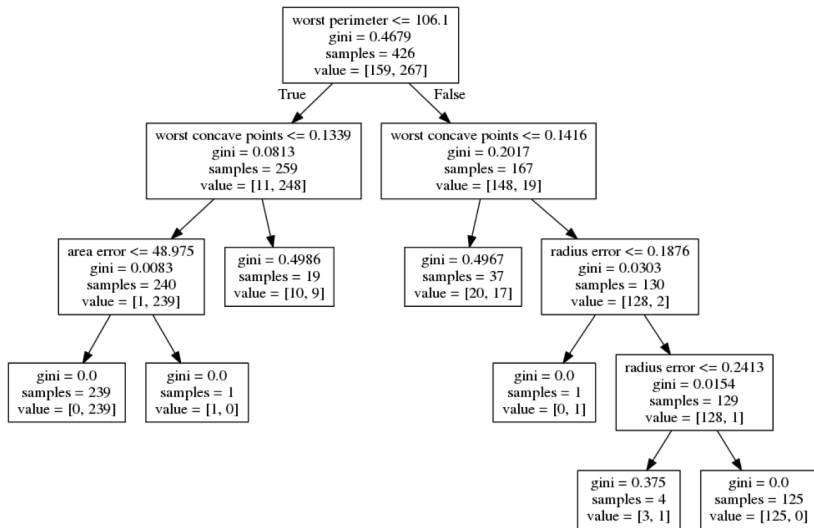
# max\_depth = 4



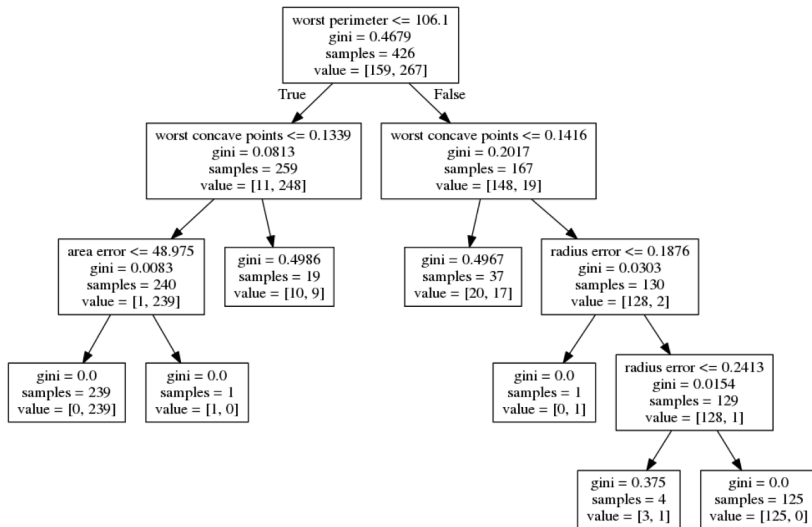
# max\_leaf\_nodes = 8



# min\_samples\_split = 50

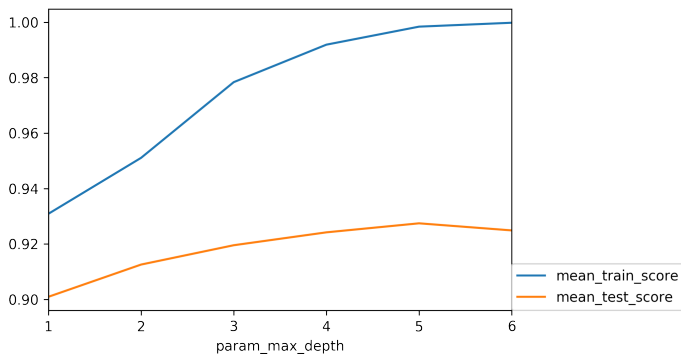


# min\_samples\_split = 50



# Decision Tree Accuracy for Different max\_depth

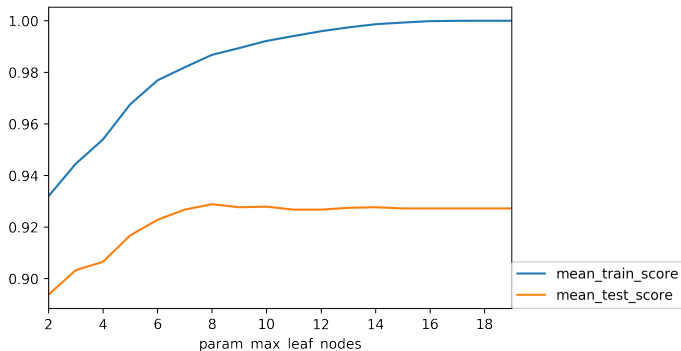
```
1 from sklearn.model_selection import GridSearchCV
2 param_grid = {'max_depth':range(1, 7)}
3 grid = GridSearchCV(DecisionTreeClassifier(random_state=0),param_grid=param_grid,
4                     cv=10)
5 grid.fit(X_train, y_train)
```





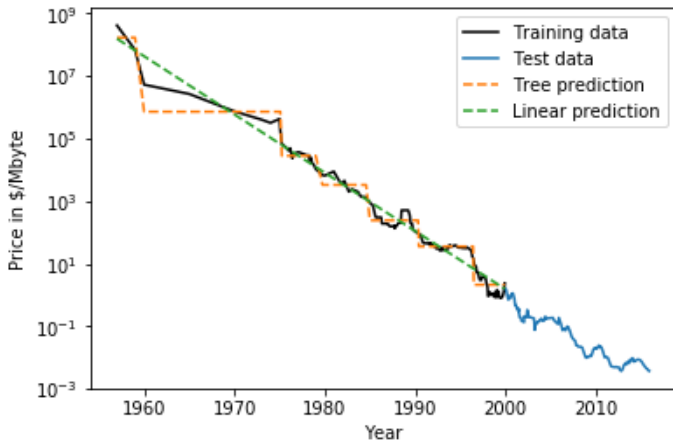
# Decision Tree Accuracy for Different max\_leaf\_nodes

```
1 from sklearn.model_selection import GridSearchCV
2 param_grid = {'max_leaf_nodes':range(2, 20)}
3 grid = GridSearchCV(DecisionTreeClassifier(random_state=0),
4                     param_grid=param_grid, cv=10)
5 grid.fit(X_train, y_train)
```

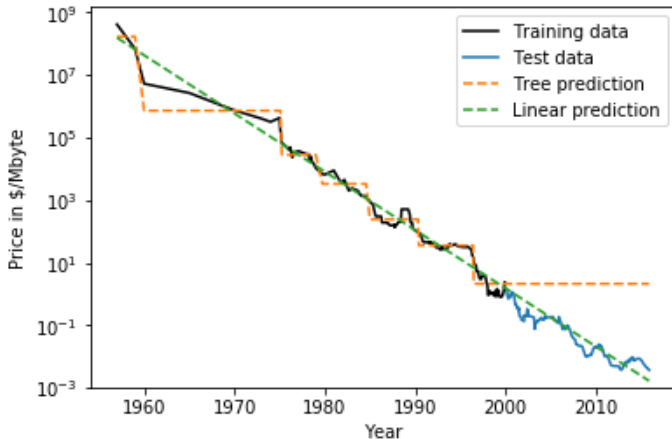


# Relation to Nearest Neighbors

- Predict average of neighbors – either by  $k$ , by epsilon ball or by leaf.
- Trees are much faster to predict.
- Both can't extrapolate



# Extrapolation

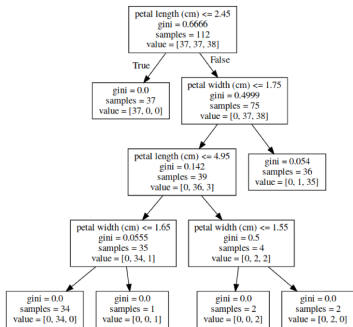


# Instability

```

1 X_train, X_test, y_train, y_test =
2 train_test_split(iris.data, iris.target,
3   stratify=iris.target, random_state=0)
4 tree = DecisionTreeClassifier(max_leaf_nodes=6)
5 tree.fit(X_train, y_train)

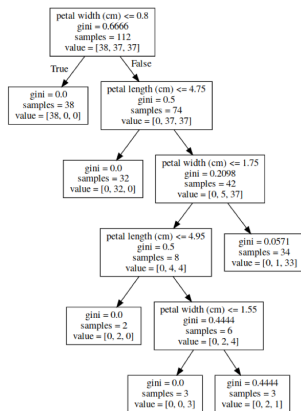
```



```

1 X_train, X_test, y_train, y_test =
2 train_test_split(iris.data, iris.target,
3   stratify=iris.target, random_state=1)
4 tree = DecisionTreeClassifier(max_leaf_nodes=6)
5 tree.fit(X_train, y_train)

```



# Feature Importance

```

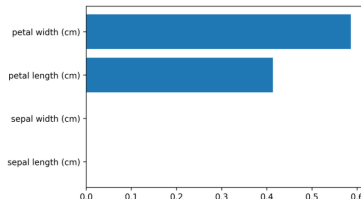
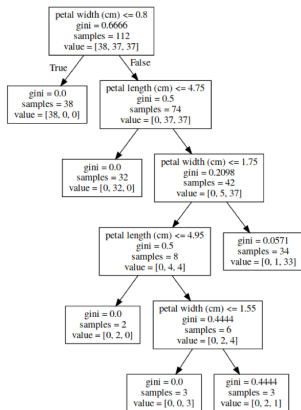
1 X_train, X_test, y_train, y_test =
  train_test_split(
2     iris.data, iris.target, stratify=iris.target,
    random_state=1)
3 tree = DecisionTreeClassifier(max_leaf_nodes=6)
4 tree.fit(X_train, y_train)

```

```

1 tree.feature_importances_
2 array([0.0, 0.0, 0.414, 0.586])

```



# Predicting Probabilities

- Fraction of class in leaf.
- Without pruning: Always 100% certain!
- Even with pruning might be too certain.

# Table of Contents

1 Decision Trees

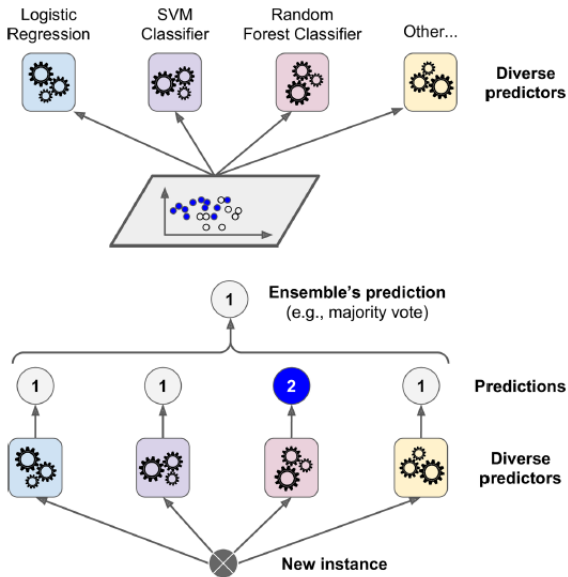
2 Ensemble Models

# Poor Man's Ensembles

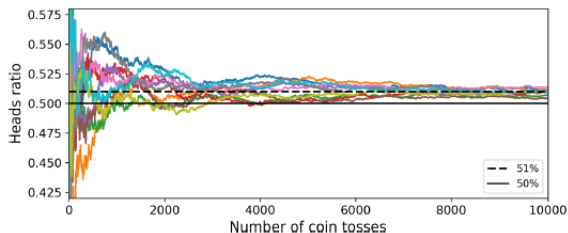
- Build different models
- Average the result
- More models are better – if they are not correlated.
- Also works with neural networks
- You can average any models as long as they provide calibrated (“good”) probabilities



# Voting Classifiers



# Voting Classifier and Law of Large Numbers

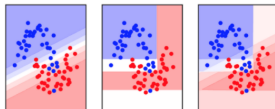


- A slightly biased coin with 51% chance of coming heads
- Toss it 1,000 times, it is 75% more likely to have heads as majority
- Toss it 10,000 times, it is 97% more likely to have heads as majority
- Keep in mind that each toss is independent from the others

# Voting Classifier

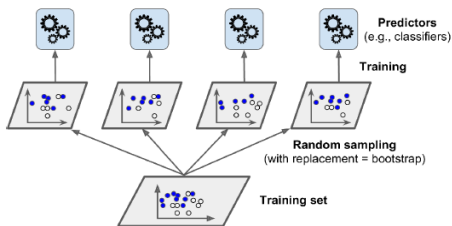
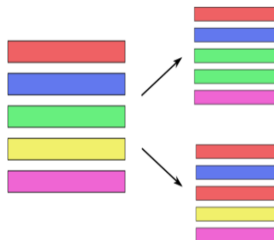
```
voting = VotingClassifier(  
    [('logreg', LogisticRegression(C=100)),  
     ('tree', DecisionTreeClassifier(max_depth=3, random_state=0))],  
    voting='soft')  
voting.fit(X_train, y_train)  
lr, tree = voting.estimators_  
voting.score(X_test, y_test), lr.score(X_test, y_test), tree.score(X_test, y_test)
```

0.88 0.84 0.80

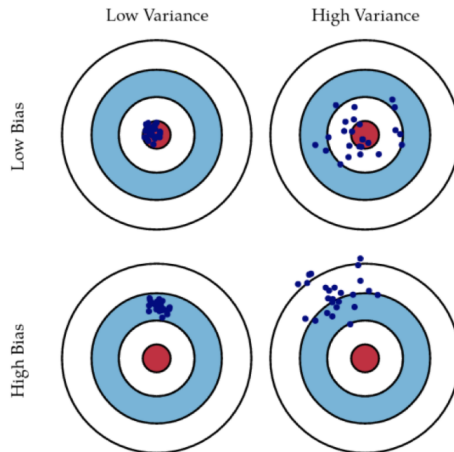


# Bagging (Bootstrap AGGregation)

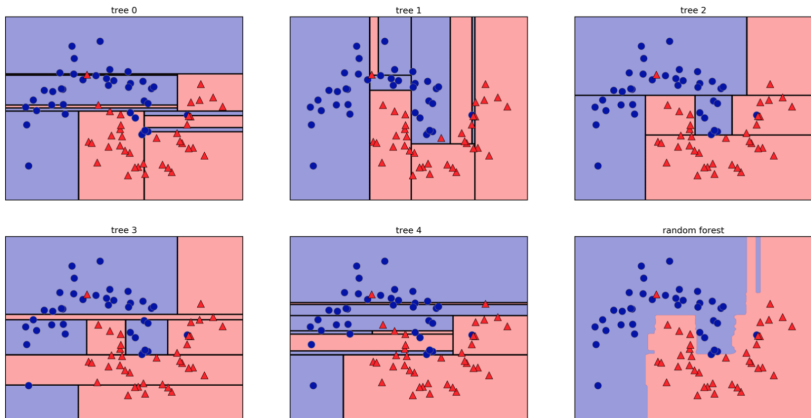
- Generic way to build “slightly different” models
- Use the same training algorithm for every predictor and train them on different random subsets of the training set



# Bias and Variance

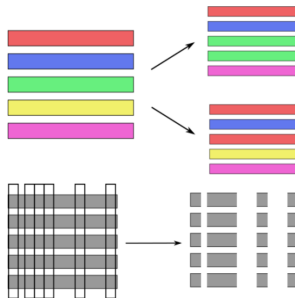


# Random Forest



# Randomize in Two Ways

- For each tree:
  - Pick bootstrap sample of data
- For each split:
  - Pick random sample of features
- More trees are always better

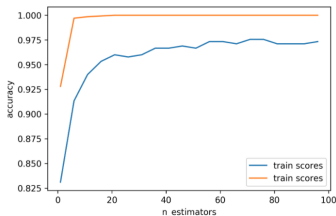


# A Simple Example for Random Forest

```

train_scores = []
test_scores = []

rf = RandomForestClassifier(warm_start=True)
estimator_range = range(1, 100, 5)
for n_estimators in estimator_range:
    rf.n_estimators = n_estimators
    rf.fit(X_train, y_train)
    train_scores.append(rf.score(X_train, y_train))
    test_scores.append(rf.score(X_test, y_test))
  
```



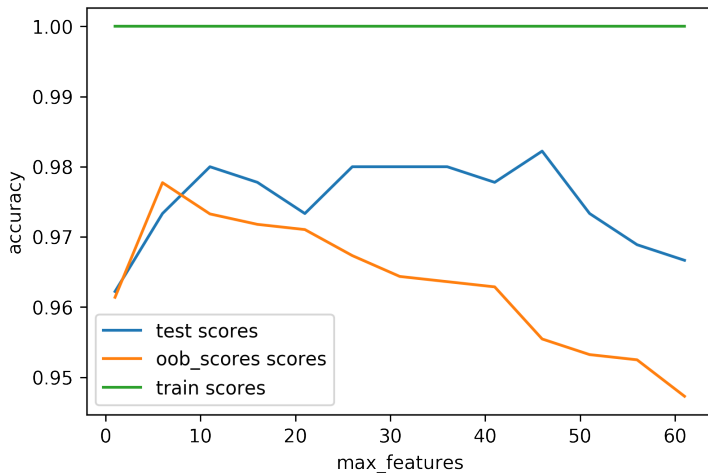


# Out of Bag Estimates

```
1 train_scores = []
2 test_scores = []
3 oob_scores = []
4 feature_range = range(1, 64, 5)
5 for max_features in feature_range:
6     rf = RandomForestClassifier(max_features=max_features, oob_score=True,
7                               n_estimators=200, random_state=0)
8     rf.fit(X_train, y_train)
9     train_scores.append(rf.score(X_train, y_train))
10    test_scores.append(rf.score(X_test, y_test))
11    oob_scores.append(rf.oob_score_)
```

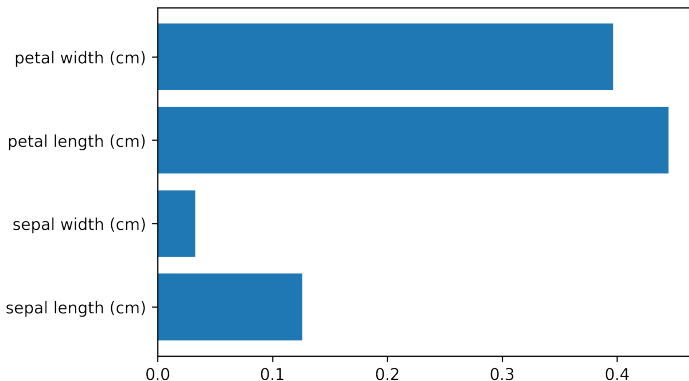
- Each tree only uses 66% of data
- Can evaluate it on the rest!
- Make predictions for out-of-bag, average, score.
- Each prediction is an average over different subset of trees

# Out of Bag Example



# Variable Importance

```
1 X_train, X_test, y_train, y_test = train_test_split(  
2     iris.data, iris.target, stratify=iris.target, random_state=1)  
3 rf = RandomForestClassifier().fit(X_train, y_train)  
4 rf.feature_importances_  
5 plt.barh(range(4), rf.feature_importances_)  
6 plt.yticks(range(4), iris.feature_names);
```



# Summary

- Decision trees are robust to data scaling, but are very unstable models
- To cope with the high variance in the decision trees, we use Random Forest models