# اصول و روش‌های داده‌کاوی (Data Mining)
## درس پنجم: روش های پایه دسته بندی و پیش بینی داده ها

**مدرس: سمیرا لویمی**

گروه مهندسی کامپیوتر – دانشگاه شهید چمران اهواز

# Data Mining:

## Concepts and Techniques

### (3rd ed.)

### — Chapter 8 —

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University

# Chapter 8. Classification: Basic Concepts

- Classification: Basic Concepts

- Decision Tree Induction

- Bayes Classification Methods

- Rule-Based Classification

- Model Evaluation and Selection

- Techniques to Improve Classification Accuracy: Ensemble Methods

- Summary

# Supervised vs. Unsupervised Learning

- **Supervised learning (classification)**

  - Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations

  - New data is classified based on the training set

- **Unsupervised learning (clustering)**

  - The class labels of training data is unknown

  - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

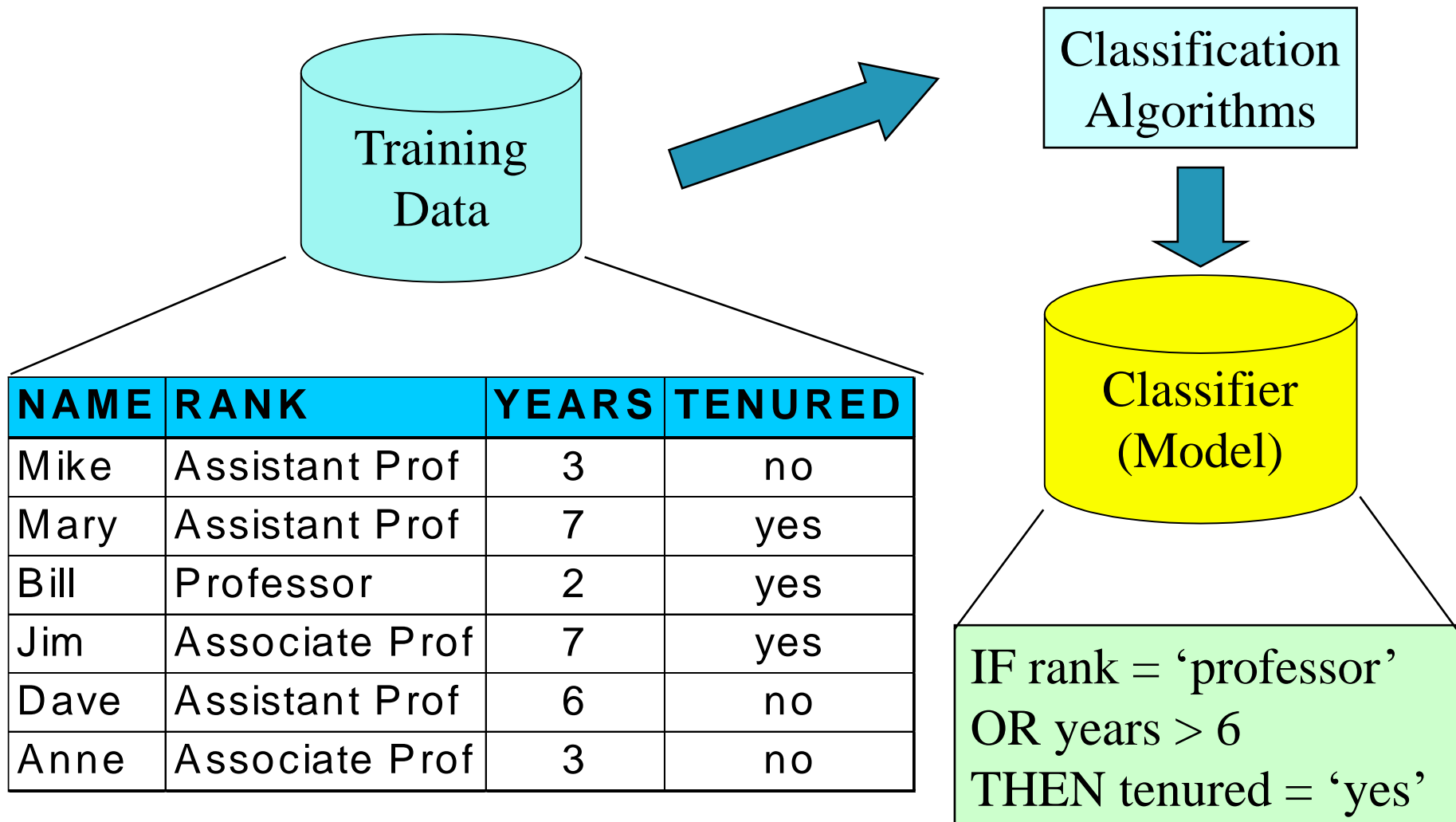# Prediction Problems: Classification vs. Numeric Prediction

- **Classification**
  - predicts categorical class labels (discrete or nominal)
  - classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data
- **Numeric Prediction**
  - models continuous-valued functions, i.e., predicts unknown or missing values
- Typical applications
  - Credit/loan approval:
  - Medical diagnosis: if a tumor is cancerous or benign
  - Fraud detection: if a transaction is fraudulent
  - Web page categorization: which category it is
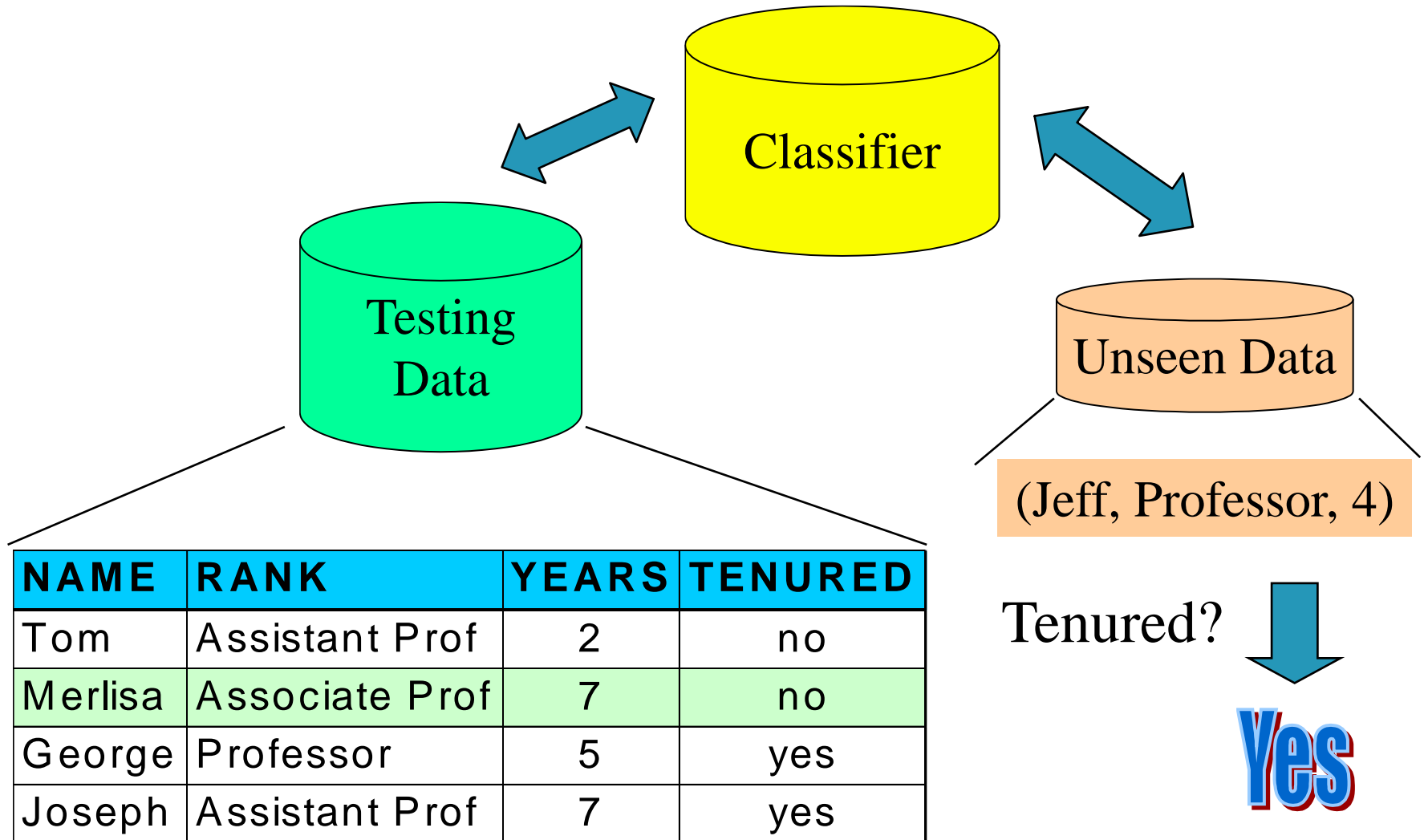
# Classification—A Two-Step Process

- **Model construction**: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
  - The set of tuples used for model construction is **training set**
  - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage**: for classifying future or unknown objects
  - **Estimate accuracy** of the model
    - The known label of test sample is compared with the classified result from the model
    - **Accuracy** rate is the percentage of test set samples that are correctly classified by the model
    - **Test set** is independent of training set (otherwise overfitting)
  - If the accuracy is acceptable, use the model to **classify new data**
- Note: If *the test set* is used to select models, it is called **validation (test) set**

# Process (1): Model Construction



Training Data

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

Classification Algorithms

Classifier (Model)

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

# Process (2): Using the Model in Prediction

Classifier

Testing Data

Unseen Data

(Jeff, Professor, 4)

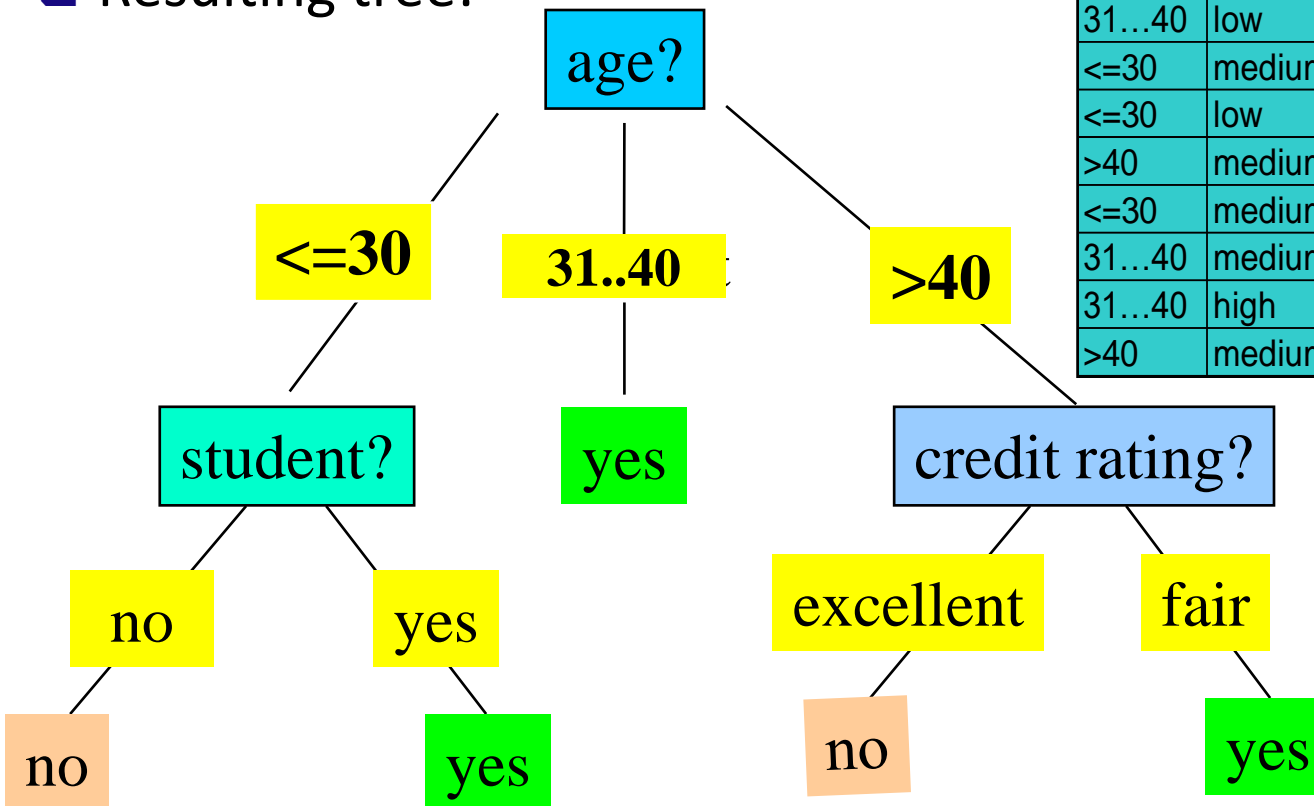| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

Tenured?

Yes

# Decision Tree Induction: An Example

- Training data set: Buys_computer
- The data set follows an example of Quinlan's ID3 (Playing Tennis)
- 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

# How to Handle Continuous-Valued Attributes?

- **Method 1:** Discretize continuous values and treat them as categorical values

  - E.g., age: < 20, 20..30, 30..40, 40..50, > 50

- **Method 2**: Determine the ***best split point*** for continuous-valued attribute A

  - Sort:, e.g. 15, 18, 21, 22, 24, 25, 29, 31, …

  - *Possible split point:* $(a_i + a_{i+1})/2$

    - e.g., (15+18)/2 = 16.5, 19.5, 21.5, 23, 24.5, 27, 30, …

  - The point with the *maximum information gain* for A is selected as the **split-point** for A

- Split:  Based on split point P

  - The set of tuples in D satisfying A ≤ P vs. those with A > P

# Brief Review of Entropy

- Entropy (Information Theory)

  - A measure of uncertainty associated with a random variable

  - Calculation: For a discrete random variable $Y$ taking $m$ distinct values $\{y_1, \ldots, y_m\}$,

    - $H(Y) = -\sum_{i=1}^{m} p_i \log(p_i)$ , where $p_i = P(Y = y_i)$

  - Interpretation:

    - Higher entropy => higher uncertainty

    - Lower entropy => lower uncertainty

- Conditional Entropy

  - $H(Y|X) = \sum_x p(x) H(Y|X = x)$

**m = 2**

# 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 $a_{i+1}$

  - The point with the *minimum expected information requirement* for A is selected as the split-point for A

- 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)

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

  - 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 Impurity/Index (CART)

- 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 $D$

- The Gini impurity considers a binary split for each attribute.

- 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 provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

# 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

- All attributes are assumed continuous-valued

- May need other tools, e.g., clustering, to get the possible split values

- Can be modified for categorical attributes

# 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

# Think!

*"So, what is the relationship between Gini impurity and information gain?"* Intuitively, both measures aim to quantify to what extent the impurity will be reduced if we split the current node based on the given attribute. Information gain, rooted in information theory, measures the impurity based on (the change of) the average amount of information needed to identify the class label of a tuple. Gini impurity is related to *mis-classification* in the following way. Based on the class label distribution in the current node, it tells how likely a randomly chosen tuple will be mis-classified if it is assigned to a random class label. Gini impurity is always used for binary split, whereas information gain allows multiway split. In terms of computation, Gini impurity is slightly more efficient than information gain, since the latter involves the logarithm computation. In practice, however, both measures often lead to very similar decision trees.

# Other Attribute Selection Measures

- <u>CHAID</u>: a popular decision tree algorithm, measure based on $\chi^2$ test for independence

- <u>C-SEP</u>: performs better than info. gain and gini index in certain cases

- <u>G-statistic</u>: has a close approximation to $\chi^2$ distribution

- <u>MDL (Minimal Description Length) principle</u> (i.e., the simplest solution is preferred):

  - The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree

- Multivariate splits (partition based on multiple variable combinations)

  - <u>CART</u>: finds multivariate splits based on a linear comb. of attrs.

- Which attribute selection measure is the best?

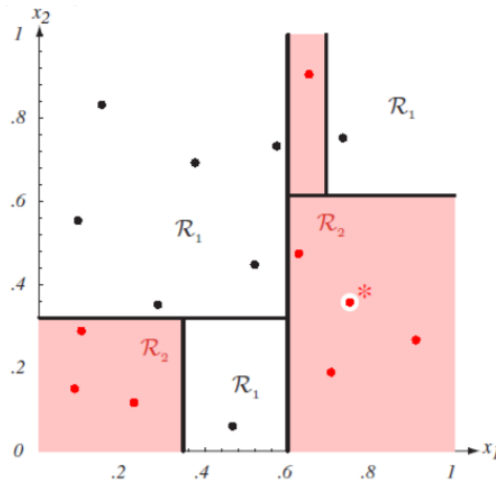  - Most give good results, none is significantly superior than others

# Pro's and Con's

- Pro's
  - Easy to explain (even for non-expert)
  - Easy to implement (many software)
  - Efficient
  - Can tolerant missing data
  - White box
  - No need to normalize data
  - Non-parametric: No assumption on data distribution, no assumption on attribute independency
  - Can work on various attribute types

# Con's

- Con's
  - Unstable. Sensitive to noise
  - Accuracy may be not good enough (depending on your data)
  - The optimal splitting is NP. Greedy algorithms are used
  - Overfitting

# Bayesian Classification: Why?

- A statistical classifier: performs *probabilistic prediction, i.e.,* predicts class membership probabilities

- Foundation: Based on Bayes' Theorem.

- Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers

- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data

- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

# Bayes' Theorem: Basics

- Total probability Theorem: $P(B) = \sum\limits_{i=1}^{M} P(B|A_i)P(A_i)$

- Bayes' Theorem: $P(H|\mathbf{X}) = \dfrac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$

  - Let **X** be a data sample ("*evidence*"): class label is unknown
  - Let H be a *hypothesis* that X belongs to class C
  - Classification is to determine P(H|**X**), (i.e., *posteriori probability):* the probability that the hypothesis holds given the observed data sample **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

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

- Informally, this can be viewed as

  posteriori = likelihood x prior/evidence

- Predicts **X** belongs to $C_i$ iff the probability $P(C_i|\mathbf{X})$ is the highest among all the $P(C_k|X)$ for all the *k* classes

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

# Classification Is to Derive the Maximum Posteriori

- Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n-D attribute vector $\mathbf{X} = (x_1, x_2, ..., x_n)$

- Suppose there are $m$ classes $C_1, C_2, ..., C_m$.

- Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i|\mathbf{X})$

- This can be derived from Bayes' theorem

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

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

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

  needs to be maximized

# Naïve Bayes Classifier

- A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):

$$P(\mathbf{X}|C_i) = \prod_{k=1}^{n} P(x_k|C_i) = P(x_1|C_i) \times P(x_2|C_i) \times ... \times P(x_n|C_i)$$

- This greatly reduces the computation cost: Only counts the class distribution

- If $A_k$ is categorical, $P(x_k|C_i)$ is the # of tuples in $C_i$ having value $x_k$ for $A_k$ divided by $|C_{i,D}|$ (# of tuples of $C_i$ in D)

- If $A_k$ is continous-valued, $P(x_k|C_i)$ is usually computed based on Gaussian distribution with a mean µ and standard deviation σ

and $P(x_k|C_i)$ is

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(\mathbf{X}|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

# 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

| 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 |

- $P(C_i)$:   $P(\text{buys\_computer} = \text{"yes"}) = 9/14 = 0.643$

  $P(\text{buys\_computer} = \text{"no"}) = 5/14 = 0.357$

- Compute $P(X|C_i)$ for each class

  $P(\text{age} = \text{"<=30"} \mid \text{buys\_computer} = \text{"yes"}) = 2/9 = 0.222$

  $P(\text{age} = \text{"<= 30"} \mid \text{buys\_computer} = \text{"no"}) = 3/5 = 0.6$

  $P(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"yes"}) = 4/9 = 0.444$

  $P(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"no"}) = 2/5 = 0.4$

  $P(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"yes}) = 6/9 = 0.667$

  $P(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"no"}) = 1/5 = 0.2$

  $P(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"yes"}) = 6/9 = 0.667$

  $P(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"no"}) = 2/5 = 0.4$

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

**$P(X|C_i)$ : $P(X|\text{buys\_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$**

  **$P(X|\text{buys\_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$**

**$P(X|C_i)*P(C_i)$ :** $P(X|\text{buys\_computer} = \text{"yes"}) * P(\text{buys\_computer} = \text{"yes"}) = 0.028$

  **$P(X|\text{buys\_computer} = \text{"no"}) * P(\text{buys\_computer} = \text{"no"}) = 0.007$**

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

# Avoiding the Zero-Probability Problem

- Naïve Bayesian prediction requires each conditional prob. be **non-zero**. Otherwise, the predicted prob. will be zero

$$P(X \mid C_i) = \prod_{k=1}^{n} P(x_k \mid C_i)$$

- Ex. Suppose a dataset with 1000 tuples, income=low (0), income= medium (990), and income = high (10)

- Use **Laplacian correction** (or Laplacian estimator)
  - *Adding 1 to each case*

    Prob(income = low) = 1/1003

    Prob(income = medium) = 991/1003

    Prob(income = high) = 11/1003
  - The "corrected" prob. estimates are close to their "uncorrected" counterparts

31

# 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., hospitals: patients: 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

# Lazy vs. Eager Learning

- Lazy vs. eager learning

  - **Lazy learning** (e.g., instance-based learning): Simply stores training data (or only minor processing) and waits until it is given a test tuple

  - **Eager learning** (the above discussed methods): Given a set of training tuples, constructs a classification model before receiving new (e.g., test) data to classify

- Lazy: less time in training but more time in predicting

- Accuracy

  - Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form an implicit global approximation to the target function

  - Eager: must commit to a single hypothesis that covers the entire instance space

# Lazy Learner: Instance-Based Methods

- Instance-based learning:
  - Store training examples and delay the processing ("lazy evaluation") until a new instance must be classified
- Typical approaches
  - *k*-nearest neighbor approach
    - Instances represented as points in a Euclidean space.
  - Locally weighted regression
    - Constructs local approximation
  - Case-based reasoning
    - Uses symbolic representations and knowledge-based inference

# The *k*-Nearest Neighbor Algorithm

- All instances correspond to points in the n-D space

- The nearest neighbor are defined in terms of Euclidean distance, dist($X_1$, $X_2$)

- Target function could be discrete- or real- valued

- For discrete-valued, *k*-NN returns the most common value among the *k* training examples nearest to $x_q$

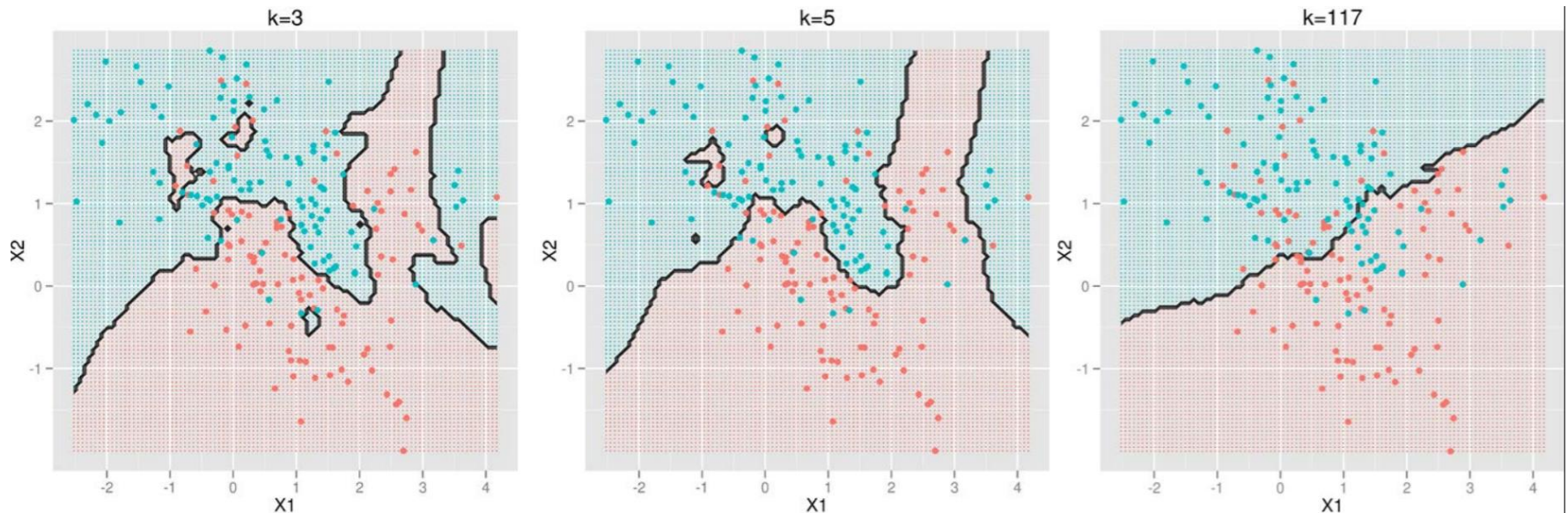- Vonoroi diagram: the decision surface induced by 1-NN for a typical set of training examples

# Discussion on the *k*-NN Algorithm

- *k*-NN for <u>real-valued prediction</u> for a given unknown tuple
  - Returns the mean values of the *k* nearest neighbors
- <u>Distance-weighted</u> nearest neighbor algorithm
  - Weight the contribution of each of the *k* neighbors according to their distance to the query $x_q$    $w = \dfrac{1}{d(x_q, x_i)^2}$
    - Give greater weight to closer neighbors
- <u>Robust</u> to noisy data by averaging *k*-nearest neighbors
- <u>Curse of dimensionality</u>: distance between neighbors could be dominated by irrelevant attributes
  - To overcome it, axes stretch or elimination of the least relevant attributes

# Selection of k for kNN

- The number of neighbors k

  - Small k: overfitting (high var., low bias)

  - Big k: bringing too many irrelevant points (high bias, low var.)



http://scott.fortmann-roe.com/docs/BiasVariance.html

# Case-Based Reasoning (CBR)

- **CBR**: Uses a database of problem solutions to solve new problems
- Store <u>symbolic description</u> (tuples or cases)—not points in a Euclidean space
- <u>Applications:</u> Customer-service (product-related diagnosis), legal ruling
- <u>Methodology</u>
    - Instances represented by rich symbolic descriptions (e.g., function graphs)
    - Search for similar cases, multiple retrieved cases may be combined
    - Tight coupling between case retrieval, knowledge-based reasoning, and problem solving
- <u>Challenges</u>
    - Find a good similarity metric
    - Indexing based on syntactic similarity measure,  and when failure, backtracking, and adapting to additional cases

# Model Evaluation and Selection

- Evaluation metrics
    - How can we measure accuracy?
    - Other metrics to consider?
- Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- Methods for estimating a classifier's accuracy
    - Holdout method
    - Cross-validation
    - Bootstrap (not covered)
- Comparing classifiers:
    - ROC Curves

# Classifier Evaluation Metrics: Confusion Matrix

- **Confusion Matrix:**

| Actual class\Predicted class | $C_1$ | $\neg C_1$ |
|---|---|---|
| $C_1$ | **True Positives (TP)** | **False Negatives (FN)** |
| $\neg C_1$ | **False Positives (FP)** | **True Negatives (TN)** |

- In a confusion matrix w. *m* classes, $CM_{i,j}$ indicates # of tuples in class *i* that were labeled by the classifier as class *j*
  - May have extra rows/columns to provide totals
- **Example of Confusion Matrix:**

| Actual class\Predicted class | play_golf = yes | play_golf = no | Total |
|---|---|---|---|
| play_golf = yes | **6954** | **46** | 7000 |
| play_golf = no | **412** | **2588** | 3000 |
| Total | 7366 | 2634 | 10000 |

# Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

| A\P | C | ¬C | |
|-----|-----|-----|-----|
| C | **TP** | **FN** | **P** |
| ¬C | **FP** | **TN** | **N** |
| | **P'** | **N'** | **All** |

Real-world truth

Predictions

- **Classifier accuracy,** or recognition rate
  - Percentage of test set tuples that are correctly classified

  **Accuracy = (TP + TN)/All**
- **Error rate:** *1 – accuracy*, or
  **Error rate = (FP + FN)/All**

- **Class imbalance problem**
  - One class may be *rare*
    - E.g., fraud, or HIV-positive
  - Significant *majority of the negative class* and minority of the positive class
  - Measures handle the class imbalance problem
  - **Sensitivity** (recall): True positive recognition rate
    - **Sensitivity = TP/P**
  - **Specificity**: True negative recognition rate
    - **Specificity = TN/N**

# Classifier Evaluation Metrics: Precision and Recall, and F-measures

| A\P | C | ¬C | |
|---|---|---|---|
| C | TP | FN | P |
| ¬C | FP | TN | N |
| | P' | N' | All |

- **Precision**: Exactness: what % of tuples that the classifier labeled as positive are actually positive?

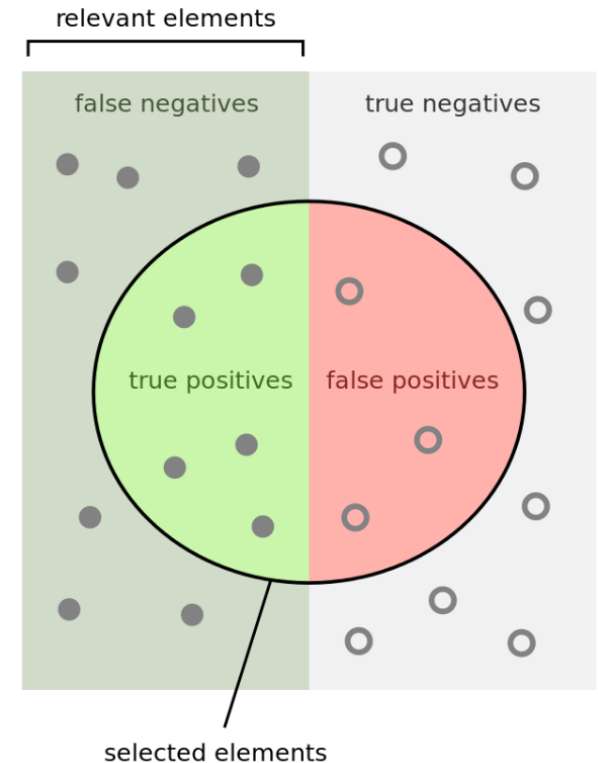$$P = \text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Precision} = \frac{}{}$$

- **Recall:** Completeness: what % of positive tuples did the classifier label as positive?

$$R = \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Recall} = \frac{}{}$$

  - Range: [0, 1]

relevant elements

false negatives        true negatives

true positives    false positives

selected elements

https://en.wikipedia.org/wiki/Precision_and_recall

# Classifier Evaluation Metrics: Precision and Recall, and F-measures

- The "inverse" relationship between precision & recall
- ***We want one number to say if a classifier is good or not***
- **F measure (**or **F-score):** <u>harmonic</u> mean of precision and recall
  - In general, it is the weighted measure of precision & recall

$$F_\beta = \frac{1}{\alpha \cdot \frac{1}{P} + (1 - \alpha) \cdot \frac{1}{R}} = \frac{(\beta^2 + 1)P * R}{\beta^2 P + R}$$

Assigning β times as much weight to recall as to precision

  - ***F1-measure (balanced F-measure)***
    - That is, when β = 1,
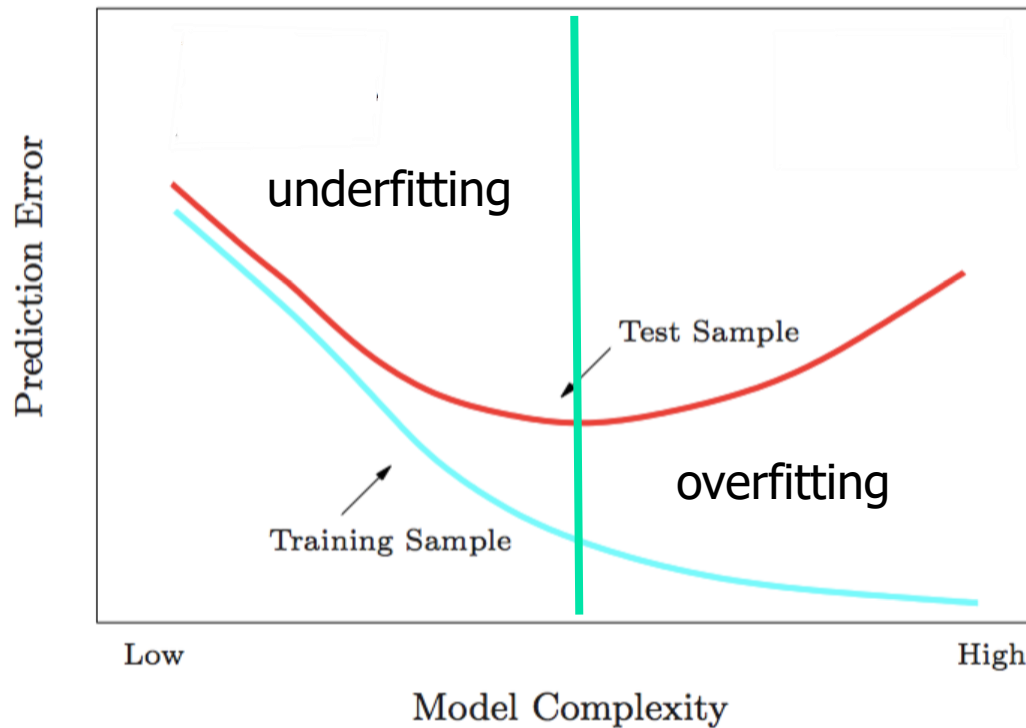
$$F_1 = \frac{2P * R}{P + R}$$

# Classifier Evaluation Metrics: Example

- Use the same confusion matrix, calculate the measure just introduced

| Actual Class\Predicted class | cancer = yes | cancer = no | Total |
|---|---|---|---|
| cancer = yes | **90** | **210** | 300 |
| cancer = no | **140** | **9560** | 9700 |
| Total | 230 | 9770 | 10000 |

- Sensitivity =

- Specificity =

- Accuracy =

- Error rate =

- Precision =

- Recall =

- F1 =

# Classifier Evaluation Metrics: Example

- Use the same confusion matrix, calculate the measure just introduced

| Actual Class\Predicted class | cancer = yes | cancer = no | Total |
|---|---|---|---|
| cancer = yes | **90** | **210** | 300 |
| cancer = no | **140** | **9560** | 9700 |
| Total | 230 | 9770 | 10000 |

- Sensitivity = TP/P = 90/300 = 30%

- Specificity = TN/N = 9560/9700 = 98.56%

- Accuracy = (TP + TN)/All = (90+9560)/10000 = 96.50%

- Error rate = (FP + FN)/All = (140 + 210)/10000 = 3.50%

- Precision = TP/(TP + FP) = 90/(90 + 140) = 90/230 = 39.13%

- Recall = TP/ (TP + FN) = 90/(90 + 210) = 90/300 = 30.00%

- F1 = 2 P × R /(P + R) = 2 × 39.13% × 30.00%/(39.13% + 30%) = 33.96%

# Training Error VS Testing Error



underfitting

overfitting

# Classifier Evaluation: Holdout

- **Holdout method**

    - Given data is randomly partitioned into two independent sets

        - Training set (e.g., 2/3) for model construction

        - Test set (e.g., 1/3) for accuracy estimation

    - Repeated random sub-sampling validation: a variation of holdout

        - Repeat holdout k times, accuracy = avg. of the accuracies
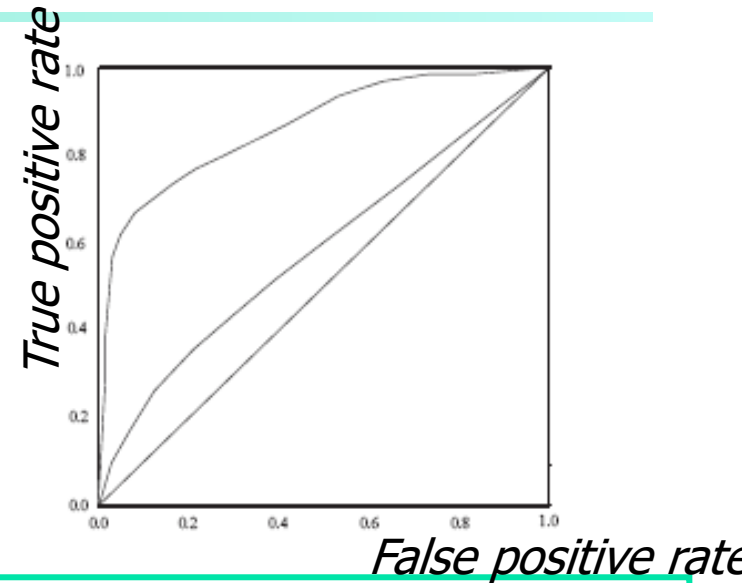
          obtained

# Classifier Evaluation: Cross-Validation

- **Cross-validation** (*k*-fold, where k = 10 is most popular)

  - Randomly partition the data into *k mutually exclusive* subsets, each approximately equal size

  - At *i*-th iteration, use $D_i$ as test set and others as training set

  - Leave-one-out: *k* folds where *k* = # of tuples, for small sized data

  - **\*Stratified cross-validation\***: folds are stratified so that class distribution, in each fold is approximately the same as that in the initial data

# Model Selection: ROC Curves

- **ROC** (Receiver Operating Characteristics) curves: for visual comparison of classification models

- Originated from signal detection theory

- Shows the trade-off between the true positive rate and the false positive rate

- The area under the ROC curve (**AUC**: Area Under Curve) is a measure of the accuracy of the model

- Rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list

- The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model

*True positive rate*

*False positive rate*

- ❑ Vertical axis represents the true positive rate (TP/P)
- ❑ Horizontal axis rep. the false positive rate (FP/N)
- ❑ The plot also shows a diagonal line
- ❑ A model with perfect accuracy will have an area of 1.0
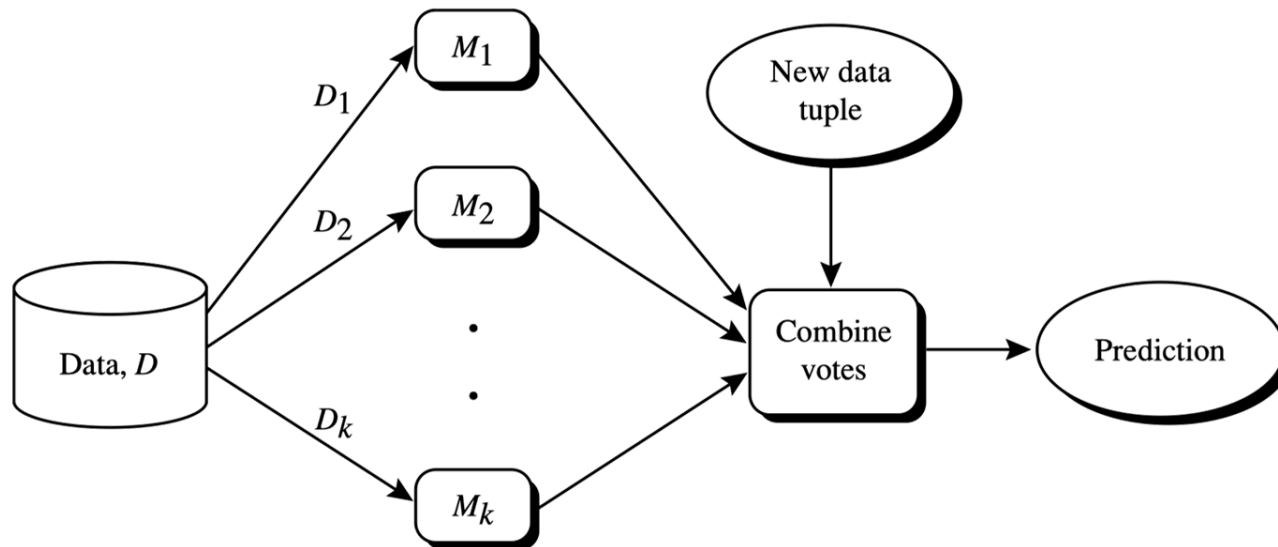
# Techniques to Improve Classification Accuracy

- Introducing Ensemble Methods

- Bagging

- Boosting

- Random Forests

- Imbalanced Classification

# 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*

# Ensemble Methods: Increasing the Accuracy

- What are the requirements to generate an improved model?
  - Example: majority voting

|  | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $M_1$ | ✓ | ✓ | ✗ |
| $M_2$ | ✗ | ✓ | ✓ |
| $M_3$ | ✓ | ✗ | ✓ |
| Voting Ensemble | ✓ | ✓ | ✓ |

Case 1:
Ensemble has positive effect

|  | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $M_1$ | ✓ | ✓ | ✗ |
| $M_2$ | ✓ | ✓ | ✗ |
| $M_3$ | ✓ | ✓ | ✗ |
| Voting Ensemble | ✓ | ✓ | ✗ |

Case 2:
Ensemble has no effect

|  | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $M_1$ | ✓ | ✗ | ✗ |
| $M_2$ | ✗ | ✓ | ✗ |
| $M_3$ | ✗ | ✗ | ✓ |
| Voting Ensemble | ✗ | ✗ | ✗ |

Case 3:
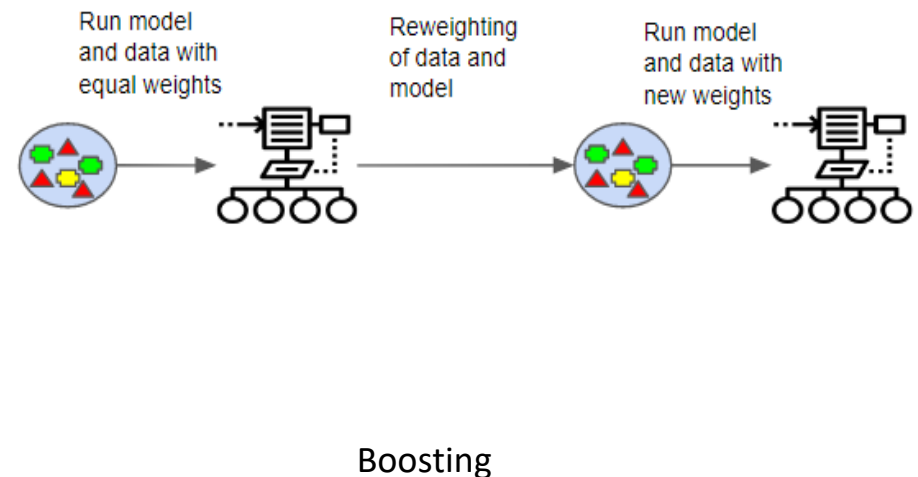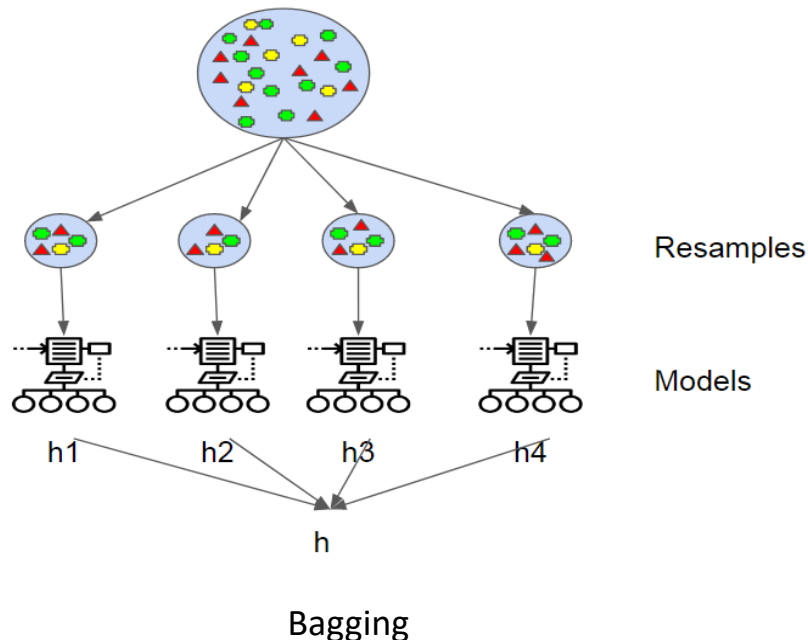Ensemble has negative effect

Base model performance →

Ensemble performance →

- ❑ Base models should be
  - ❑ Accurate
  - ❑ Diverse

# Ensemble Methods: Increasing the Accuracy

- Popular ensemble methods
  - Bagging:  Trains each model using a subset of the training set, and models learned in <u>parallel</u>
  - Boosting:  Trains each new model instance to emphasize the training instances that previous models mis-classified, and models learned <u>in order</u>


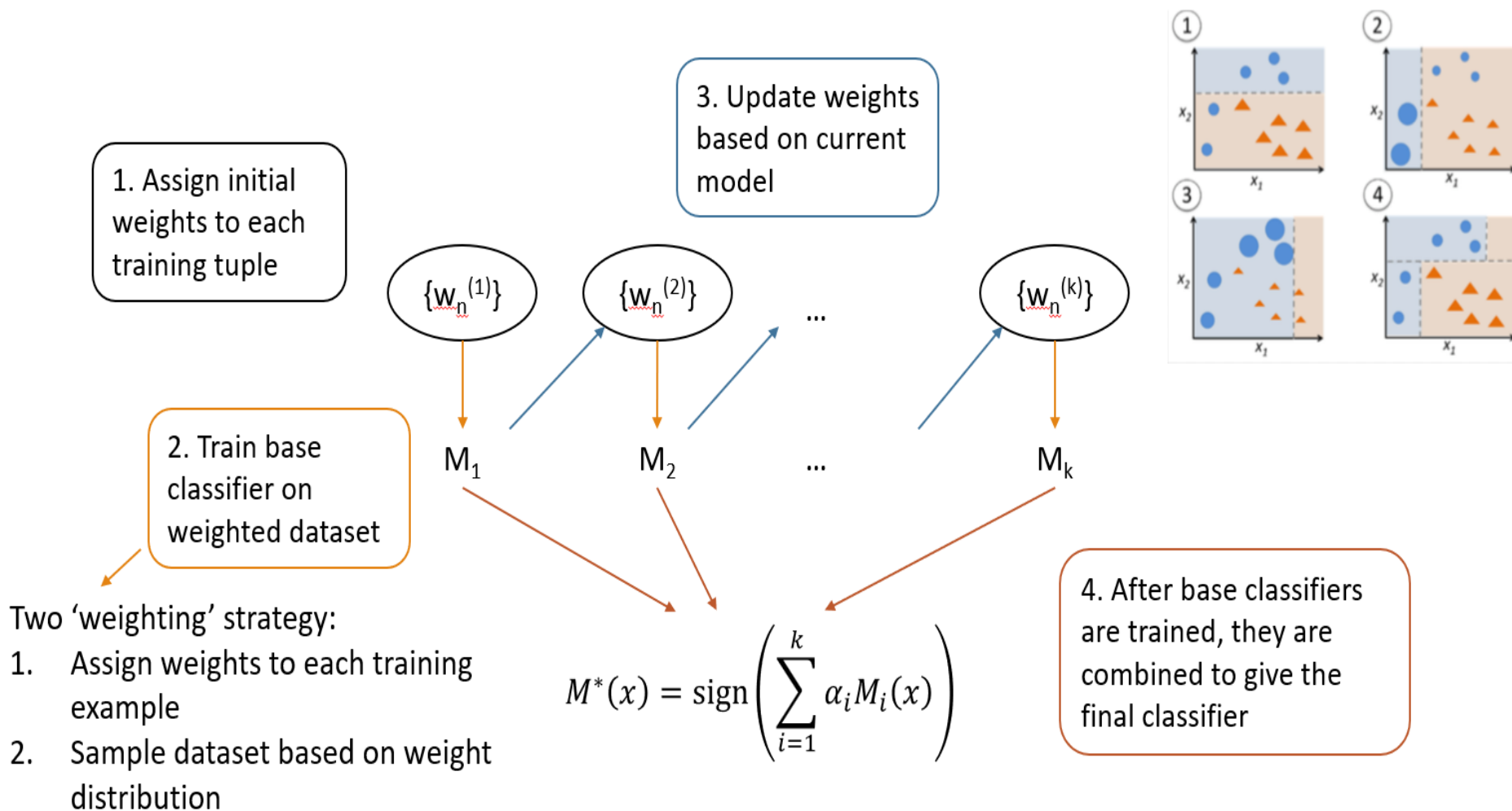
Bagging

Boosting

# Bagging: Bootstrap Aggregation

- Analogy: Diagnosis based on multiple doctors' majority vote

- Training
    - For i = 1 to k
        - create bootstrap sample, $D_i$, by sampling D with replacement;
        - use $D_i$ and the learning scheme to derive a model, $M_i$ ;

- Classification: classify an unknown sample **X**
    - let each of the k models classify X and return the majority vote

- Prediction:
    - To predict continuous variables, use average prediction instead of vote

# Boosting

- Analogy: Consult several doctors, based on a combination of weighted diagnoses—weight assigned based on the previous diagnosis accuracy

- How boosting works?
  - A series of k classifiers are iteratively learned
  - After a classifier $M_i$ is learned, set the subsequent classifier, $M_{i+1}$, to **pay more attention to the training tuples that were misclassified** by $M_i$
  - The final **M* combines the votes** of each individual classifier, where the weight of each classifier's vote is a function of its accuracy

- Boosting algorithm can be extended for numeric prediction

# Adaboost (Freund and Schapire, 1997)

1. Assign initial weights to each training tuple

3. Update weights based on current model



$\{w_n^{(1)}\}$    $\{w_n^{(2)}\}$    ...    $\{w_n^{(k)}\}$

2. Train base classifier on weighted dataset

$M_1$    $M_2$    ...    $M_k$

Two 'weighting' strategy:
1. Assign weights to each training example
2. Sample dataset based on weight distribution

$$M^*(x) = \text{sign}\left(\sum_{i=1}^{k} \alpha_i M_i(x)\right)$$

4. After base classifiers are trained, they are combined to give the final classifier

# Adaboost (Freund and Schapire, 1997)

- **Input**: Training set $D = \{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\}$
- Initialize all weights $\{w_n^{(1)}\}$ to 1/N
- For round i = 1 to k,
  - Fit a classifier $M_i$ based on weighted error function
  
  $$J_m = \sum_{n=1}^{N} w_n^{(i)} I(M_i(x_n) \neq y_n)$$
  
  - Evaluate error rate $\epsilon_i = J_m / \sum w_n^{(i)}$ (stop iteration if $\epsilon_i$ < threshold)
  
    and the base model $M_i$'s vote $\alpha_i = \frac{1}{2} \ln\left(\frac{1-\epsilon_i}{\epsilon_i}\right)$
  
  - Update weights
    $$w_n^{(i+1)} = w_n^{(i)} \exp\{\alpha_i \cdot I(M_i(x_n) \neq y_n)\}$$

- The final model is given by voting based on $\{\alpha_n\}$

# Random Forest: Basic Concepts

- Random Forest (first proposed by L. Breiman in 2001)

    - Bagging with **decision trees** as base models

    - *Data bagging*

        - Use a **subset of training data** by sampling with replacement for each tree

    - *Feature bagging* ⟵

        Advantage of decision trees – more diversity

        - At each node use a **random selection of attributes** as candidates and split by the best attribute among them

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

# Random Forest

- 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. The CART methodology is used to grow the trees to maximum size

  - 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)

- Comparable in accuracy to Adaboost, but more robust to errors and outliers

- Insensitive to the number of attributes selected for consideration at each split, and faster than typical bagging or boosting
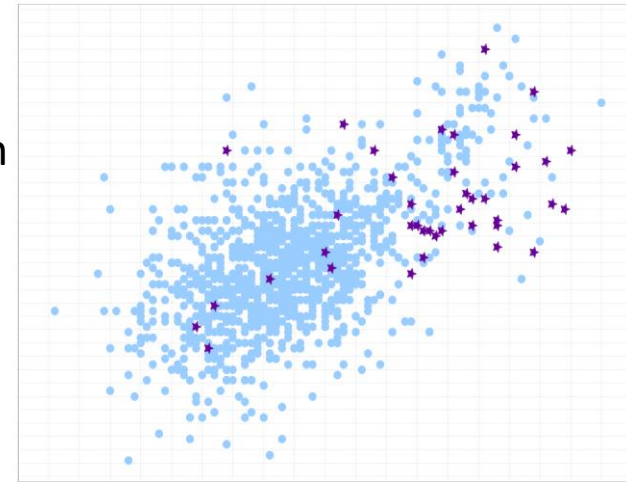
# Ensemble Methods Recap

- Random forest and XGBoost are the most commonly used algorithms for tabular data
- Pros
  - Good performance for tabular data, requires no data scaling
  - Can scale to large datasets
  - Can handle missing data to some extent
- Cons
  - Can overfit to training data if not tuned properly
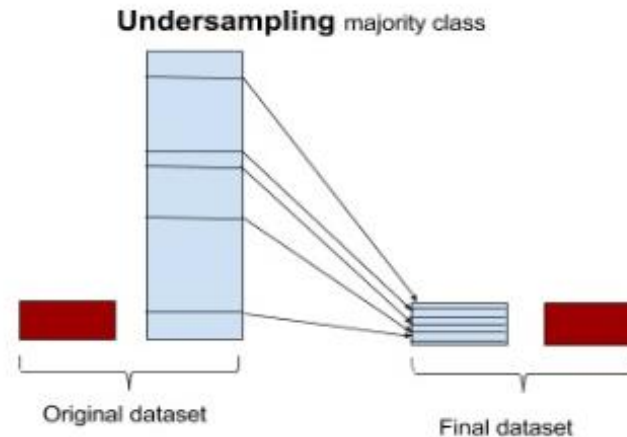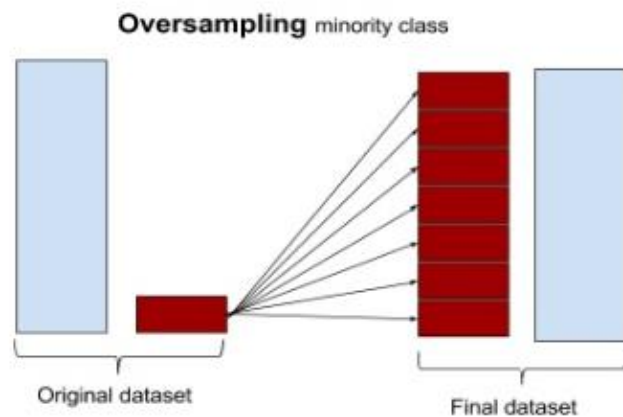  - Lack of interpretability (compared to decision trees)

# Classification of Class-Imbalanced Data Sets

- Traditional methods assume a balanced distribution of classes and equal error costs. But in real world situations, we may face imbalanced data sets, which has rare positive examples but numerous negative ones.

  - Medical diagnosis: Medical screening for a condition

    is usually performed on a large population of people

    without the condition, to detect a small minority

    with it (e.g., HIV prevalence in the USA is ~0.4%)

  - Fraud detection: About 2% of credit card accounts

  are defrauded per year. (Most fraud detection domains are heavily imbalanced.)

  - Product defect, accident (oil-spill), disk drive failures, etc.

# Classification of Class-Imbalanced Data Sets

❑ Typical methods on imbalanced data (Balance the training set)
   ❑ **Oversampling**: Oversample the minority class.
   ❑ **Under-sampling**: Randomly eliminate tuples from majority class
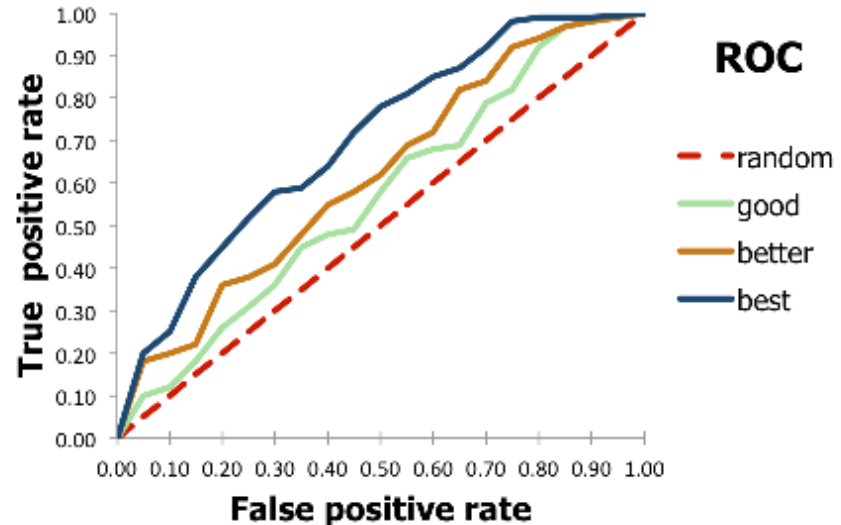   ❑ **Synthesizing:** Synthesize new minority classes

# Classification of Class-Imbalanced Data Sets

❑ Typical methods on imbalanced data (At the algorithm level)

❑ **Class weight adjusting**: Since false negative costs more than false positive, we can give larger weight to false negative

❑ **Ensemble techniques**: Ensemble multiple classifiers introduced in the following chapter

# Evaluate imbalanced data classifier

- Can we use Accuracy to evaluate imbalanced data classifier?

- Accuracy simply counts the number of errors. If a data set has 2% positive labels and 98% negative labels, a classifier that map all inputs to negative class will get an accuracy of 98%!

- ROC Curve

# Summary

- Classification: Model construction from a set of training data

- Effective and scalable methods

  - Decision tree induction, Bayes classification methods, lazy learners, linear classifier

  - No single method has been found to be superior over all others for all data sets

- Evaluation metrics: Accuracy, sensitivity, specificity, precision, recall, $F$ measure

- Model evaluation: Holdout, cross-validation, bootstrapping, ROC curves (AUC)

- Techniques to improve classification accuracy: Ensemble Methods (bagging, boosting, random forests), imbalanced classification

# References (1)

- C. Apte and S. Weiss. **Data mining with decision trees and decision rules**. Future Generation Computer Systems, 13, 1997

- C. M. Bishop,  **Neural Networks for Pattern Recognition**.  Oxford University Press, 1995

- L. Breiman, J. Friedman, R. Olshen, and C. Stone. **Classification and Regression Trees**. Wadsworth International Group, 1984

- C. J. C. Burges. **A Tutorial on Support Vector Machines for Pattern Recognition**. *Data Mining and Knowledge Discovery*, 2(2): 121-168, 1998

- P. K. Chan and S. J. Stolfo. **Learning arbiter and combiner trees from partitioned data for scaling machine learning**. KDD'95

- H. Cheng, X. Yan, J. Han, and C.-W. Hsu, **Discriminative Frequent Pattern Analysis for Effective Classification**, ICDE'07

- H. Cheng, X. Yan, J. Han, and P. S. Yu, **Direct Discriminative Pattern Mining for Effective Classification**, ICDE'08

- W. Cohen.  **Fast effective rule induction**. ICML'95

- G. Cong, K.-L. Tan, A. K. H. Tung, and X. Xu.  **Mining top-k covering rule groups for gene expression data**.  SIGMOD'05

# References (2)

- A. J. Dobson. **An Introduction to Generalized Linear Models**. Chapman & Hall, 1990.

- G. Dong and J. Li. **Efficient mining of emerging patterns: Discovering trends and differences**. KDD'99.

- R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley, 2001

- U. M. Fayyad. **Branching on attribute values in decision tree generation**. AAAI'94.

- Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting**. J. Computer and System Sciences, 1997.

- J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets**. VLDB'98.

- J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction**. SIGMOD'99.

- T. Hastie, R. Tibshirani, and J. Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction.** Springer-Verlag, 2001.

- D. Heckerman, D. Geiger, and D. M. Chickering. **Learning Bayesian networks: The combination of knowledge and statistical data**. Machine Learning, 1995.

- W. Li, J. Han, and J. Pei, **CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules**, ICDM'01.

# References (3)

- T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of  thirty-three old and new classification algorithms.**  Machine Learning, 2000.

- J. Magidson.  **The Chaid approach to segmentation modeling:  Chi-squared automatic interaction detection**. In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, Blackwell Business, 1994.

- M. Mehta, R. Agrawal, and J. Rissanen. **SLIQ : A fast scalable classifier for data mining**. EDBT'96.

- T. M. Mitchell. **Machine Learning**. McGraw Hill, 1997.

- S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey**, Data Mining and Knowledge Discovery 2(4): 345-389, 1998

- J. R. Quinlan. **Induction of decision trees**. *Machine Learning*, 1:81-106, 1986.

- J. R. Quinlan and R. M. Cameron-Jones. **FOIL: A midterm report**. ECML'93.

- J. R. Quinlan. **C4.5: Programs for Machine Learning**. Morgan Kaufmann, 1993.

- J. R. Quinlan.  **Bagging, boosting, and c4.5**. AAAI'96.

# References (4)

- R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning**. VLDB'98.

- J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining**. VLDB'96.

- J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning**. Morgan Kaufmann, 1990.

- P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining**. Addison Wesley, 2005.

- S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems**. Morgan Kaufman, 1991.

- S. M. Weiss and N. Indurkhya. **Predictive Data Mining**. Morgan Kaufmann, 1997.

- I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques**, 2ed. Morgan Kaufmann, 2005.

- X. Yin and J. Han. **CPAR: Classification based on predictive association rules**. SDM'03

- H. Yu, J. Yang, and J. Han. **Classifying large data sets using SVM with hierarchical clusters**. KDD'03.