# Lecture 06

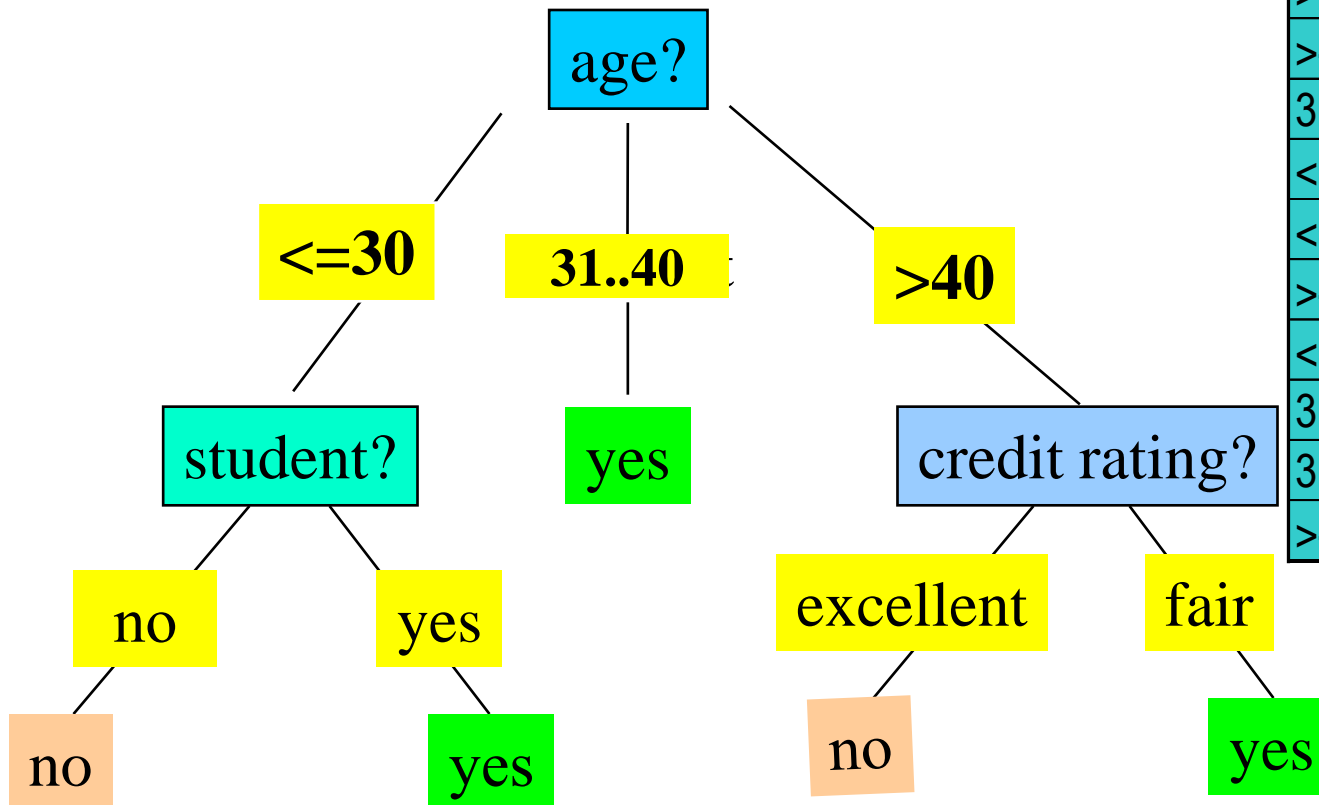# Classification

# Chapter 8. Classification: Basic Concepts

❑ Classification: Basic Concepts

❑ Decision Tree Induction

❑ Bayes Classification Methods

❑ Model Evaluation and Selection

❑ Techniques to Improve Classification Accuracy: Ensemble Methods

❑ Summary

# Decision Tree Induction: An Example

- Training data set: Buys_computer
- Resulting tree:

| age | income | student | credit_rating | buys_computer |
|---|---|---|---|---|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

age?

<=30    31..40    >40

student?

yes

credit rating?

no    yes

excellent    fair

no    yes

no

yes

# Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
  - Tree is constructed in a **top-down recursive divide-and-conquer manner**
  - At start, all the training examples are at the root
  - Attributes are categorical (if continuous-valued, they are discretized in advance)
  - Examples are partitioned recursively based on selected attributes
  - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- Conditions for stopping partitioning
  - All samples for a given node belong to the same class
  - There are no remaining attributes for further partitioning—**majority voting** is employed for classifying the leaf
  - There are no samples left

# Attribute Selection Measure: Information Gain (ID3/C4.5)

- ❑ Select the attribute with the highest information gain

- ❑ Let $p_i$ be the probability that an arbitrary tuple in D belongs to class $C_i$, estimated by $|C_{i, D}|/|D|$

- ❑ Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$$

- ❑ Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$

- ❑ Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

# Attribute Selection: Information Gain

- ❑ Class P: buys_computer = "yes"
- ❑ Class N: buys_computer = "no"

$$Info(D) = I(9,5) = -\frac{9}{14}\log_2(\frac{9}{14}) - \frac{5}{14}\log_2(\frac{5}{14}) = 0.940$$

| age | $p_i$ | $n_i$ | $I(p_i, n_i)$ |
|------|------|------|------|
| <=30 | 2 | 3 | 0.971 |
| 31…40 | 4 | 0 | 0 |
| >40 | 3 | 2 | 0.971 |

| age | income | student | credit_rating | buys_computer |
|------|--------|---------|---------------|---------------|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

$$Info_{age}(D) = \frac{5}{14}I(2,3) + \frac{4}{14}I(4,0)$$
$$+ \frac{5}{14}I(3,2) = 0.694$$

$\frac{5}{14}I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly,

$$Gain(income) = 0.029$$
$$Gain(student) = 0.151$$
$$Gain(credit\_rating) = 0.048$$

# Computing Information-Gain for Continuous-Valued Attributes

❑ Let attribute A be a continuous-valued attribute

❑ Must determine the **best split point** for A

  ❑ Sort the value A in increasing order

  ❑ Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*

    ❑ $(a_i+a_{i+1})/2$ is the midpoint between the values of $a_i$ and $a_{i+1}$

❑ Split:

  ❑ D1 is the set of tuples in D satisfying A ≤ split-point, and D2 is the set of tuples in D satisfying A > split-point

# Gain Ratio for Attribute Selection (C4.5)

❏ Information gain measure is biased towards attributes with a large number of values

❏ C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$SplitInfo_A(D) = -\sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$$

❏ GainRatio(A) = Gain(A)/SplitInfo(A)

❏ Ex. $SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$

❏ gain_ratio(income) = 0.029/1.557 = 0.019

❏ The attribute with the maximum gain ratio is selected as the splitting attribute

# Gini Index

❑ If a data set $D$ contains examples from $n$ classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^{n} p_j^2$$

  where $p_j$ is the relative frequency of class $j$ in $D$

❑ If a data set $D$ is split on A into two subsets $D_1$ and $D_2$, the $gini$ index $gini(D)$ is defined as

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

❑ Reduction in Impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

❑ The attribute which provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity i.e. Δgini(A)) is chosen to split the node

# Computation of Gini Index

❑ Ex.  D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

❑ Suppose the attribute income partitions D into 10 in $D_1$: {low, medium} and 4 in $D_2$

$$gini_{income \in \{low,medium\}}(D) = \left(\frac{10}{14}\right)Gini(D_1) + \left(\frac{4}{14}\right)Gini(D_2)$$

$$= \frac{10}{14}\left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$$

$$= 0.443$$

$$= Gini_{income \in \{high\}}(D).$$

Gini$_{\{low,high\}}$ is 0.458; Gini$_{\{medium,high\}}$ is 0.450.  Thus, split on the {low,medium} (and {high}) since it has the lowest Gini index

# Comparing Attribute Selection Measures

❑ The three measures, in general, return good results but

❑ **Information gain**:

 ❑ biased towards multivalued attributes

❑ **Gain ratio**:

 ❑ tends to prefer unbalanced splits in which one partition is much smaller than the others

❑ **Gini index**:

 ❑ biased to multivalued attributes

 ❑ has difficulty when # of classes is large

 ❑ tends to favor tests that result in equal-sized partitions and purity in both partitions

# Overfitting and Tree Pruning

❑ <u>Overfitting</u>:  An induced tree may overfit the training data

   ❑ Too many branches, some may reflect anomalies due to noise or outliers

   ❑ Poor accuracy for unseen samples

❑ Two approaches to avoid overfitting

   ❑ <u>Prepruning</u>: *Halt tree construction early*-do not split a node if this would result in the goodness measure falling below a threshold

     ❑ Difficult to choose an appropriate threshold

   ❑ <u>Postpruning</u>: *Remove branches* from a "fully grown" tree—get a sequence of progressively pruned trees

     ❑ Use a set of data different from the training data to decide which is the "best pruned tree"- validation data

# Classification in Large Databases

❑ Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed

❑ Why is decision tree induction popular?

  ❑ relatively faster learning speed (than other classification methods)

  ❑ convertible to simple and easy to understand classification rules

  ❑ can use SQL queries for accessing databases

  ❑ comparable classification accuracy with other methods

# Chapter 8. Classification: Basic Concepts

❑ Classification: Basic Concepts

❑ Decision Tree Induction

❑ Bayes Classification Methods

❑ Model Evaluation and Selection

❑ Techniques to Improve Classification Accuracy: Ensemble Methods

❑ Summary

# Bayesian Classification: Why?

❑ <u>A statistical classifier</u>: performs *probabilistic prediction, i.e.,* predicts class membership probabilities

❑ <u>Foundation:</u> Based on Bayes' Theorem.

❑ <u>Performance:</u> A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers

❑ <u>Incremental</u>: Each training example can incrementally increase/decrease the probability that a hypothesis is correct

# Bayes' Theorem: Basics

❑ Bayes' Theorem:

$$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$$

- ❑ P(H) (*prior probability*): the initial probability
  - ❑ E.g., **X** will buy computer, regardless of age, income, …
- ❑ P(**X**): probability that sample data is observed
- ❑ P(**X**|H) (likelihood): the probability of observing the sample **X**, given that the hypothesis holds
  - ❑ E.g., Given that **X** will buy computer, the prob. that X is 31..40, medium income

# Prediction Based on Bayes' Theorem

❑ Given training data **X**, *posteriori probability of a hypothesis* H, P(H|**X**), follows the Bayes' theorem

$$P(H \mid \mathbf{X}) = \frac{P(\mathbf{X} \mid H) P(H)}{P(\mathbf{X})} = P(\mathbf{X} \mid H) \times P(H) / P(\mathbf{X})$$

❑ Practical difficulty:  It requires initial knowledge of many probabilities, involving significant computational cost

# Classification Is to Derive the Maximum Posteriori

❑ Since P(X) is constant for all classes, only

$$P(C_i|\mathbf{X}) = \frac{P(\mathbf{X}|C_i)P(C_i)}{P(\mathbf{X})}$$

needs to be maximized

$$P(C_i|\mathbf{X}) = P(\mathbf{X}|C_i)P(C_i)$$

# Naïve Bayes Classifier: Training Dataset

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data to be classified:

X = (age <=30, Income = medium, Student = yes, Credit_rating = Fair)

| age | income | student | credit_rating | _com |
|---|---|---|---|---|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

# Naïve Bayes Classifier: An Example

❑  $P(C_i)$:  P(buys_computer = "yes") = 9/14 = 0.643

  P(buys_computer = "no") = 5/14 = 0.357

❑  Compute $P(X|C_i)$ for each class

P(age = "<=30" | buys_computer = "yes") = 2/9 = 0.222

P(age = "<= 30" | buys_computer = "no") = 3/5 = 0.6

P(income = "medium" | buys_computer = "yes") = 4/9 = 0.444

P(income = "medium" | buys_computer = "no") = 2/5 = 0.4

P(student = "yes" | buys_computer = "yes) = 6/9 = 0.667

P(student = "yes" | buys_computer = "no") = 1/5 = 0.2

P(credit_rating = "fair" | buys_computer = "yes") = 6/9 = 0.667

P(credit_rating = "fair" | buys_computer = "no") = 2/5 = 0.4

| age | income | student | credit_rating | com |
|------|--------|---------|---------------|-----|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

❑  **X = (age <= 30 , income = medium, student = yes, credit_rating = fair)**

**P(X|C_i) :** P(X|buys_computer = "yes") = 0.222 x 0.444 x 0.667 x 0.667 = 0.044

  P(X|buys_computer = "no") = 0.6 x 0.4 x 0.2 x 0.4 = 0.019

**P(X|C_i)*P(C_i) :** P(X|buys_computer = "yes") * P(buys_computer = "yes") = 0.028

  P(X|buys_computer = "no") * P(buys_computer = "no") = 0.007

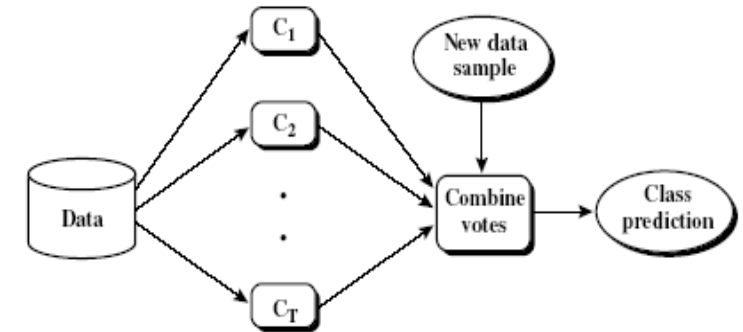**Therefore,  X belongs to class ("buys_computer = yes")**

$$P(C_i|\mathbf{X}) = P(\mathbf{X}|C_i)P(C_i)$$

# Naïve Bayes Classifier: Comments

❑ Advantages

  ❑ Easy to implement

  ❑ Good results obtained in most of the cases

❑ Disadvantages

❑ Assumption: class conditional independence, therefore loss of accuracy

❑ Practically, dependencies exist among variables

  ❑ E.g., Patient's Profile: age, family history, etc.

    Symptoms: fever, cough etc., Disease: lung cancer, diabetes, etc.

  ❑ Dependencies among these cannot be modeled by Naïve Bayes Classifier

❑ How to deal with these dependencies? Bayesian Belief Networks

# Ensemble Methods: Increasing the Accuracy

- ❏ Ensemble methods
  - ❏ Use a combination of models to increase accuracy
  - ❏ Combine a series of k learned models, $M_1$, $M_2$, ..., $M_k$, with the aim of creating an improved model M*
- ❏ Popular ensemble methods
  - ❏ Bagging: averaging the prediction over a collection of classifiers
  - ❏ Boosting: weighted vote with a collection of classifiers
  - ❏ Ensemble: combining a set of heterogeneous classifiers
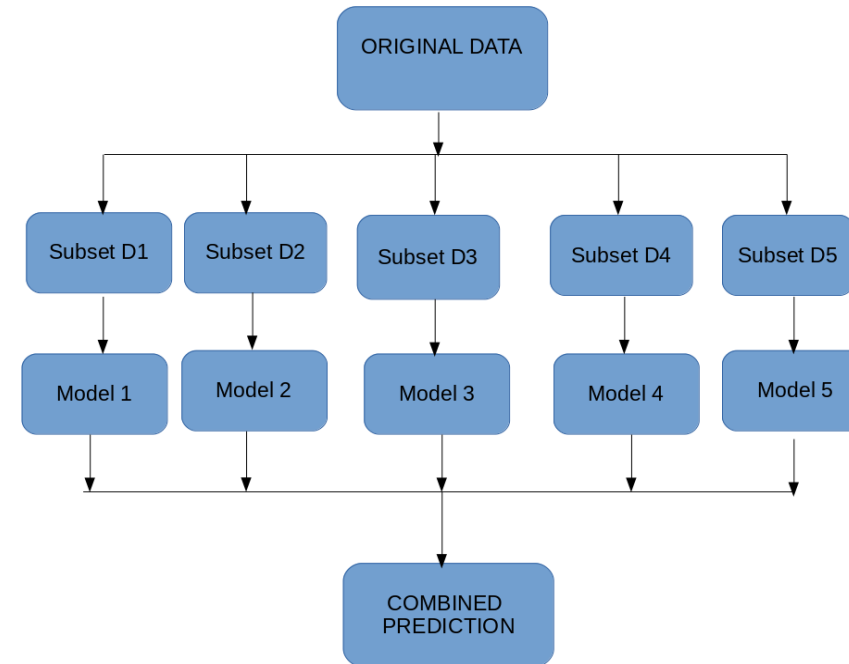
# Bagging

**Bagging** works as follows:-

1.Multiple subsets are created from the original dataset, selecting observations with replacement.

2.A base model (weak model) is created on each of these subsets.

3.The models run in parallel and are independent of each other.

4.The final predictions are determined by combining the predictions from all the models.
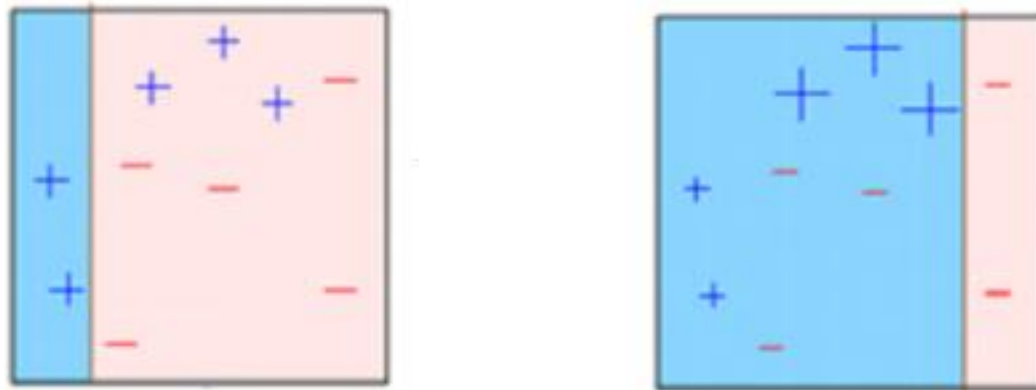


Accuracy: Proven improved accuracy in prediction
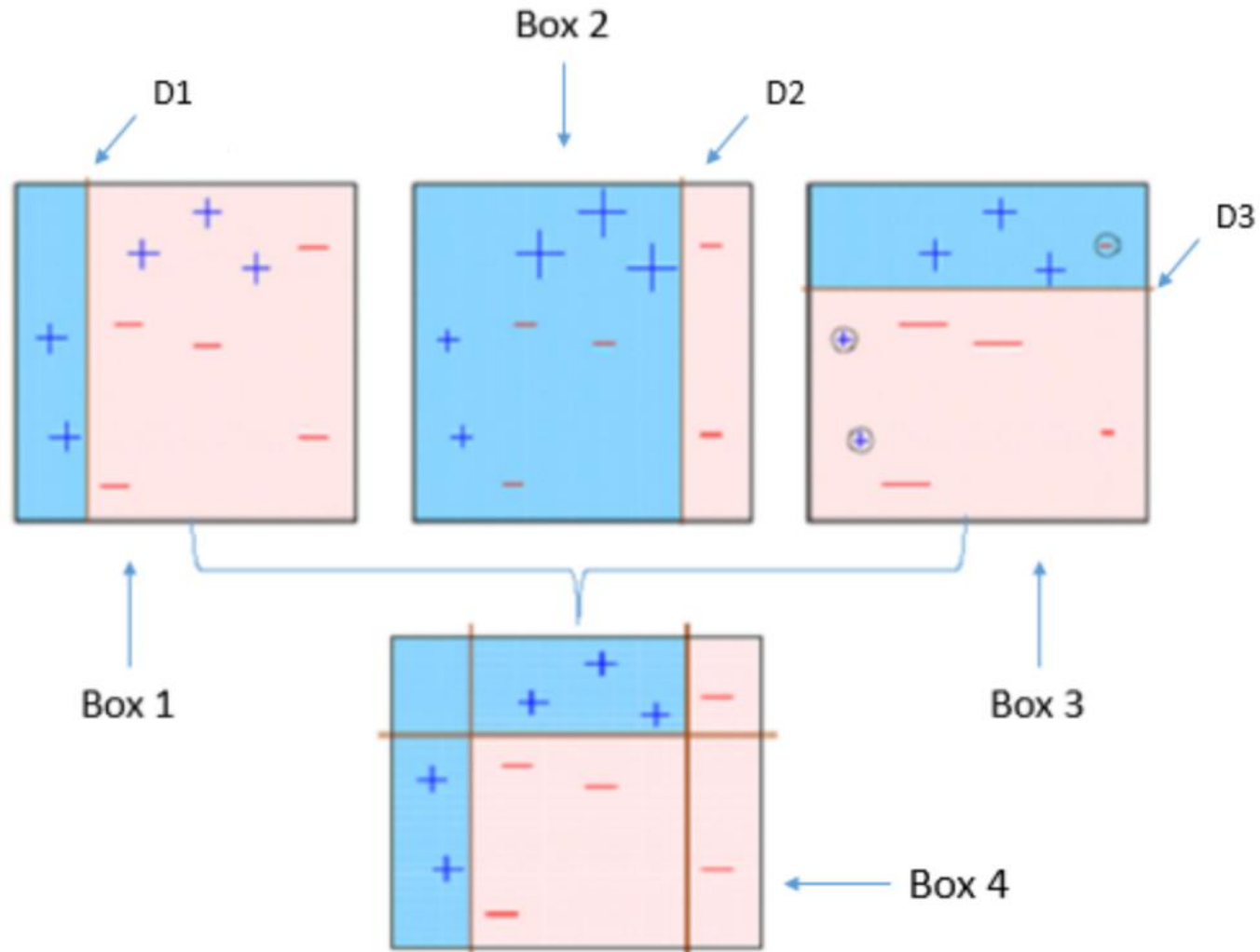- Often significantly better than a single classifier derived from D
- For noise data: not considerably worse, more robust

# Boosting

- Boosting is a **sequential** process, where each subsequent model attempts to correct the errors of the previous model. The succeeding models are dependent on the previous model.
- When an input is misclassified by a hypothesis, its weight is increased so that next hypothesis is more likely to classify it correctly. By combining the whole set at the end converts weak learners into better performing model.

# Boosting

# Random Forest

❑ Random Forest:

   ❑ Each classifier in the ensemble is a *decision tree* classifier and is generated using a random selection of attributes at each node to determine the split

   ❑ During classification, each tree votes and the most popular class is returned

❑ Two Methods to construct Random Forest:

   ❑ Forest-RI (*random input selection*): Randomly select, at each node, F attributes as candidates for the split at the node.

   ❑ Forest-RC (*random linear combinations*): Creates new attributes (or features) that are a linear combination of the existing attributes (reduces the correlation between individual classifiers)

❑ More robust to errors and outliers

❑ Insensitive to the number of attributes selected for consideration at each split, and faster than bagging or boosting

# Classification of Class-Imbalanced Data Sets

❑ Class-imbalance problem: Rare positive example but numerous negative ones, e.g., medical diagnosis, fraud, oil-spill, fault, etc.

❑ Traditional methods assume a balanced distribution of classes and equal error costs: not suitable for class-imbalanced data

❑ Typical methods in two-class classification:

    ❑ **Oversampling**: re-sampling of data from positive class

    ❑ **Under-sampling**: randomly eliminate tuples from negative class

    ❑ **Threshold-moving**: move the decision threshold, t, so that the rare class tuples are easier to classify, and hence, less chance of costly false negative errors

    ❑ **Ensemble techniques**: Ensemble multiple classifiers introduced above

❑ Still difficult for class imbalance problem on multiclass tasks