

# **Data Wrangling and Data Analysis**

## **Regression, Classification and Evaluation Techniques**

**Hakim Qahtan**

Department of Information and Computing Sciences

Utrecht University



# Topics for This Week

- Regression
- Classification
- Evaluation



# Data Analytics Models

- Prediction vs Inference
- Regression vs Classification
- Supervised vs Unsupervised



# Prediction vs. Inference

$$\hat{y} = \hat{f}(x)$$

- Prediction
  - Main interest is in the outcome value  $\hat{y}$
  - Use the model to predict the outcomes for new data points.
- Inference
  - Main interest is in  $\hat{f}$
  - Use training data to learn a model
  - Association of predictors
  - Types of relations



# Regression vs. Classification

- Regression

always continuous/numerical data

- Algorithms attempt to estimate the mapping function  $f$  from the input variables  $x$  to numerical or continuous output variables  $y$
- Given a dataset about house prices – predict the price of a given house

- Classification

- Algorithms attempt to estimate the mapping function  $f$  from the input variables  $x$  to discrete or categorical output variables  $y$
- Houses dataset – predict if the selling price is more or less than the recommended price

binary/classes - putting it into category



# Supervised vs. Unsupervised

- Supervised models
  - Training a model using data which is well "labeled"
  - Supervised learning algorithm learns from labeled training data to predict outcomes for unforeseen data
- Unsupervised models
  - The model works on its own to discover information
  - Deals with unlabeled data and looks for unknown patterns in data
- There are also Semi-Supervised Models
  - Small amount of labeled data and large amount of unlabeled data
  - E.g. A teacher provides set of solved and unsolved problem for exam preparation



# Regression

---

# Regression

- Given the values of inputs  $X$  and the corresponding output  $Y$  belongs to the set of real values  $R$ , predict output accurately for new input.
- Formally:
  - Given:
    - A set of  $N$  observations  $\{x_n\}_{n=1\dots N}$  with their corresponding target values  $\{y_n\}_{n=1\dots N}$
  - Goal:
    - Predict the value of  $y_{n+1}$  for a give  $x_{n+1}$
- Predictive technique where the target variable to be estimated is continuous





# Regression (Cont.)

- Let  $D$  denote a dataset containing  $N$  observations,  
$$D = \{(x_i, y_i) | i = 1, 2, \dots, N\}$$
  - $x_i$  corresponds to the values of attributes of the  $i$ -th observation.
    - These are called **explanatory variables** and can be discrete or continuous.
  - $y_i$  corresponds to the target variable.
- Target: find a function that can minimize the error between the predicted and the actual values
  - The error can be measured as the **sum of absolute or squared error**

$$\begin{aligned} \text{sum absolute error (SAE)} &= \sum_i |y_i - f(x_i)| \\ \text{sum squared error (SSE)} &= \sum_i (y_i - f(x_i))^2 \end{aligned}$$

The sum of squares error, or SSE. The error is the difference between the observed value and the predicted value.

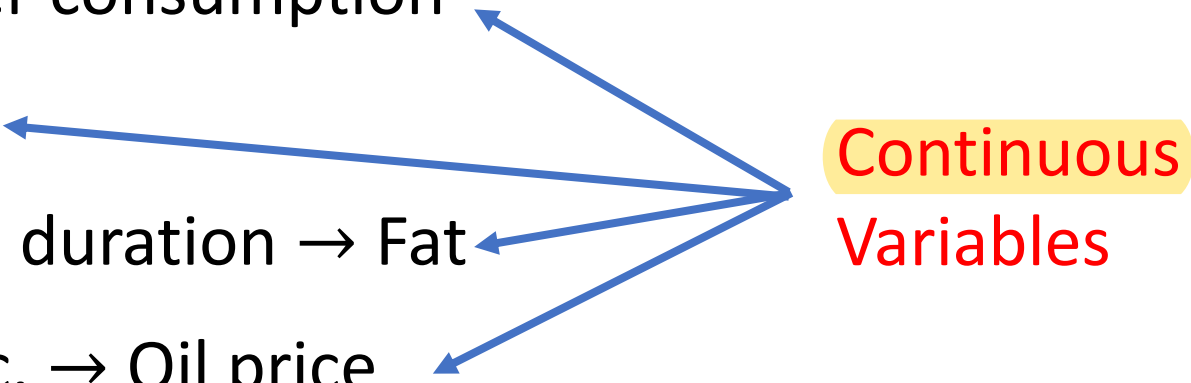


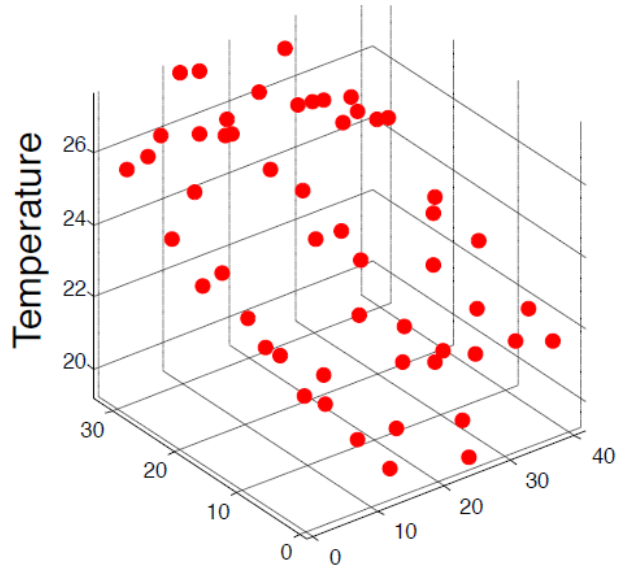
# What is Regression?

- Regression: making predictions about real-world quantities.
  - What would be the price of a product after producing a new version?
  - How the discount will affect the sales volume?
  - How the weather will affect the sales of the restaurants?
  - How many students are expected to show up in a lecture?
  - ...
- Regression is a supervised machine learning model



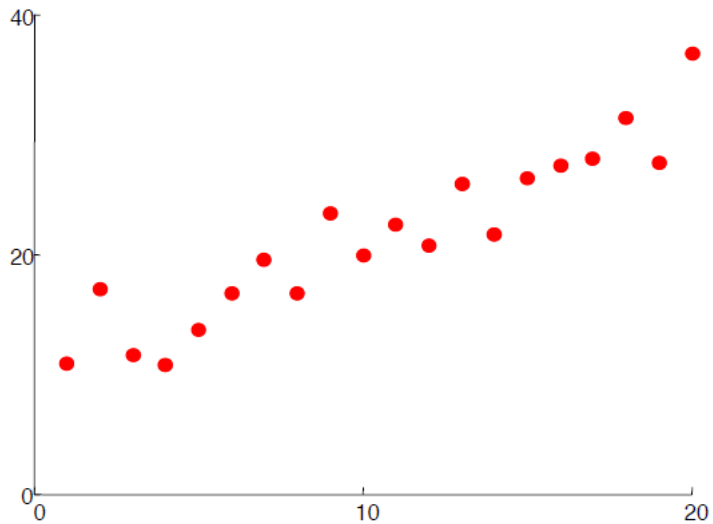
# Regression – More Examples

- Processes, memory → Power consumption
  - Protein structure → Energy
  - Heart-beat rate, age, speed, duration → Fat
  - Oil supply, consumption, etc. → Oil price
  - ...
  - Definition: Regression is the task of learning a target function  $f$  that maps each attribute set  $x$  into a continuous-valued output  $y$ .
- 
- The diagram illustrates the concept of regression by showing four examples of input-output pairs. Each example is represented by a bullet point. Four blue arrows originate from the right side of each bullet point and converge towards a central yellow box containing the text 'Continuous Variables' in red. This visualizes that all these different tasks are instances of regression, where the goal is to predict a continuous output from a set of input attributes.

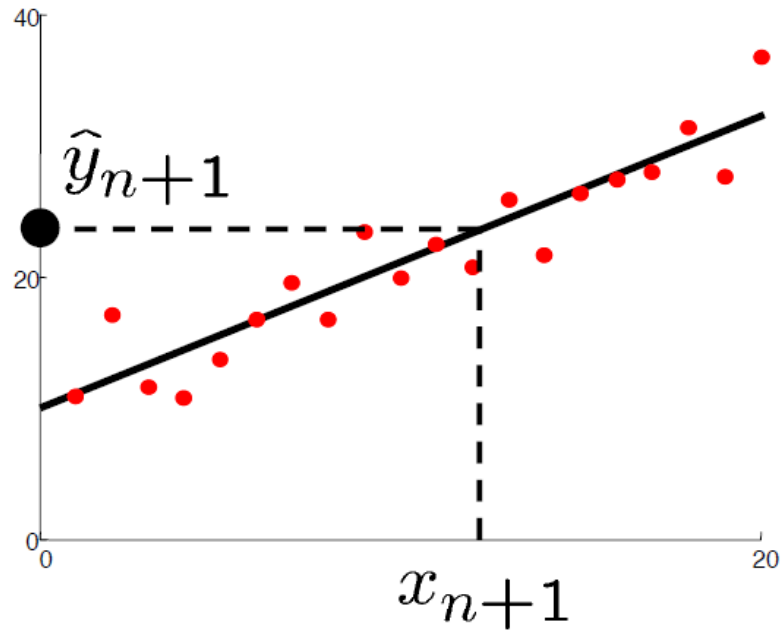


# Regression – Examples

- Given  $\{x_n, y_n\}_{n=1 \dots N}$
- Predict the value of  $y_{n+1}$  for a given  $x_{n+1}$

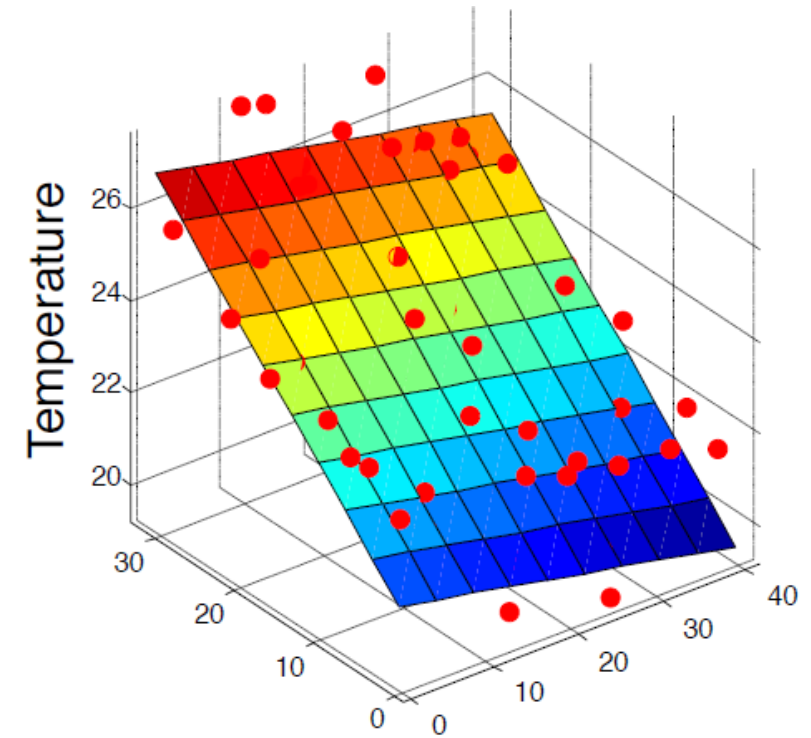


# Regression – Examples



- Prediction  $\hat{y}_i = w_0 + w_1 x_i$

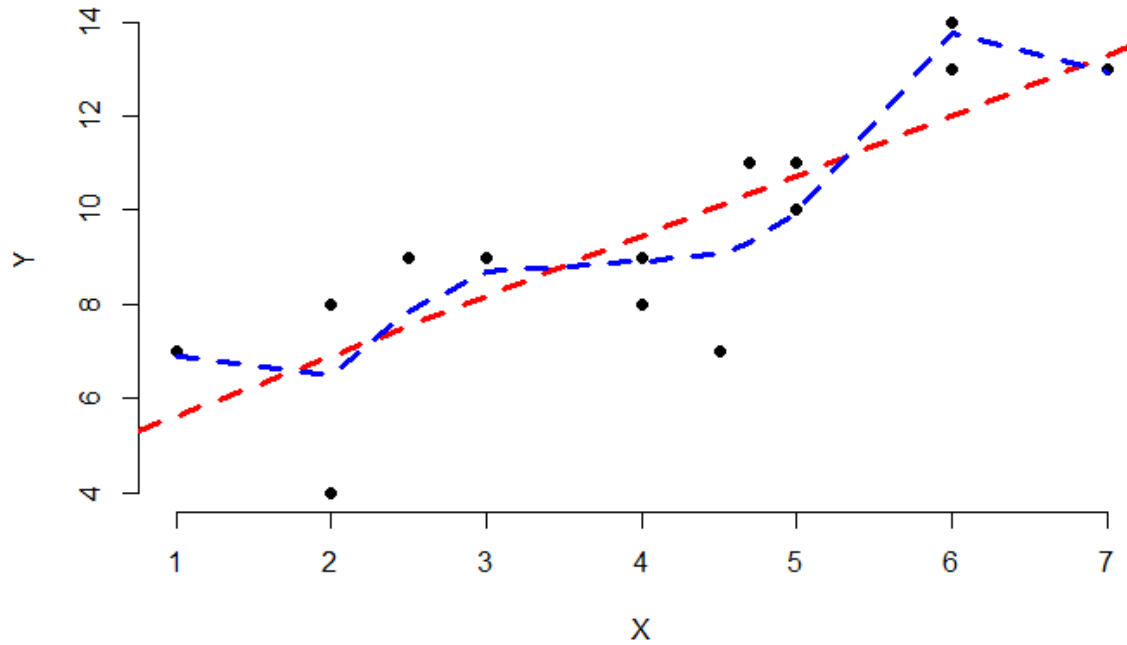
- This is called linear regression as the function is linear in the parameters  $w_0$ ,  $w_1$  and  $w_2$ .



- Prediction  $\hat{y}_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2}$

$$= (1 \quad x_{i,1} \quad x_{i,2}) \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$$

# Simple Linear Regression



Given a set of points  $(x_i, y_i)$  such as the points in the scatterplot, find the best fitting line

$$f(x_i) = \omega_0 + \omega_1 x_i$$

such that:

$$SSE = \sum_i (y_i - f(x_i))^2$$

$$= \sum_i (y_i - \omega_0 - \omega_1 x_i)^2$$

is minimized



# Simple Linear Regression (Cont.)

- The above optimization problem can be solved by:
  1. Taking the partial derivatives of  $SSE$  with respect to  $\omega_0$  and  $\omega_1$
  2. Setting  $\frac{\partial SSE}{\partial \omega_0}$  and  $\frac{\partial SSE}{\partial \omega_1}$  to 0
  3. Solving the system of linear equations

$$\text{Since: } SSE = \sum_i (y_i - \omega_0 - \omega_1 x_i)^2$$

$$\text{Then } \frac{\partial SSE}{\partial \omega_0} = -2 \sum_i (y_i - \omega_0 - \omega_1 x_i) = 0$$

$$\text{And } \frac{\partial SSE}{\partial \omega_1} = -2 \sum_i x_i (y_i - \omega_0 - \omega_1 x_i) = 0$$



# Simple Linear Regression (Cont.)

- The equations can be summarized by the normal equation:

$$\begin{pmatrix} N & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{pmatrix} \begin{pmatrix} \omega_0 \\ \omega_1 \end{pmatrix} = \begin{pmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{pmatrix}$$





# Example

- Consider the following dataset

$x$	1	2	2.5	2	4	4	4	4.5	4.7	5	5	6	6	7
$y$	7	4	9	8	9	8	9	7	11	11	10	13	14	13

$$\sum_i x_i \overset{\text{sum of all } x}{=} 57.7$$

$$\sum_i y_i = 133$$

$$\sum_i x_i^2 \overset{1^2+2^2+2.5^2+\dots 7^2}{=} 276.59$$

$$\sum_i x_i y_i = 589.7$$

## Example (Cont.)

$x$	1	2	2.5	2	4	4	4	4.5	4.7	5	5	6	6	7
$y$	7	4	9	8	9	8	9	7	11	11	10	13	14	13

$$\sum_i x_i = 57.7 \quad \sum_i x_i^2 = 276.59 \quad \sum_i y_i = 133 \quad \sum_i x_i y_i = 589.7$$

$$\begin{pmatrix} 14 & 57.7 \\ 57.7 & 276.59 \end{pmatrix} \begin{pmatrix} \omega_0 \\ \omega_1 \end{pmatrix} = \begin{pmatrix} 133 \\ 598.7 \end{pmatrix}$$

By solving the equations we get:

$$\omega_0 = 4.1282 \text{ and } \omega_1 = 1.3033$$

Hence:

$$f(x_i) = 1.3033 x_i + 4.1282$$

# Simple Linear Regression (Cont.)

- A general solution for the normal equation can be found as follows:

$$\omega_0 = \bar{y} - \omega_1 \bar{x}$$

and

$$\omega_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \quad (\text{explanatory var.} - \text{average value})$$

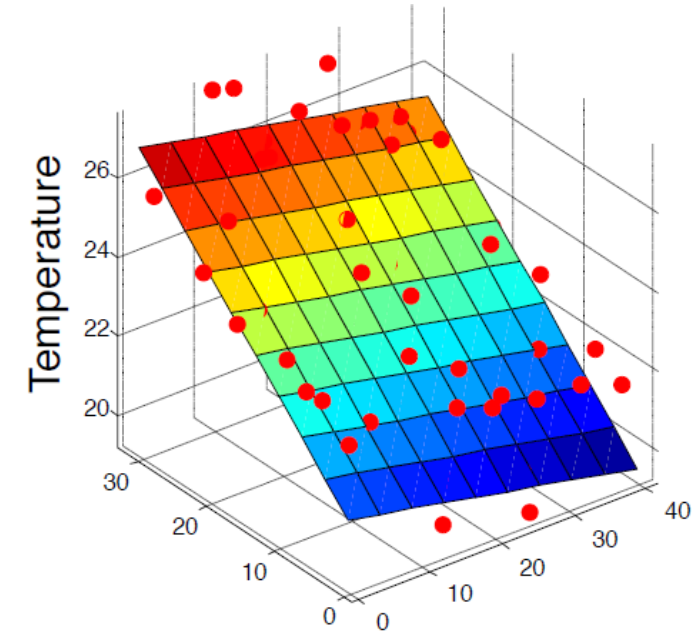
where

$\bar{x}, \bar{y}$  are the mean (average) values for the vectors  $x, y$



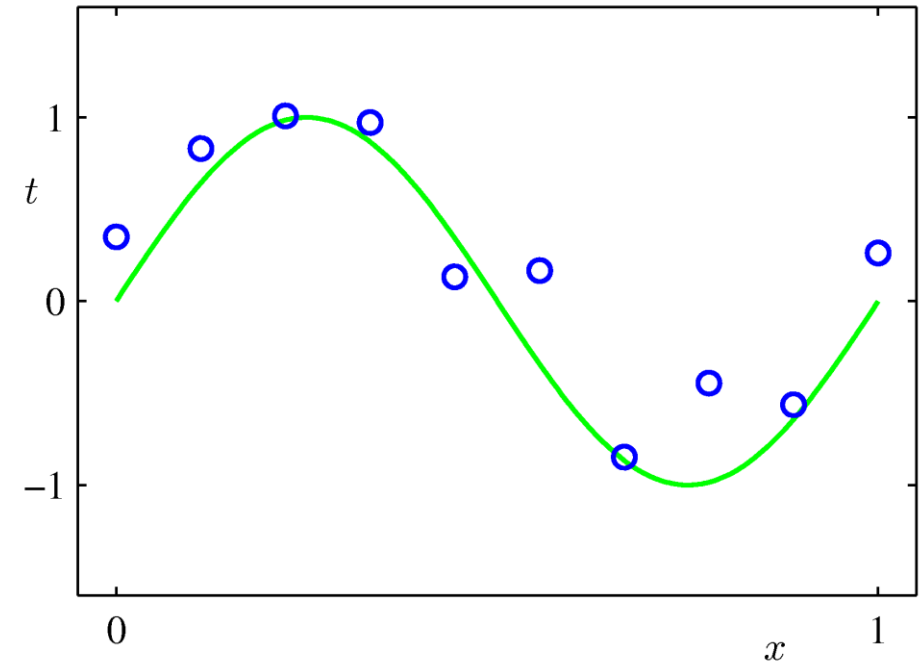
# Multiple Linear Regression

- Fitting d-dimensional hyperplane to the d variables
- Simple, yet powerful
- If the relation between the response and independent variables is non-linear, we can use non-linear transformation of the variables
- E.g.  $x_i^2$  instead of  $x_i$



# Polynomial Regression

- Suitable when the relationship between the response and the independent variables is non-linear
- Higher order polynomials complicate the model
- May cause model overfitting
- Increase the computational complexity
- Case Study: Predicting the price of a new housing market – check the provided notebook



# Evaluating the Goodness of a Fit

- To evaluate how well data points fit to a line, we use the  $R^2$  which is defined as:

$$R^2 = \frac{\sum_i (f(x_i) - \bar{y})^2}{\sum_i (y_i - \bar{y})^2}$$

- $R^2$  takes values in the interval  $[0,1]$  where the values close to 1 means that the data fits well to the regression line
- When adding more explanatory variables, the value of  $R^2$  increases so it is adjusted using the formula

$$\text{Adjusted } R^2 = 1 - \left( \frac{N - 1}{N - d} \right) (1 - R^2)$$

to adjust the increased value we put the number of parameters in consideration by adding  $d$

where  $N$  is the number of data points and  $d + 1$  is the number of parameters of the regression model



# Classification

---

# Classification

- Given an input vector of the data  $x$ , the goal is to assign  $x$  to one of the  $K$  discrete classes.
  - Classes are taken to be disjoint.
- Classification techniques:
  - Logistic Regression
  - Bayesian classification
  - Learning from the neighbors
  - Decision tree-based methods
  - Rule-Based classification
  - Neural networks
  - Support vector machines





# Classification vs. 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

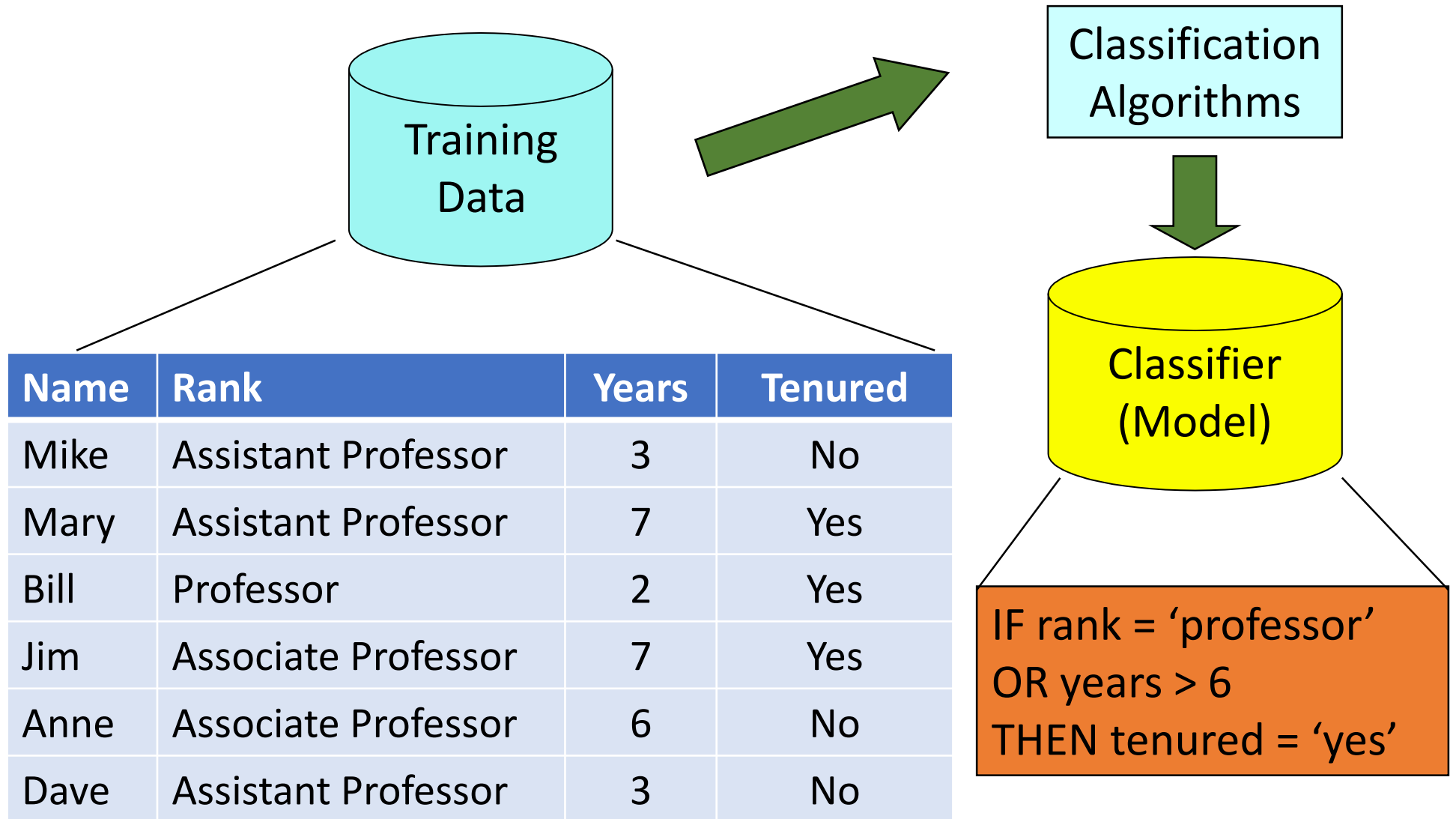


# Classification – A Two-Step Process

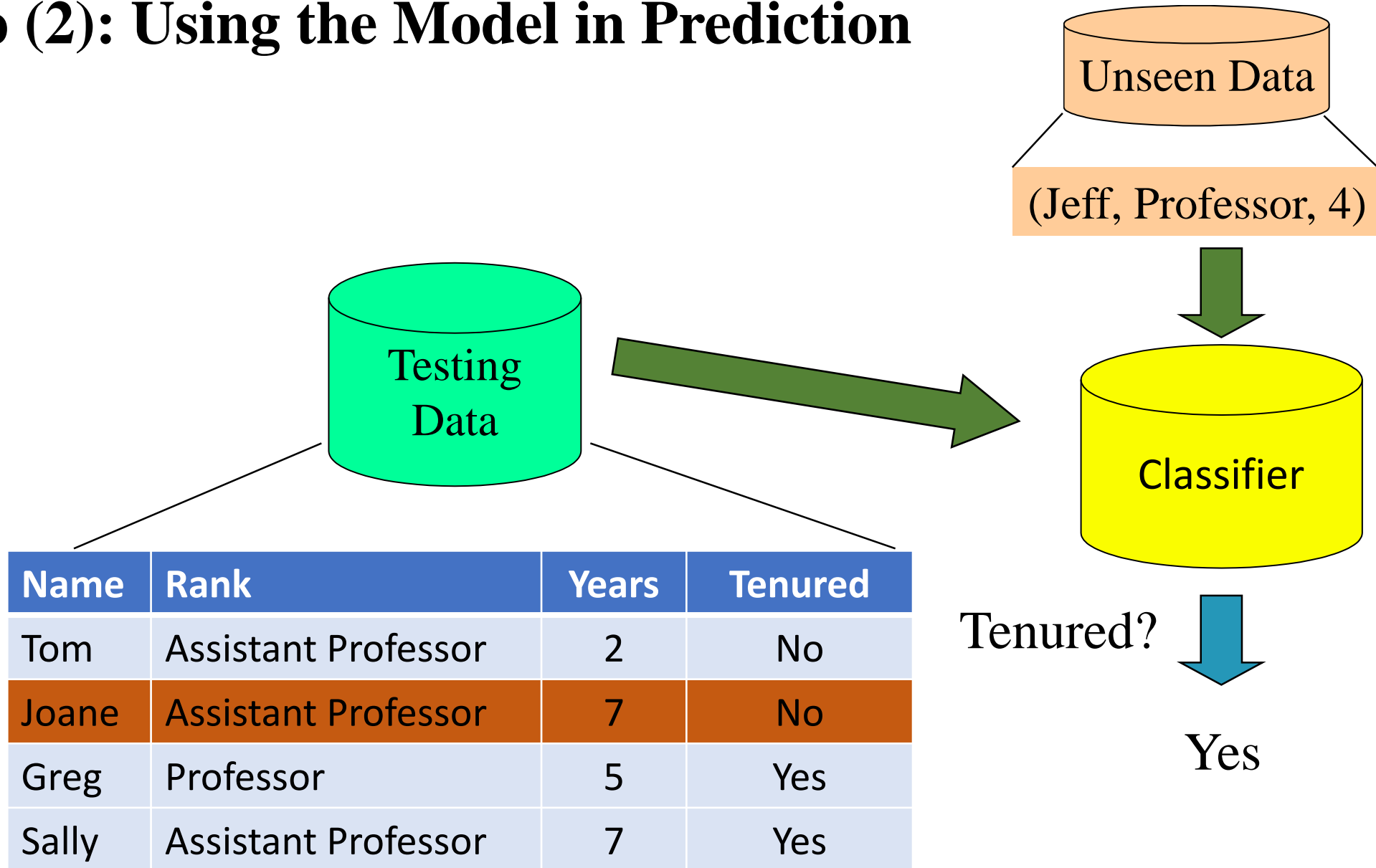
- Model construction: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
  - The set of tuples used for model construction is training set
  - The model is represented as classification rules, decision 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



# Step (1): Model Construction



## Step (2): Using the Model in Prediction



# **Classification**

## **Logistic Regression**

---

# Classification – Logistic Regression

- Useful when the target is binary
- Logistic regression is a type of probabilistic statistical classification model
- It measures the relationship between the dependent (target) binary variable and the independent explanatory variables
- We have a binary target variable  $Y$ , and we want to model the conditional probability  $P(Y = 1 \mid X = x)$  as a function  $p(x)$  of the explanatory variables  $x$ .
- Any unknown parameters (recall  $\omega_0$  and  $\omega_1$ ) are estimated by maximum likelihood.



# Logistic Regression (Cont.)

- Let  $p(x)$  be a linear function
  - We are estimating a probability, which must be between 0 and 1
  - Linear functions are unbounded, so this approach doesn't work
- Better idea: Set the odds ratio to a linear function:

$$\log(odds) = \text{logit}(p) = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

Solving for  $p$ :

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

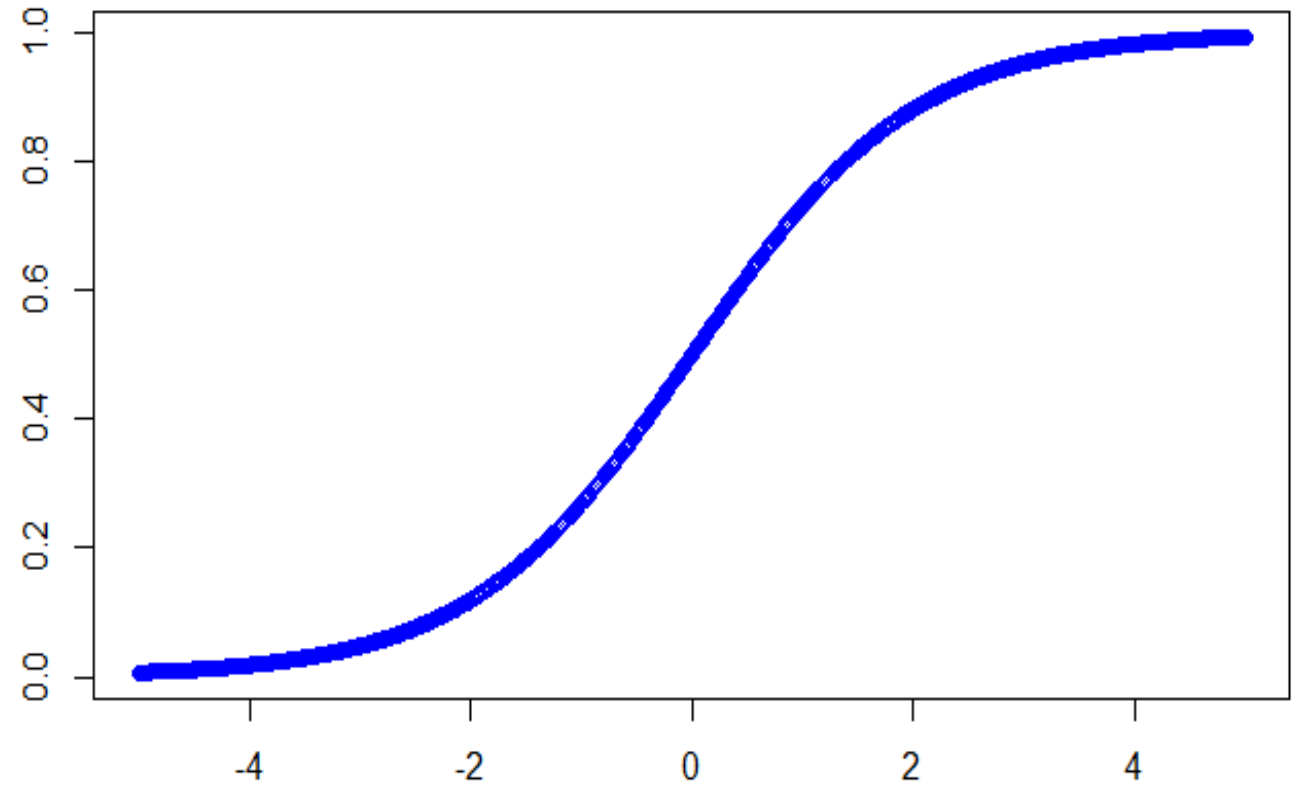
- This is called the logistic (logit) function and it takes values in the interval  $[0,1]$



# Logistic Curve

- A sigmoid function that assumes values in the range  $[0,1]$

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$
$$= \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

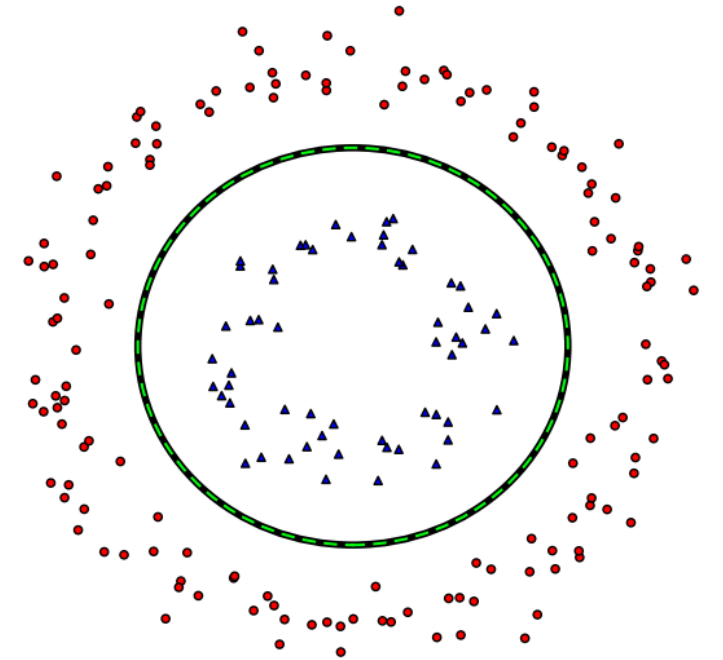
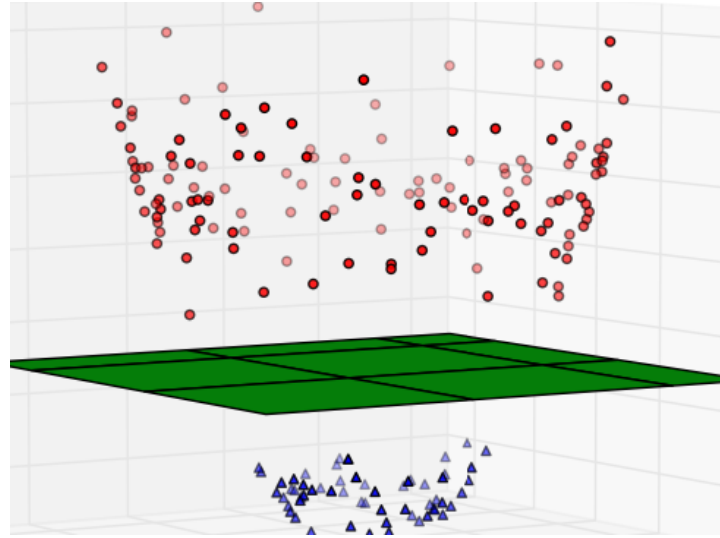
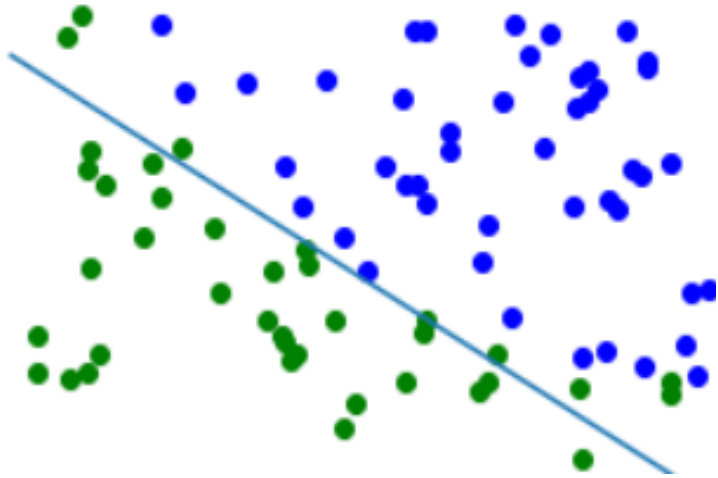




# Logistic Regression

- To minimize misclassification rates, we predict:
  - $Y = 1$  when  $p(x) \geq 0.5$  and  $Y = 0$  when  $p(x) < 0.5$
  - So  $Y = 1$  when  $\beta_0 + \beta_1 x$  is non-negative and 0 otherwise
- Logistic regression gives us a linear classifier where the decision boundary separating the two classes is the solution of  $\beta_0 + \beta_1 x = 0$





# Decision Boundaries

# Logistic Regression

- The parameters  $\beta_0, \beta_1, \dots$  are estimated using a technique called Maximum Likelihood Estimation (MLE)
  - Unlike the least squares methods used for Linear regression, finding a closed form for the coefficients using MLE is not possible. Instead, an iterative process (e.g., Newton's method) is used.
  - This process begins with a tentative solution, revises it slightly to see if it can be improved, and repeats this revision until improvement is very small, at which point the process is said to have converged.



# **Classification**

## **Bayesian Classification**

---

# Bayesian Classification

- 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



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

# Bayes' Theorem

- Total probability Theorem:  $p(B) = \sum_{i=1}^M p(B|A_i)p(A_i)$
- Bayes' Theorem:  $p(H|X) = p(X|H) \times p(H)/p(X)$ 
  - 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  $\mathbf{X}$ , posteriori probability of a hypothesis  $H$ ,  $p(H|\mathbf{X})$ , follows the Bayes' theorem

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

- Informally, this can be viewed as

$$\textit{posteriori} = \textit{likelihood} \times \textit{prior/evidence}$$

- Predicts  $\mathbf{X}$  belongs to  $C_i$  iff. the probability  $p(C_i|\mathbf{X})$  is the highest among all the  $p(C_k|\mathbf{X})$  for all the  $k$  classes
- Practical difficulty: it requires initial knowledge of many probabilities, involving significant computational cost





# Derive the Maximum Posteriori

- Let  $D$  be a training set of tuples and their associated class labels, and each tuple is represented by an  $n$  attribute vector  $\mathbf{X} = (x_1, x_2, \dots, x_n)$
- Suppose there are  $k$  classes  $C_1, C_2, \dots, C_k$ .
- Classification is to derive the maximum posteriori, i.e., the maximal  $p(C_i|\mathbf{X})$
- This can be derived from Bayes' theorem

$$p(C_i|\mathbf{X}) = p(\mathbf{X}|C_i) \times p(C_i)/p(\mathbf{X})$$

- Since  $p(\mathbf{X})$  is constant for all classes, only

$$p(C_i|\mathbf{X}) = p(\mathbf{X}|C_i) \times p(C_i)$$

needs to be maximized



# Naïve Bayes Classifier

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

$$p(\mathbf{X}|C_i) = \prod_{j=1}^n p(x_j|C_i) = p(x_1|C_i) \times p(x_2|C_i) \times \cdots \times p(x_n|C_i)$$

- This greatly reduces the computation cost: only counts the class distribution
- If  $A_m$  is categorical,  $p(x_m|C_i)$  is the # of tuples in  $C_i$  having value  $x_m$  for  $A_m$  divided by  $|C_{i,D}|$  (# of tuples of  $C_i$  in  $D$ )
- If  $A_m$  is continuous-valued,  $p(x_m|C_i)$  is usually computed based on Gaussian distribution with a mean  $\mu$  and standard deviation  $\sigma$ :  $g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$   
and  $p(x_m|C_i)$  is  $p(\mathbf{X}|C_i) = g(x_m, \mu_{C_i}, \sigma_{C_i})$



# Naïve Bayes Classifier: Example

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

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

Compute  $p(\mathbf{X}|C_i)$  for each class

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

$$p(\text{age} = \text{"<= 30"} \mid \text{buys\_computer} = \text{"no"}) = 0.6$$

$$p(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"yes"}) = 0.444$$

$$p(\text{income} = \text{"medium"} \mid \text{buys\_computer} = \text{"no"}) = 0.4$$

$$p(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"yes"}) = 0.667$$

$$p(\text{student} = \text{"yes"} \mid \text{buys\_computer} = \text{"no"}) = 0.2$$

$$p(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"yes"}) = 0.667$$

$$p(\text{credit\_rating} = \text{"fair"} \mid \text{buys\_computer} = \text{"no"}) = 0.4$$

- $\mathbf{X} = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit\_rating} = \text{fair})$
- $p(\mathbf{X}|C_i) : p(\mathbf{X}|\text{buys\_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$   
 $p(\mathbf{X}|\text{buys\_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$
- $p(\mathbf{X}|C_i) \times p(C_i) : p(\mathbf{X}|\text{buys\_computer} = \text{"yes"}) \times p(\text{buys\_computer} = \text{"yes"}) = 0.028$   
 $p(\mathbf{X}|\text{buys\_computer} = \text{"no"}) \times p(\text{buys\_computer} = \text{"no"}) = 0.007$
- Therefore,  $\mathbf{X}$  belongs to class ("buys\_computer = yes")



# Avoiding the Zero-Probability Problem

- Naïve Bayesian prediction requires each conditional probability be non-zero. Otherwise, the predicted probability will be zero since:

$$p(\mathbf{X}|C_i) = \prod_{j=1}^n p(x_j|C_i)$$

- Ex. Suppose a dataset with 1000 tuples, income=low (0), income=medium (990), and income = high (10)
- Use Laplacian correction (or Laplacian estimator)
  - Add 1 to each case

$$p(\text{income} = \text{low}) = 1/1003, p(\text{income} = \text{medium}) = 991/1003, p(\text{income} = \text{high}) = 11/1003$$

- The “corrected” prob. estimates are close to their “uncorrected” counterparts



# Naïve Bayes Classifier: Final Remarks

- 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



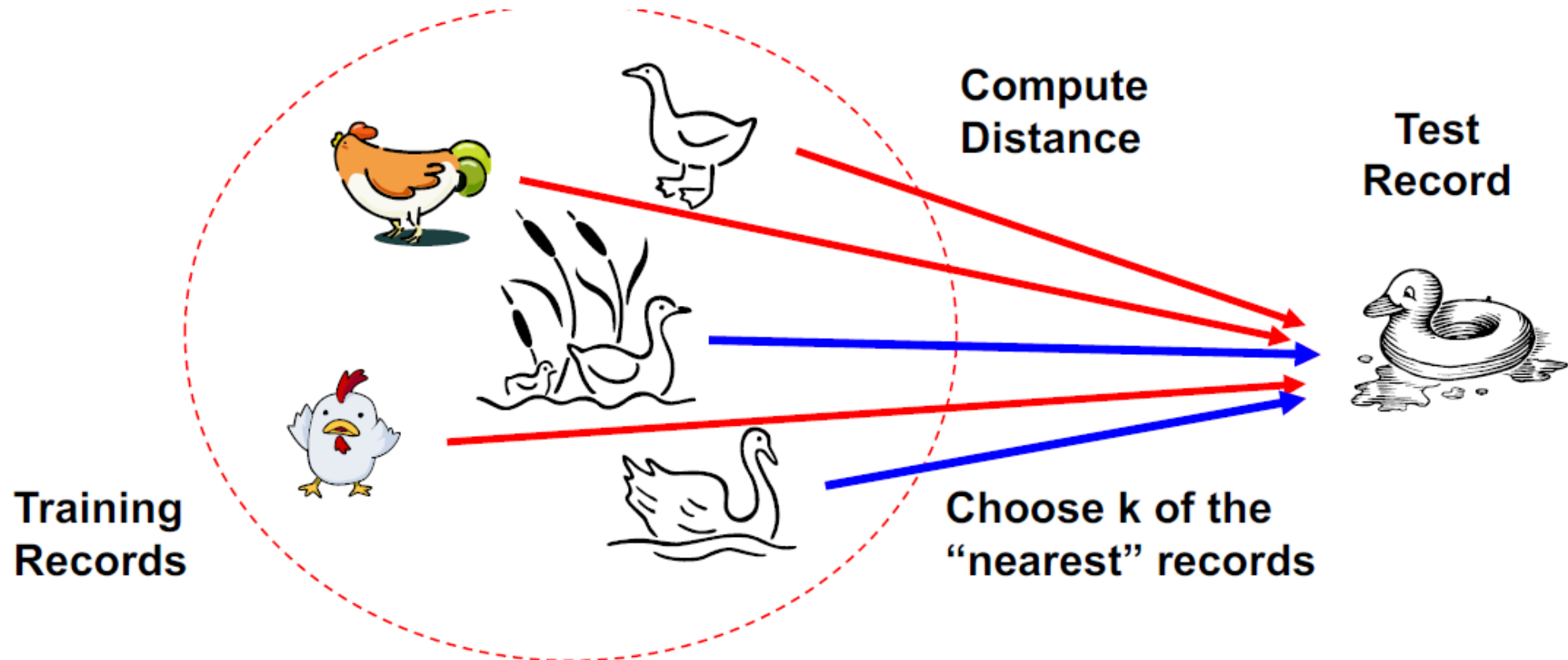
# **Classification**

## **Nearest Neighbors Classifiers**

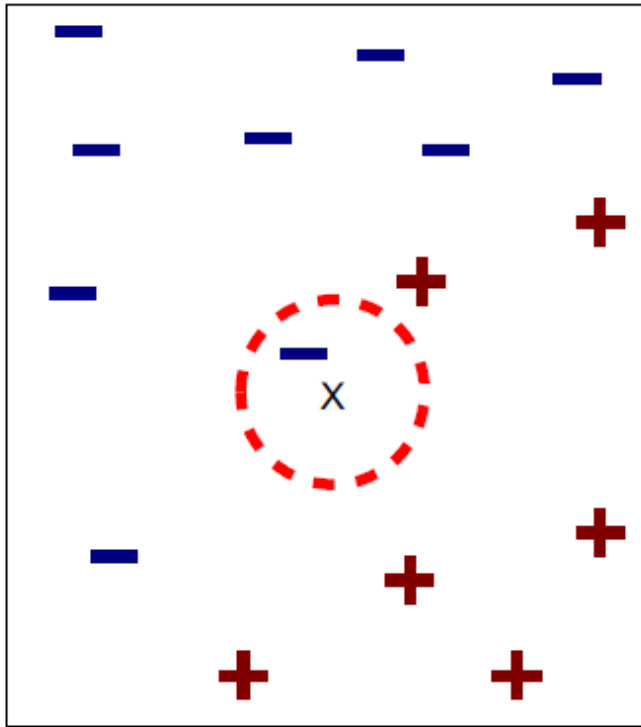
---

# Nearest Neighbors Classifiers

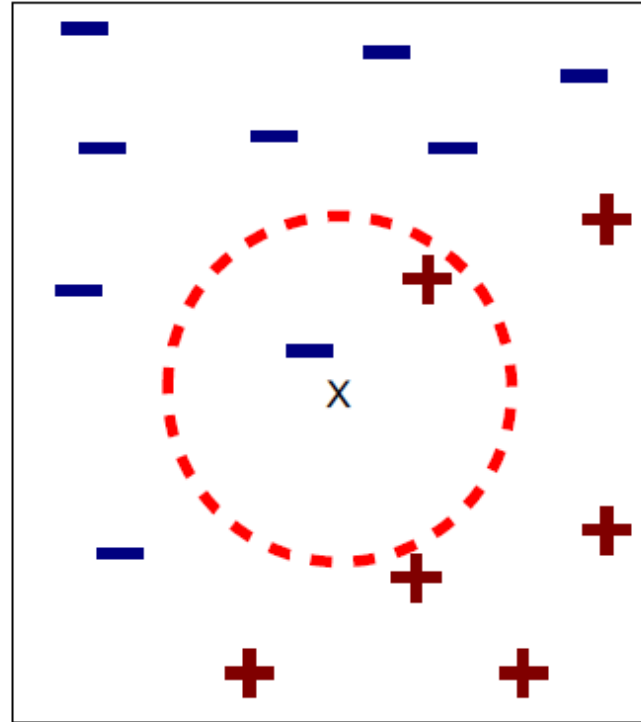
Basic idea: if it walks like a duck, quacks like a duck, then it is probably a duck



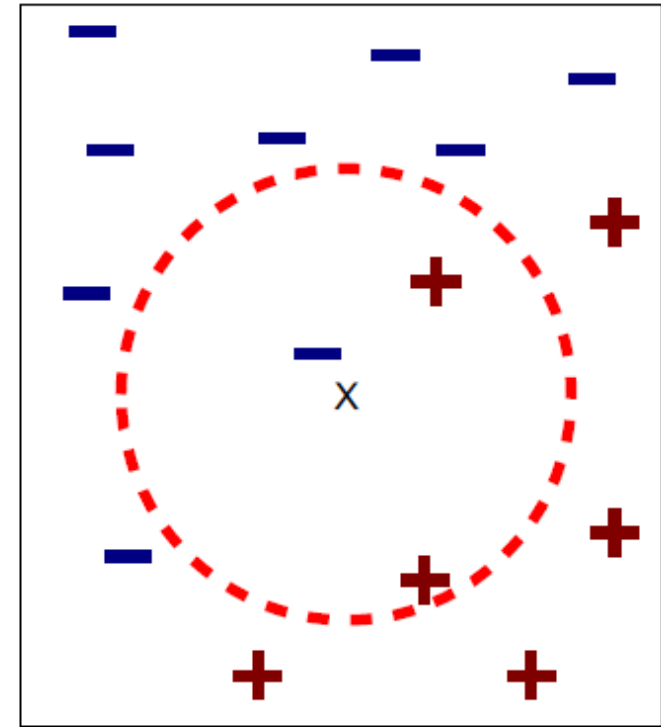
# Nearest Neighbors Classifiers (Cont.)



1- nearest neighbor



2- nearest neighbors

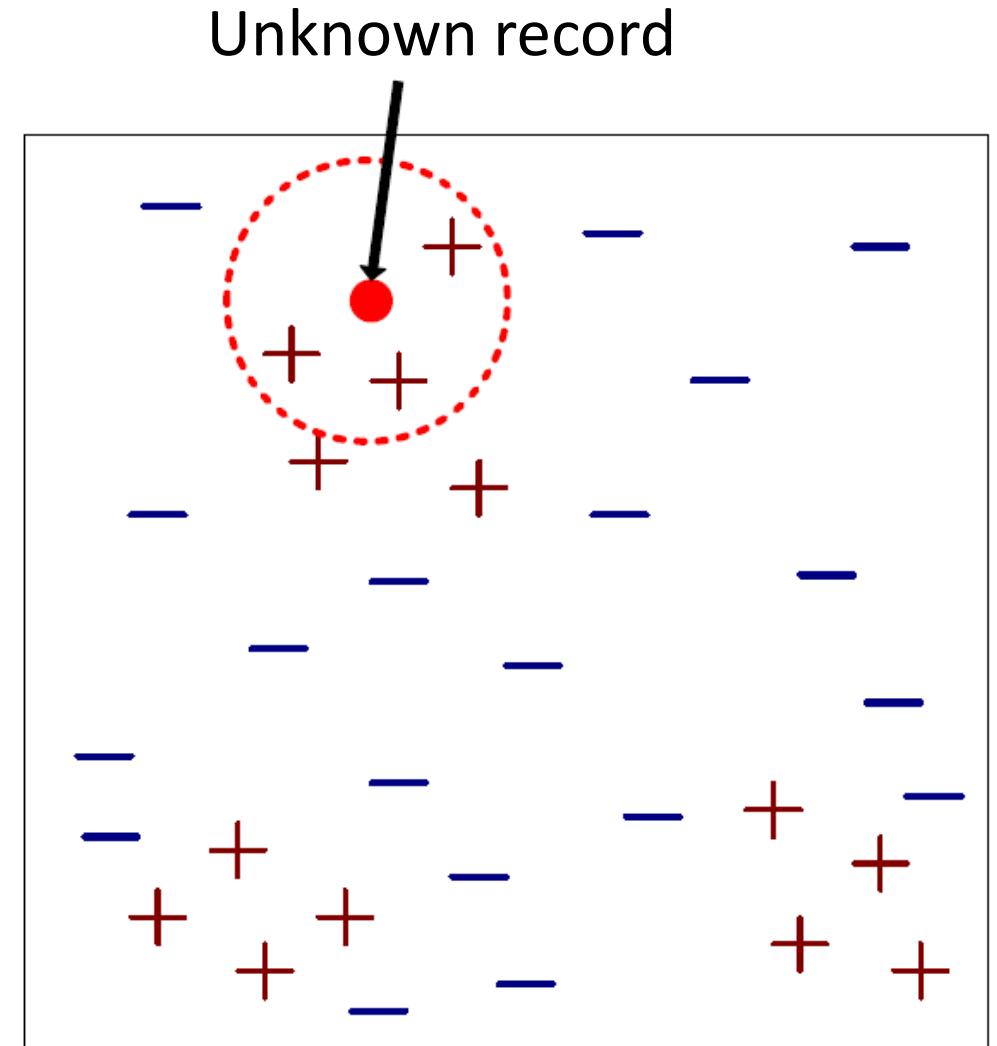


3- nearest neighbors



# Nearest Neighbors Classifiers (Cont.)

- Three requirements:
  - Set of records (training set)
  - Distance metric
  - The number of neighbors to be considered  $k$
- Classifying unknown record  $x$ :
  - Compute the distance from  $x$  to the other training records
  - Identify the  $k$ -Nearest Neighbors ( $kNN(x)$ )
  - Use class labels of the  $kNN$  records to determine the class of  $x$  .. How?



# K Nearest Neighbors Classification

- To determine the class of unknown record  $x$  from the classes of its neighbors:
  - Use the majority vote – is this enough?
  - Weight the vote according to the distance
  - Examples of weight factors,  $w = d^{-2}$ ,  $w = e^{-d^2}$ , etc.
- How to choose  $k$ :
  - $k$  too small, sensitive to noise
  - $k$  too large, neighborhood may include points from other classes



# K Nearest Neighbors Classification (Cont.)

- Scaling issues
  - Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes
- Example:
  - Height of a person may vary from 1.5m to 1.8m
  - Weight of a person may vary from 90lb to 300lb
  - Income of a person may vary from €10K to €1M
- Solution: Normalize the vectors to unit length



# K Nearest Neighbors Classification (Cont.)

- k-NN classifiers are lazy learners
  - It does not build models explicitly
  - Robust to noisy data by averaging k-nearest neighbors unlike eager learners such as decision tree induction and rule-based systems
  - Classifying unknown records are relatively expensive



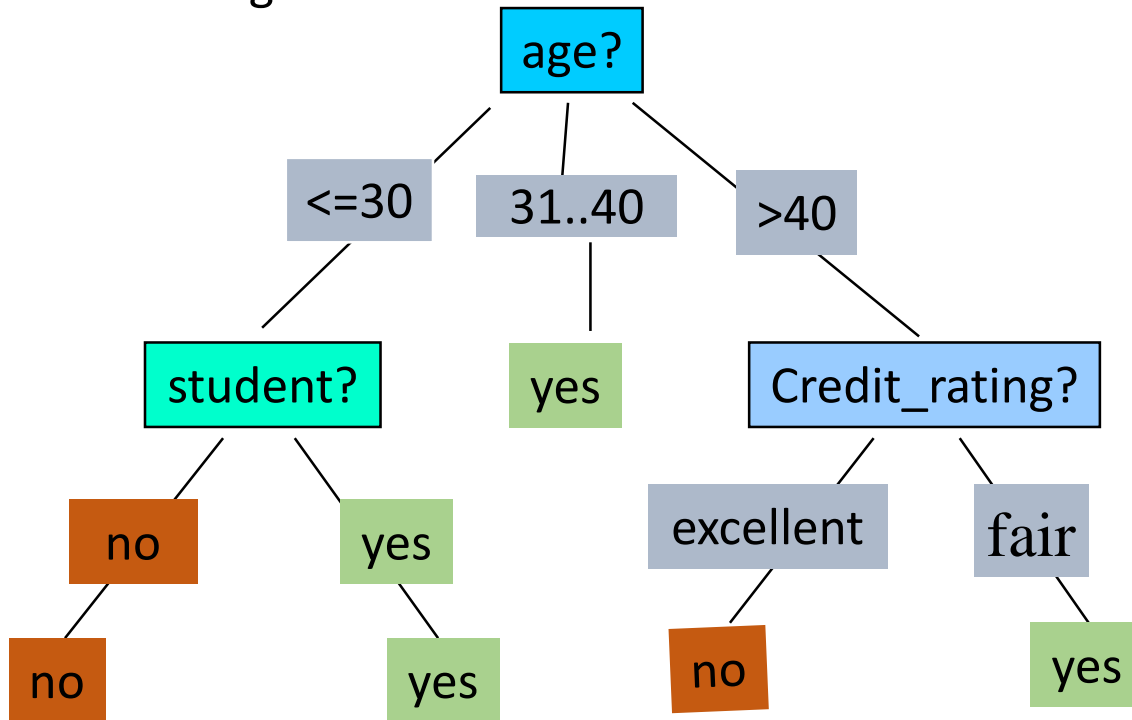
# **Classification**

## **Decision Tree Induction**

---

# Decision Tree Induction

- Training dataset: Buys\_computer
- The dataset 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



# 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



# Review of Entropy

- Entropy (from information theory):
  - A measure of uncertainty associated with a random variable
  - Computing the entropy: for a discrete random variable  $Y$  that takes  $m$  distinct values  $\{y_1, y_2, \dots, y_m\}$ :

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

- Interpretation:
  - Higher entropy implies higher uncertainty and vice versa
- Conditional entropy:  $H(Y|X) = \sum_x p(x)H(Y|X = x)$





# Attribute Selection Measure: Information Gain

- Select the attribute with the highest information gain
- Let  $p_i$  be the probability that an arbitrary tuple in a dataset  $D$  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  $A$

$$gain(A) = info(D) - info_A(D)$$



# Information Gain – Example

- Class  $C_1$ : buys\_computer = “yes”
- Class  $C_2$ : buys\_computer = “no”

age	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
31...40	4	0	0
$> 40$	3	2	0.971

$$info(D) = I(9,5) = -\frac{9}{14} \log_2 \left( \frac{9}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right) = 0.94$$

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

$$= \frac{5}{14} (0.971) + 0 + \frac{5}{14} (0.971) = 0.694$$

$\frac{5}{14} I(2,3)$  means that “age  $\leq 30$ ” has 5 out of 14 samples with 2 from class  $C_1$  and 3 from class  $C_2$

age	income	student	credit_rating	buys_computer
$\leq 30$	high	no	fair	no
$\leq 30$	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
$\leq 30$	medium	no	fair	no
$\leq 30$	low	yes	fair	yes
$> 40$	medium	yes	fair	yes
$\leq 30$	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
$> 40$	medium	no	excellent	no



# Information Gain – Example

$$\text{gain}(\text{age}) = \text{info}(D) - \text{info}(D_{\text{age}}) = 0.94 - 0.694 = 0.246$$

- Similarly:

$$\text{gain}(\text{income}) = 0.029$$

$$\text{gain}(\text{student}) = 0.151$$

$$\text{gain}(\text{credit\_rating}) = 0.048$$



# Information-Gain for Continuous-Valued Attributes

- Let attribute  $A$  be a continuous-valued attribute
  - Must determine the best split point for  $A$
  - Sort the values in  $A$  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  $A \leq \text{split-point}$ , and  $D2$  is the set of tuples in  $D$  satisfying  $A > \text{split-point}$



# Gain Ratio for Attribute Selection

- 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)

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

$$GainRatio(A) = Gain(A) / SplitInfo_A(D)$$

- 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 dataset  $D$  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 dataset  $D$  is split on  $A$  into two subsets  $D1$  and  $D2$ , the gini index  $gini(D)$  is defined as

$$gini_A(D) = \frac{|D1|}{|D|} gini(D1) + \frac{|D2|}{|D|} gini(D2)$$

- Reduction in Impurity:

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

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



# Computing 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  $D1: \{low, medium\}$  and 4 in  $D2$

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right) gini(D1) + \left(\frac{4}{14}\right) gini(D2)$$

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right) \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \left(\frac{4}{14}\right) \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right) = 0.443$$

$$gini_{income \in \{medium, high\}} = 0.450$$

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





# Overfitting and Tree Pruning

- Overfitting: 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
  - Prepruning: 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
  - Postpruning: 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”



# Enhancements to Basic Decision Tree Induction

- Allow for continuous-valued attributes
  - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
- Handle missing attribute values
  - Assign the most common value of the attribute
  - Assign probability to each of the possible values
- Attribute construction
  - Create new attributes based on existing ones that are sparsely represented
  - This reduces fragmentation, repetition, and replication



# Classification in Large Databases

- Scalability: Classifying datasets 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



# Classification

## Rule-Based Classification

---

- how many students will attend the lecture, what will the temperature be in August etc. for prediction models
- questions like: what class does it fall into with XX values
- what attributes to split on with decision trees
- ...



# Using IF-THEN Rules for Classification

- Represent the knowledge in the form of IF-THEN rules

Rule: IF age = youth AND student = yes THEN buys\_computer = yes

- Rule antecedent/precondition vs. rule consequent

- Assessment of a rule: coverage and accuracy

- $n_{covers}$  = # of tuples covered by R
- $n_{correct}$  = # of tuples correctly classified by R

$coverage(R) = n_{covers} / |D|$  /\* *D: training dataset* \*/

$accuracy(R) = n_{correct} / n_{covers}$

labels for rules D= training dataset |D|  
cardinality of D - when doing a rule  
extraction we need the labels for the  
records of the dataset

- If more than one rule are triggered, need conflict resolution

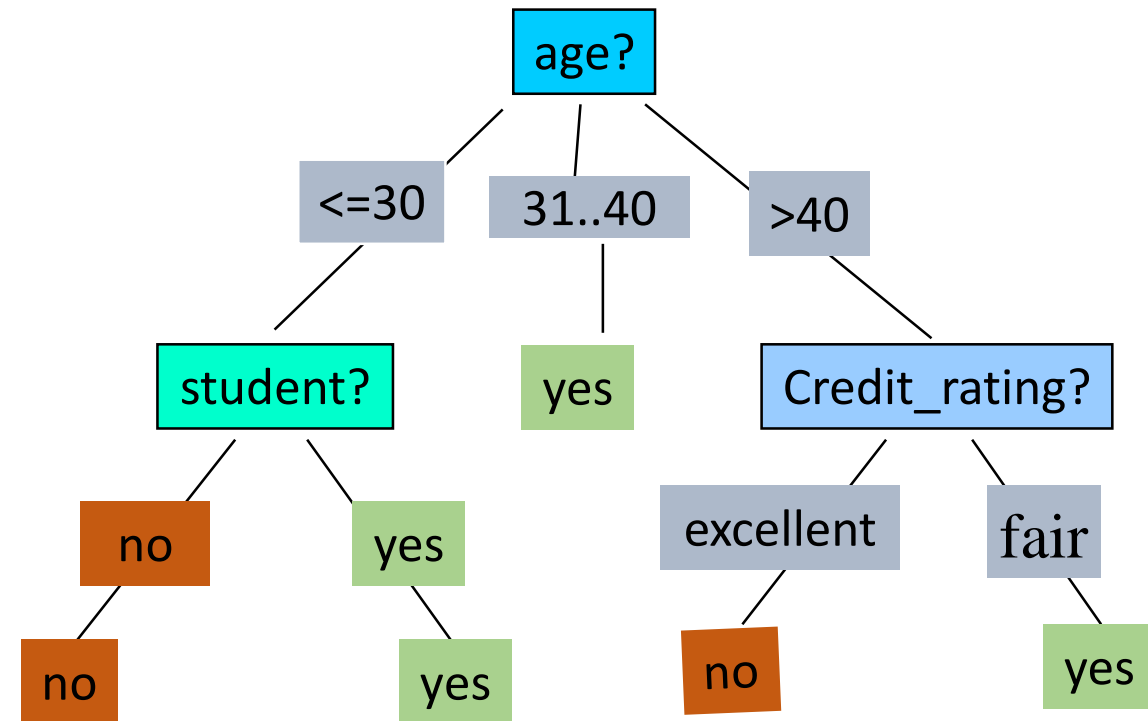
- Size ordering: assign the highest priority to the triggering rules that has the “toughest” requirement (i.e., with the most attribute tests)
- Class-based ordering: decreasing order of prevalence or misclassification cost per class
- Rule-based ordering (decision list): rules are organized into one long priority list, according to some measure of rule quality or by experts



# Rule Extraction from a Decision Tree

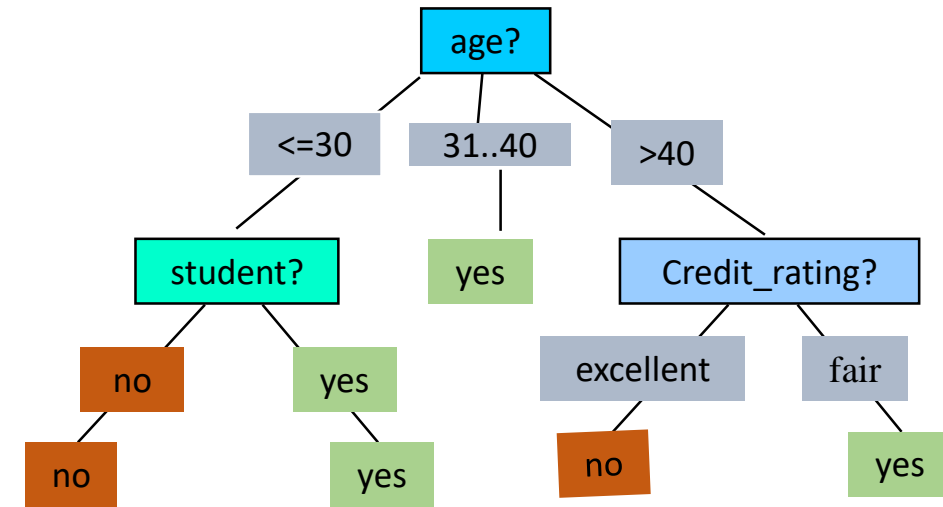
helps with the problem that a rule cannot satisfy two decisions

- Rules are easier to understand than large trees
- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction: the leaf holds the class prediction
- Rules are mutually exclusive and exhaustive



# Rule Extraction from a Decision Tree (Cont.)

- Example:



Rule extraction from our buys\_computer decision-tree

IF age = young AND student = no	THEN buys_computer = no
IF age = young AND student = yes	THEN buys_computer = yes
IF age = mid-age	THEN buys_computer = yes
IF age = old AND credit_rating = excellent	THEN buys_computer = no
IF age = old AND credit_rating = fair	THEN buys_computer = yes



# Rule Induction: Sequential Covering Method

- Sequential covering algorithm: Extracts rules directly from training data
- Typical sequential covering algorithms: FOIL, AQ, CN2, RIPPER
- Rules are learned sequentially, each for a given class  $C_i$  will cover many tuples of  $C_i$  but none (or a few) of the tuples of other classes
- Steps:
  - Rules are learned one at a time
  - Each time a rule is learned, the tuples covered by the rules are removed
  - Repeat the process on the remaining tuples until termination condition, e.g., when no more training examples or when the quality of a rule returned is below a user-specified threshold
- Comparison with decision-tree induction: learning a set of rules simultaneously



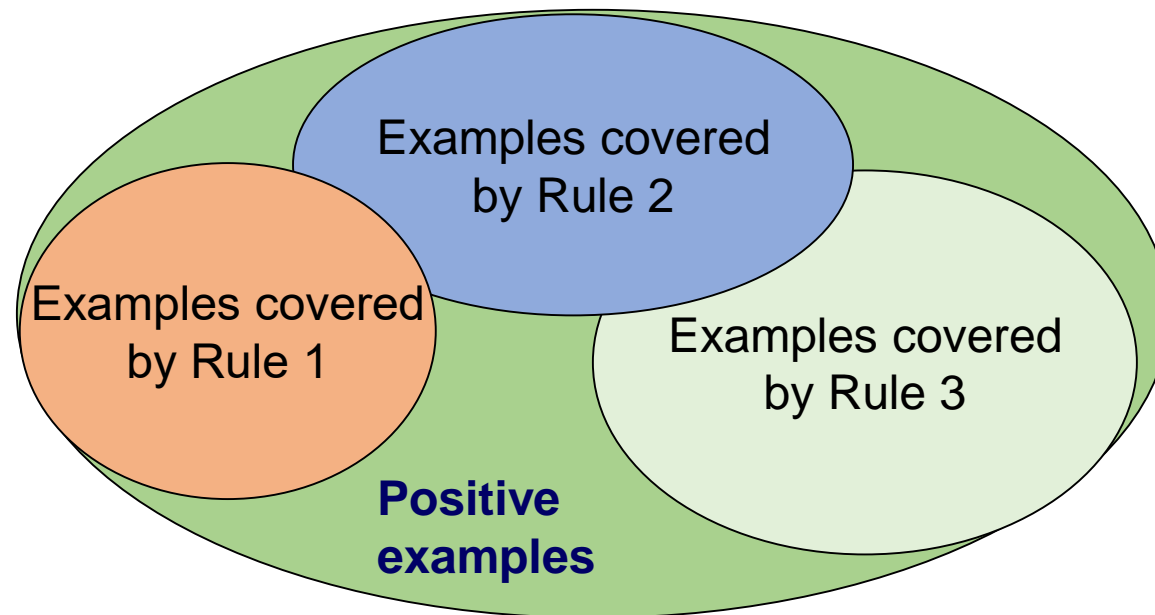


# Sequential Covering Algorithm

WHILE (enough target tuples left)

    generate a rule

    remove positive tuples that satisfy the rule

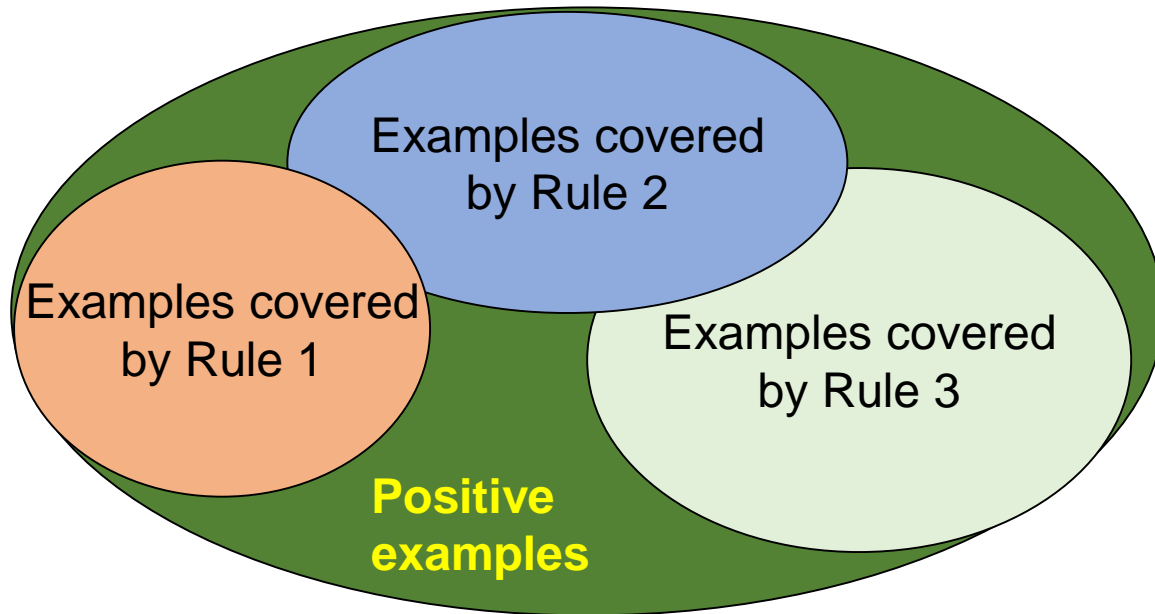


# Sequential Covering Algorithm

WHILE (enough target tuples left)

    generate a rule

    remove positive tuples that satisfy the rule

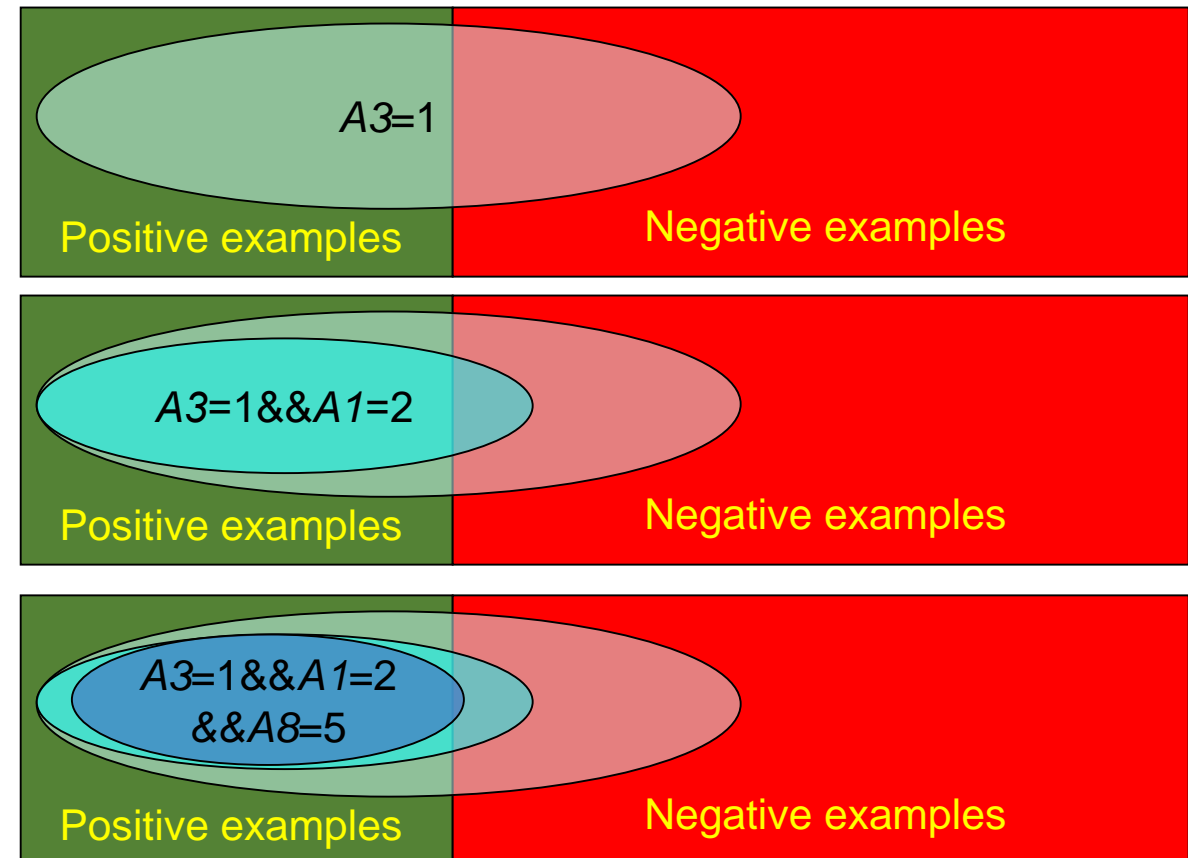


WHILE(TRUE)

    find the best predicate  $p$

    if  $\text{foil-gain}(p) > \text{threshold}$  then add  $p$  to current rule

    else break



# Learning a Rule

positive example= correctly predicted  
if binary - always one class pos/neg  
if more classes - positives for each rule

if a rule covers more  
positives and less negatives  
then we chose the one with  
less negatives

- Start with the most general rule possible: condition = empty
- Adding new attributes by adopting a greedy depth-first strategy
  - Picks the one that most improves the rule quality
- Rule-Quality measures: consider both coverage and accuracy
  - Foil-gain (in FOIL & RIPPER): assesses info\_gain by extending condition

$$FOIL\_GAIN = pos' \times \left( \log_2 \frac{pos'}{pos' + neg'} - \log_2 \frac{pos}{pos + neg} \right)$$

- favors rules that have high accuracy and cover many positive tuples
- Rule pruning based on an independent set of test tuples

$$FOIL\_PRUNE(R) = \frac{pos - neg}{pos + neg}$$

$pos/neg$  and  $pos'/neg'$  are # of positive/negative tuples covered by  $R$  and  $R'$ .

If FOIL\_Prune is higher for the pruned version of  $R$ , prune  $R$



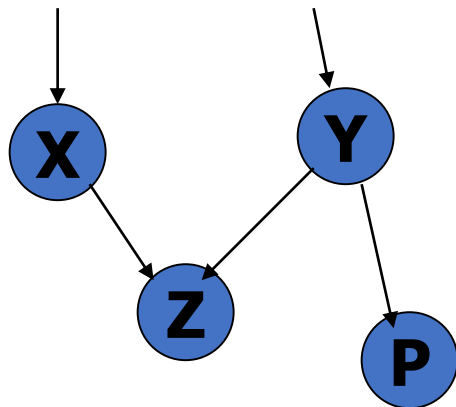
# **Classification**

## **Neural Networks**

---

# Bayesian Belief Networks

- Bayesian belief networks (also known as Bayesian networks, probabilistic networks): allow class conditional independencies between subsets of variables
- A (directed acyclic) graphical model of causal relationships
  - Represents dependency among the variables
  - Gives a specification of joint probability distribution



- Nodes: random variables
- Links: dependency
- $X$  and  $Y$  are the parents of  $Z$ , and  $Y$  is the parent of  $P$
- No dependency between  $Z$  and  $P$
- Has no loops/cycles

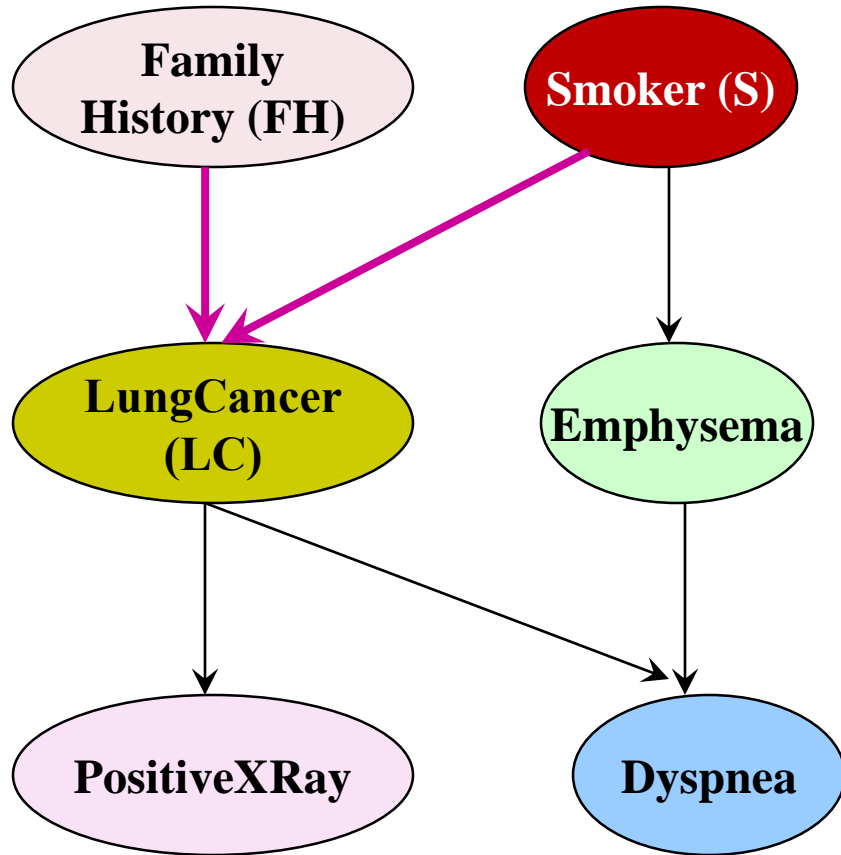
variable with direction  
no cycles are allowed

if X and Y are satisfied - what's the prob. of Z

table shows each possible combination and prob

# Bayesian Belief Networks: Example

family and history - higher probability of 0.9



	(FH, S)	(FH, ~S)	(~FH, S)	(~FH, ~S)
LC	0.8	0.5	0.7	0.1
~LC	0.2	0.5	0.3	0.9

Conditional Probability Table (CPT)

- shows the conditional probability for each possible combination of its parents
- Derivation of the probability of a particular combination of values of  $X$ , from CPT:

$$p(x_1, x_2, \dots, x_n) = \prod_{j=1}^n p(x_j | \text{Parents}(Y_i))$$

representing value

what is the prob. of x for the parents  $Y$  (like family history & smoker)

Bayesian Belief Networks



# Training Bayesian Networks

- **Scenario 1:** Given both the network structure and all variables observable: *compute only the CPT entries*
- **Scenario 2:** Network structure known, some variables hidden: *gradient descent* (greedy hill-climbing) method, i.e., search for a solution along the steepest descent of a criterion function
  - Weights are initialized to random probability values
  - At each iteration, it moves towards what appears to be the best solution at the moment, without backtracking
  - Weights are updated at each iteration & converge to local optimum
- **Scenario 3:** Network structure unknown, all variables observable: search through the model space to *reconstruct network topology*
- **Scenario 4:** Unknown structure, all hidden variables: **No good algorithms** known for this purpose



# Classification by Backpropagation

- Backpropagation: A **neural network** learning algorithm
- Started by psychologists and neurobiologists to develop and test computational analogues of neurons
- A neural network: A set of connected input/output units where each connection has a **weight** associated with it      set of knots to represent input layer and output layers
- During the learning phase, the **network learns by adjusting the weights** so as to be able to predict the correct class label of the input tuples
- Also referred to as **connectionist learning** due to the connections between units  
learning means finding the weight for each layer by training the data until the weights are stable





# Neural Network as a Classifier

- Weakness

- Long training time
- Require a number of parameters typically best determined empirically, e.g., the network topology or “structure.”     you don't know how to adjust the parameters and have to (re)construct
- Poor interpretability: Difficult to interpret the symbolic meaning behind the learned weights and of “hidden units” in the network     you only get a class at the end and don't know what happened (not like a tree where you can follow line/see a rule etc.) - can't see the weights in the hidden units

- Strength

- High tolerance to noisy data     and outliers
- Ability to classify untrained patterns
- Well-suited for continuous-valued inputs and outputs
- Successful on an array of real-world data, e.g., hand-written letters
- Algorithms are inherently parallel
- Techniques have recently been developed for the extraction of rules from trained neural networks  
e.g. predicting cancer accurately - but you cannot tell a patient that he'll now have heart failure because a computer says so. but we can use it for instance to tell a company why they need to replace an expensive part of an airplane and show them exactly why they need to do it



# A Multi-Layer Feed-Forward Neural Network

Output vector

initially random values for the weights that are being adjusted during the process

learning factors

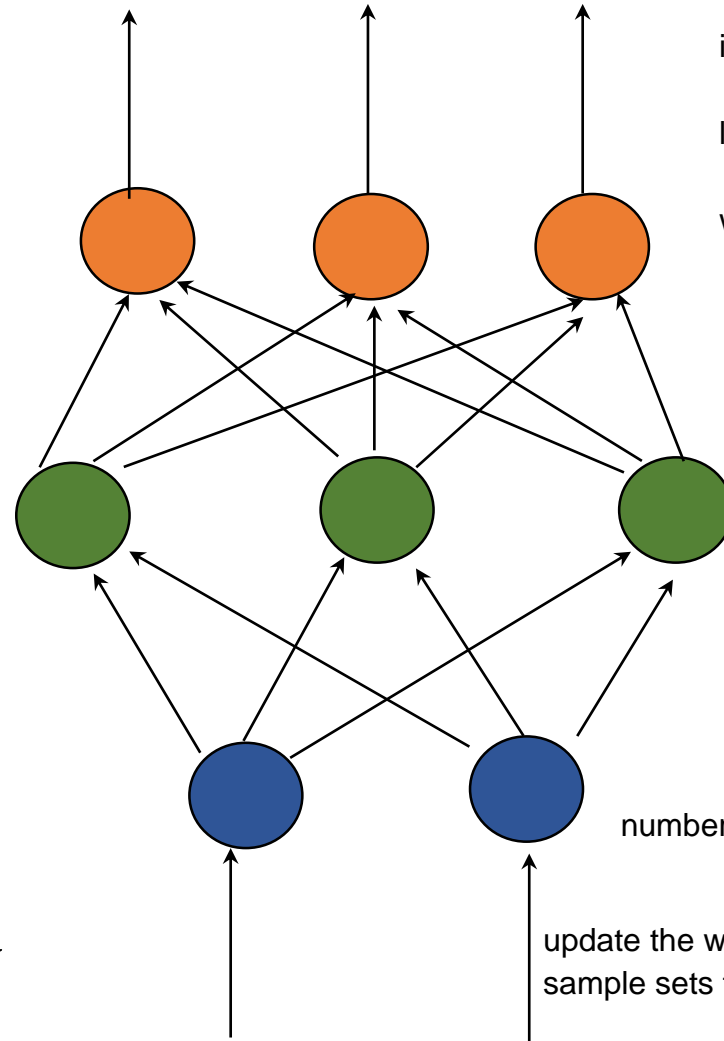
Output layer

weight will keep being adjusted

Hidden layer

Input layer

Input vector:  $X$



$$w_j^{(k+1)} = w_j^{(k)} + \lambda \left( y_i - \hat{y}_i^{(k)} \right) x_{ij}$$

input paramter

number of knots are the same as the number of variables in the dataset/tuple

update the weight after we run the training sample. So amounts of adjudications is the amount of sample sets that run through the data.



# How a Multi-Layer Neural Network Works?

- The **inputs** to the network correspond to the attributes measured for each training tuple
- Inputs are fed simultaneously into the units making up the **input layer**
- They are then weighted and fed simultaneously to a **hidden layer**
- The number of hidden layers is arbitrary, although usually only one
- The weighted outputs of the last hidden layer are input to units making up the **output layer**, which emits the network's prediction
- The network is **feed-forward**: None of the weights cycles back to an input unit or to an output unit of a previous layer
- From a statistical point of view, networks perform **nonlinear regression**: Given enough hidden units and enough training samples, they can closely approximate any function



# Defining a Network Topology

- Decide the **network topology**: Specify # of units in the *input layer*, # of *hidden layers* (if  $> 1$ ), # of units in *each hidden layer*, and # of units in the *output layer*
- Normalize the input values for each attribute measured in the training tuples to [0.0—1.0]
- One **input** unit per domain value, each initialized to 0
- **Output**, if for classification and more than two classes, one output unit per class is used
- Once a network has been trained and its accuracy is **unacceptable**, repeat the training process with a *different network topology* or a *different set of initial weights*



# Backpropagation

in theory works better because it fixes false positives  
but takes a long time

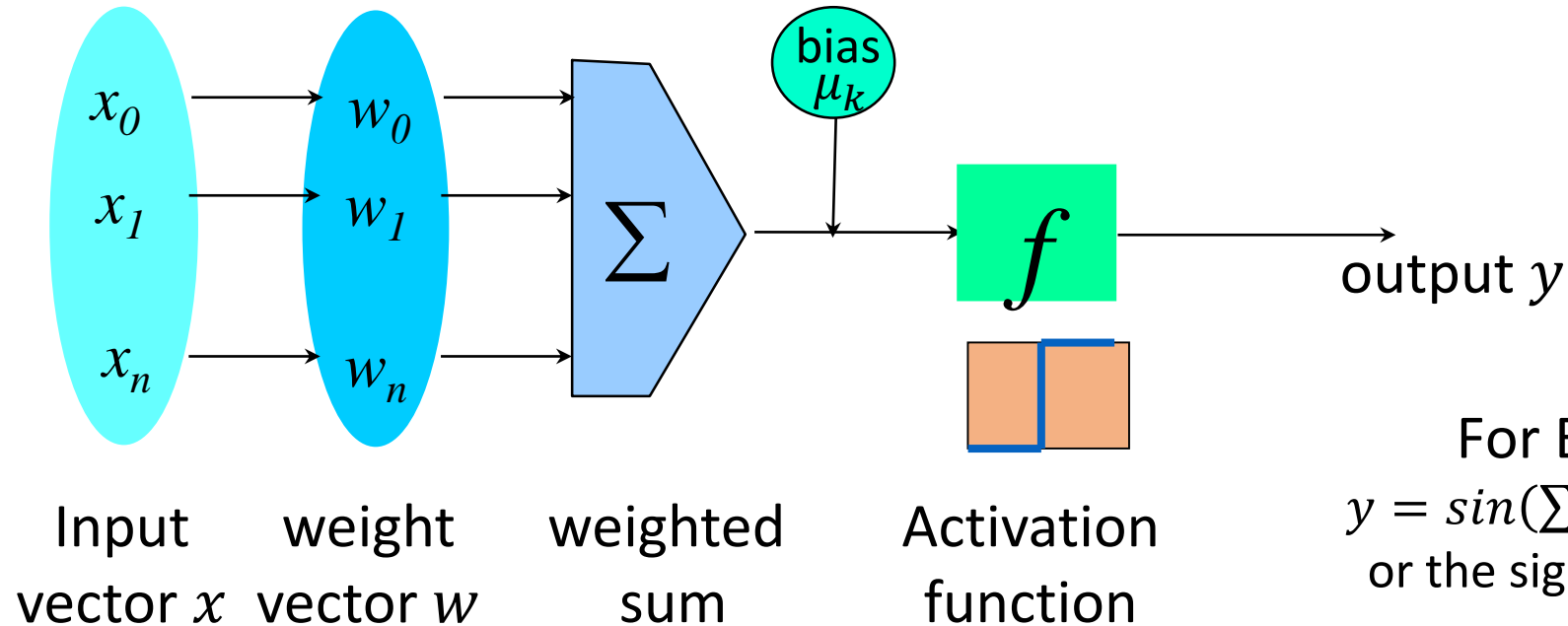
- Iteratively process a set of training tuples & compare the network's prediction with the actual known target value
- For each training tuple, the weights are modified to **minimize the mean squared error** between the network's prediction and the actual target value
- Modifications are made in the “**backwards**” direction: from the output layer, through each hidden layer down to the first hidden layer, hence “**backpropagation**”
- Steps
  - Initialize weights to small random numbers, associated with biases
  - Propagate the inputs forward (by applying activation function)
  - Backpropagate the error (by updating weights and biases)
  - Terminating condition (when error is very small, etc.)



For Example

$$y = \text{sign}\left(\sum_{i=0}^n w_i x_i - \mu_k\right)$$

# Neuron: A Hidden/Output Layer Unit



For Example:

$$y = \sin\left(\sum_{i=0}^n w_i x_i - \mu_k\right)$$

or the sigmoid function

- An  $n$ -dimensional input vector  $\mathbf{x}$  is mapped into variable  $y$  by means of the scalar product and a nonlinear function mapping
- The inputs to unit are outputs from the previous layer. They are multiplied by their corresponding weights to form a weighted sum, which is added to the bias associated with unit. Then a nonlinear activation function is applied to it

# Efficiency and Interpretability

- **Efficiency** of backpropagation: Each epoch (one iteration through the training set) takes  $O(|D| \times w)$ , with  $|D|$  tuples and  $w$  weights, but # of epochs can be exponential to  $n$ , the number of inputs, in worst case
- For easier comprehension: **Rule extraction** by network pruning
  - Simplify the network structure by removing weighted links that have the least effect on the trained network
  - Then perform link, unit, or activation value clustering
  - The set of input and activation values are studied to derive rules describing the relationship between the input and hidden unit layers
- **Sensitivity analysis**: assess the impact that a given input variable has on a network output. The knowledge gained from this analysis can be represented in rules



# **Classification**

## **Support Vector Machines**

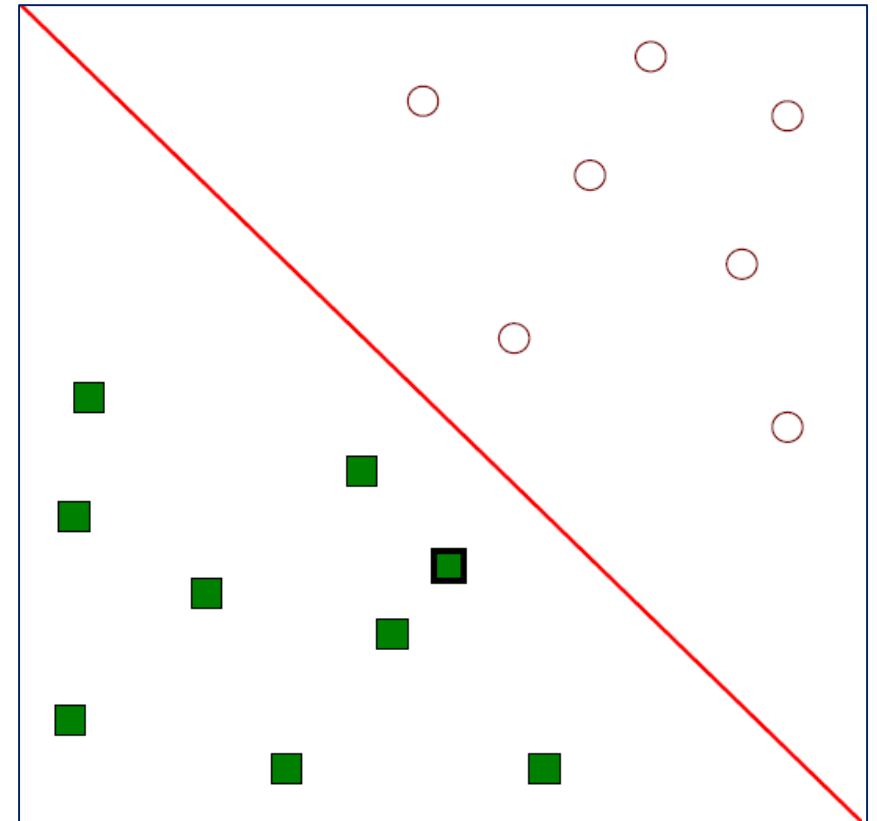
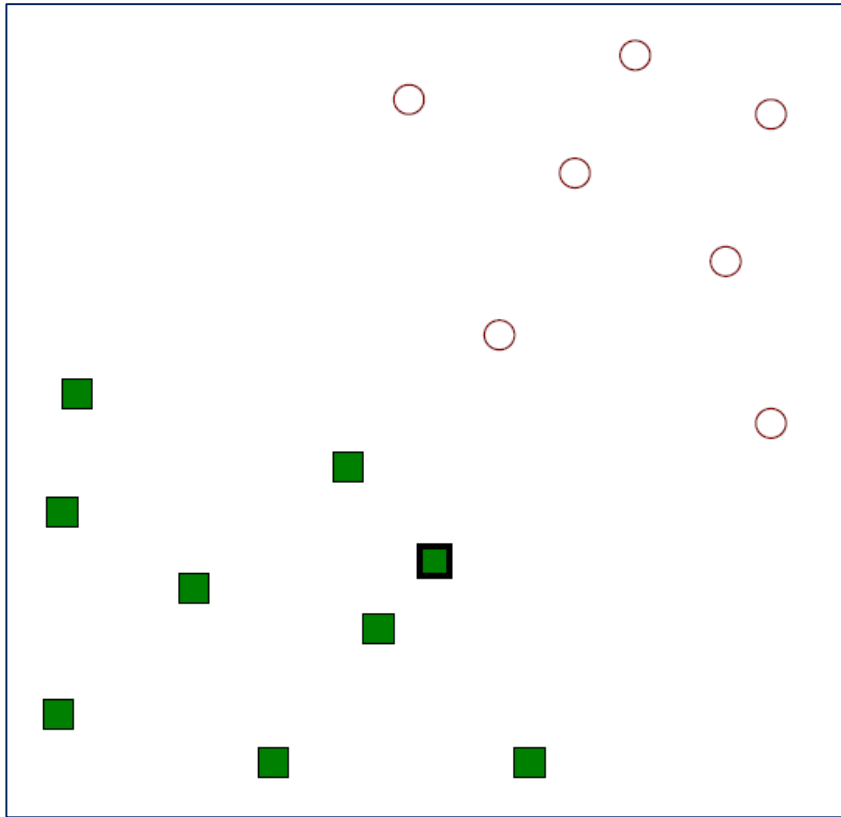
---



# Support Vector Machines (SVM)

Find linear hyperplane (decision boundary) that will separate the data

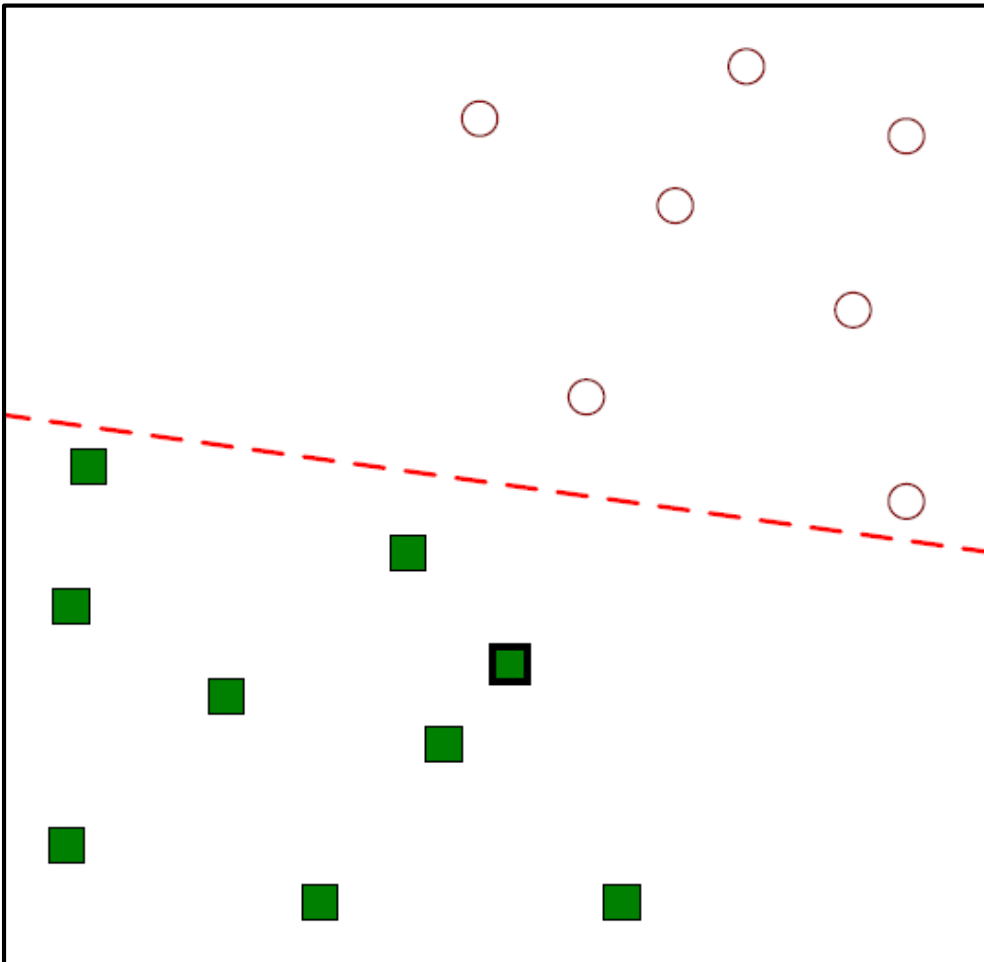
One possible separators



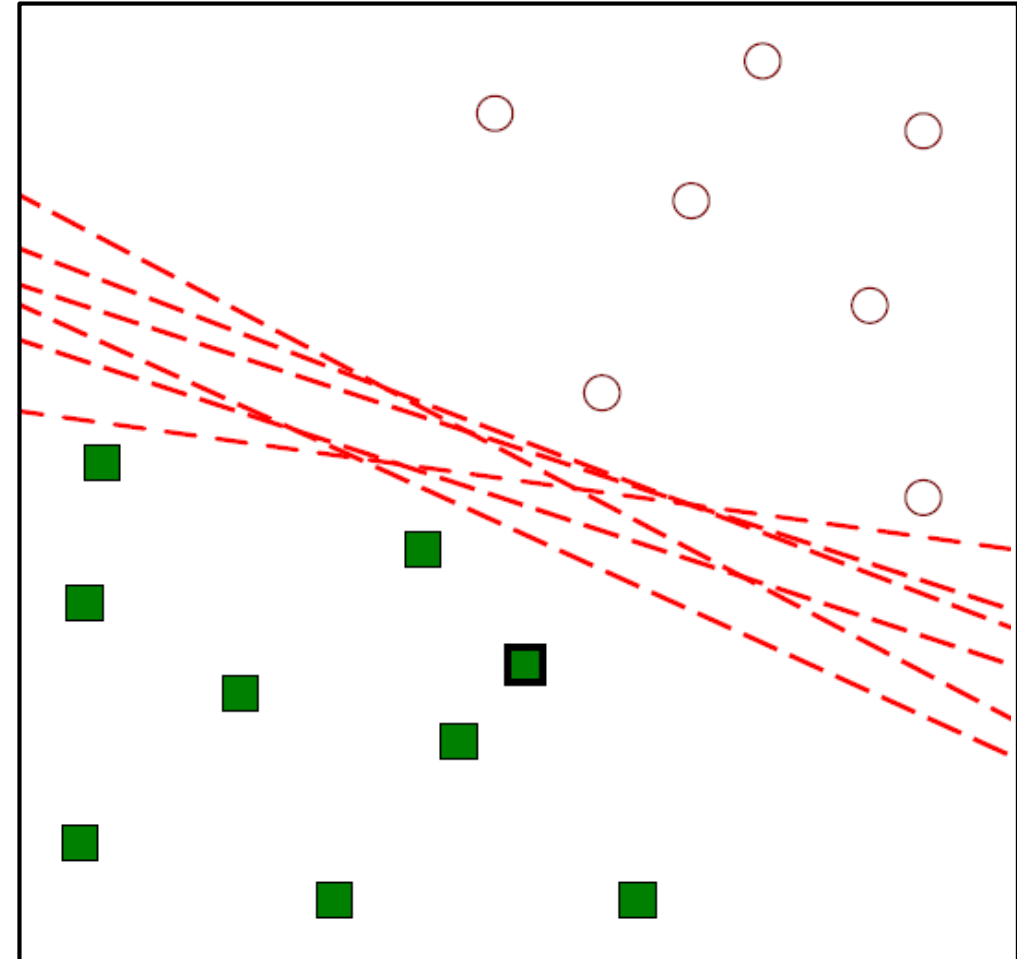
# Support Vector Machines (SVM) (Cont.)

hyperplanes are the margin

Another possible separator



Other possible separators



# Support Vector Machines (SVM) (Cont.)

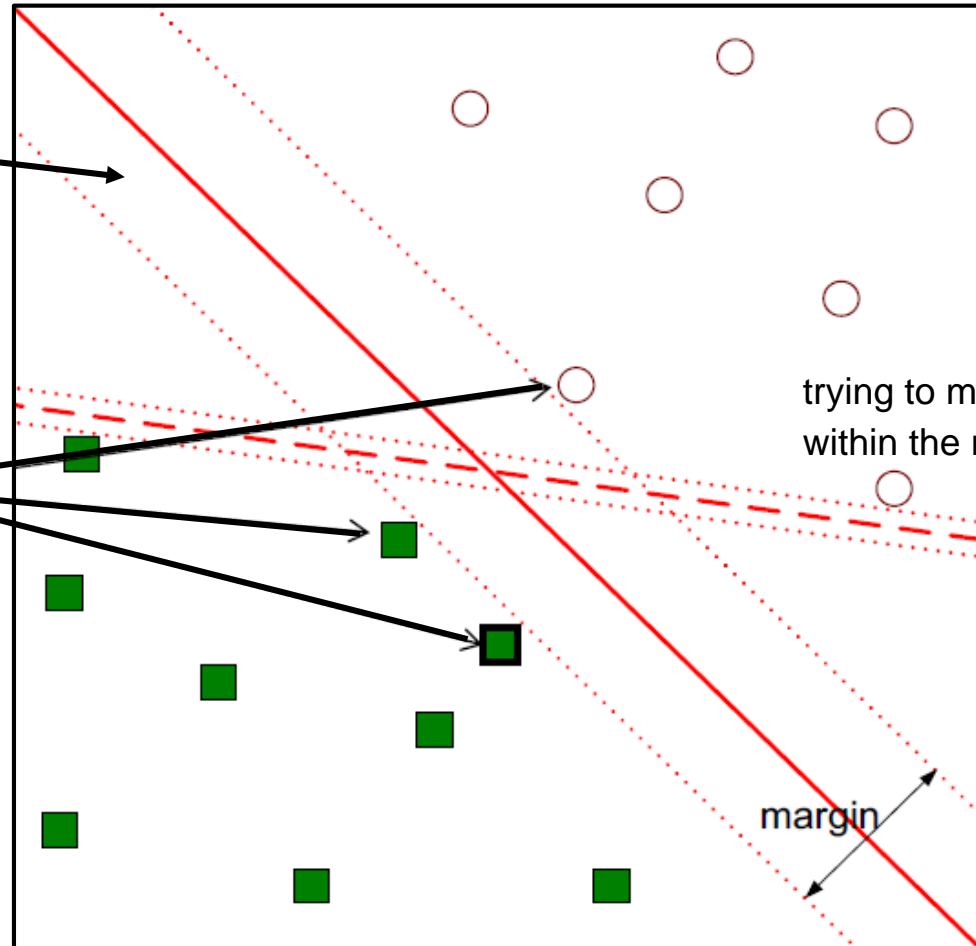
Find the hyperplane that maximizes the margin

objects close to the separator (support vector) tell the boundaries of the margins

Better separator

## Support Vectors

are the points that the margin is pushed up against



trying to maximize the margin but no data point within the margin

the dashed lines are the margins for the hyperplane. we want to add a margin/ have a margin to

seperators so that there's the least margin so that the points are outside

# SVM – Linearly Separable

- A separating hyperplane can be written as

$$W \bullet X + b = 0 \quad \text{not zero for a specific reason just to use the equation}$$

- where  $W = \{w_1, w_2, \dots, w_n\}$  is a weight vector and  $b$  a scalar (bias)
- For 2-D it can be written as

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

- The hyperplane defining the sides of the margin:

$$H_1: w_0 + w_1 x_1 + w_2 x_2 \geq 1 \quad \text{for } y_i = +1, \text{ and } \text{assume that class label is positive}$$

$$H_2: w_0 + w_1 x_1 + w_2 x_2 \leq -1 \quad \text{for } y_i = -1 \text{ assume that class label is negative}$$

- Any training tuples that fall on hyperplanes  $H_1$  or  $H_2$  (i.e., the sides defining the margin) are support vectors
- This becomes a constrained (convex) quadratic optimization problem



# SVM are Effective on High Dimensional Data

- The complexity of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data
- The support vectors are the essential or critical training examples — they lie closest to the decision boundary
- If all other training examples are removed and the training is repeated, the same separating hyperplane would be found = deterministic algorithm - if we run it again = same result
- The number of support vectors found can be used to compute an (upper) bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality we can estimate the vectors with the support vector machine
- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high



# SVM vs Neural Networks

usually used for two classes

- SVM

produces same output every time we run it

- Deterministic algorithm
- Nice generalization properties
- Hard to learn – learned in batch mode using quadratic programming techniques
- Using kernels can learn very complex functions

simple - one can run it immediately and get results right away  
doesn't have to be binary

- Neural Network

- Nondeterministic algorithm
- Generalizes well but doesn't have strong mathematical foundation
- Can easily be learned in incremental fashion
- To learn complex functions—use multilayer perceptron (nontrivial)

can have different results- because input is initialized randomly: different architecture and different values for the input because starting weights aren't the same in beginning



# Evaluation and Model Selection

---

# Evaluation: the Key to Success

- How predictive is the model we have learned?
- Error on the training data is not a good indicator of performance on future data
  - Otherwise 1-NN would be the optimum classifier!
- Simple solution that can be used if a large amount of (labeled) data is available:
  - Split data into training and test set usually two thirds for training and one third for testing
- However: (labeled) data is usually limited
  - More sophisticated techniques need to be used





# Issues in Evaluation

- Statistical reliability of estimated differences in performance
- Choice of performance measure:
  - Number of correct classifications
  - Accuracy of probability estimates
  - Error in numeric predictions
- Costs assigned to different types of errors
  - Many practical applications involve costs

how costly is a false positive to a true negative? what do we include into the model for the cost of some values being wrong (e.g. with creditcard fraud we'd rather include false alarms than exclude actual fraud)

initially we got to figure out a performance measure



# Training and Testing

- Natural performance measure for classification problems: error rate
  - Success: instance's class is predicted correctly
  - Error: instance's class is predicted incorrectly
  - Error rate: proportion of errors made over the whole set of instances
- **Resubstitution error**: error rate obtained by evaluating model on training data

allowing certain error



# Training and Testing (Cont.)

- Test set: independent instances that have played no part in formation of classifier
  - Assumption: both training data and test data are representative samples of the underlying problem
- Test and training data may differ in nature
  - Example: classifiers built using customer data from two different towns A and B
  - To estimate performance of classifier from town A in completely new town, test it on data from B



# Parameter Tuning

- It is important that the test data is not used in any way to create the classifier
- Some learning schemes operate in two stages:
  - Stage 1: build the **basic structure** very simple parameters
  - Stage 2: **optimize parameter** settings
- The **test data cannot be used for parameter tuning!** test data only to validate the model itself!
- Proper procedure uses three sets: training data, validation data, and test data
  - Validation data is used to optimize parameters

not easy to say what's more or less accurate. what does an application need? make a decision dependent on that!

# Making the Most of the Data

- Once evaluation is complete, all the data can be used to build the final classifier
- Generally, the larger the training data the better the classifier (but returns diminish)
- The larger the test data the more accurate the error estimate
- Holdout procedure: method of splitting original data into training and test set
  - Dilemma: ideally both training set and test set should be large!



# Predicting Performance

- Assume the estimated error rate is 25%. How close is this to the true error rate?
- Depends on the amount of test data
- Prediction is just like tossing a (biased!) coin
- “Head” is a “success”, “tail” is an “error”
- In statistics, a succession of independent events like this is called a Bernoulli process
- Statistical theory provides us with confidence intervals for the true underlying proportion

how confident am I that the model will be a certain percentage accurate



# Confidence Intervals

- We can say:  $p$  lies within a certain specified interval with a certain specified confidence
- Example:  $S = 750$  successes in  $N = 1000$  trials
- Estimated success rate: 75%
- How close is  $p$  to true success rate?
- Answer: with 80% confidence  $p$  is located in  $[73.2, 76.7]$
- Another example:  $S = 75$  and  $N = 100$
- Estimated success rate: 75%
- With 80% confidence  $p$  in  $[69.1, 80.1]$



# Mean and Variance

- Mean and variance for a Bernoulli trial:  $p, p(1-p)$
- Expected success rate  $X = S/N$
- Mean and variance for  $X$ :  $p, p(1-p)/N$
- For large enough  $N$ ,  $X$  follows a Normal distribution
- $c\%$  confidence interval  $[-z \leq X \leq z]$  for a random variable  $X$  is determined using:

$$p(-z \leq X \leq z) = c$$

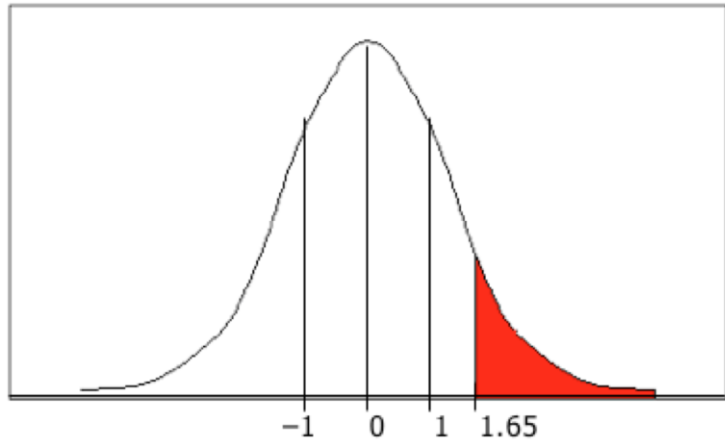
- For a symmetric distribution such as the normal distribution we have:  
$$p(-z \leq X \leq z) = 1 - 2 \times p(X \geq z)$$





# Confidence Limits

- Confidence limits for the normal distribution with 0 mean and a variance of 1:



$p(X > z)$	$z$
0.1%	3.09
0.5%	2.58
1%	2.33
5%	1.65
10%	1.28
20%	0.84
40%	0.25

- Thus:  $p(-1.65 \leq X \leq 1.65) = 90\%$
- We assume the number of trials  $N$  is 1000
- To use this we have to transform our random variable  $X$  to have 0 mean and unit variance

# Transforming $X$

- Transformed value for  $X$ :  $\frac{X-p}{\sqrt{p(1-p)/N}}$   
(i.e., subtract the mean and divide by the standard deviation)

- Resulting equation:

$$p \left( -z \leq \frac{X - p}{\sqrt{p(1 - p)/N}} \leq z \right) = c$$

- Solving for  $p$  yields an expression for the confidence limits:

$$p = \left( X + \frac{z^2}{2N} \pm z \sqrt{\frac{X}{N} - \frac{X^2}{N} + \frac{z^2}{4N^2}} \right) / \left( 1 + \frac{z^2}{N} \right)$$



# Examples

- $z = 1.28$   
 $X = 75\%$ ,  $N = 1000$  the  $p \in [0.732, 0.767]$   
 $X = 75\%$ ,  $N = 100$  the  $p \in [0.691, 0.801]$   
 $X = 75\%$ ,  $N = 10$  the  $p \in [0.549, 0.881]$
- Note that: the Normal distribution assumption is only valid for large  $N$  (i.e.,  $N > 100$ )



# Holdout Estimation

- What should we do if we only have a single dataset?
- The holdout method reserves a certain amount for testing and uses the remainder for training, after shuffling
  - Usually: one third for testing, the rest for training
- Problem: the samples might not be representative
  - Example: class might be missing in the test data
- Advanced version uses stratification
  - Ensures that each class is represented with approximately equal proportions in both subsets

correct ratios in the training and testing data. so that if there's 10% of a certain class in the data it should be 10% of that class in the testing data etc.



# Repeated Holdout Method

- Holdout estimate can be made more reliable by repeating the process with different subsamples
  - In each iteration, a certain proportion is randomly selected for training (possibly with stratification)
  - The error rates on the different iterations are averaged to yield an overall error rate
- This is called the repeated holdout method
- Still not optimum: the different test sets overlap
  - Can we prevent overlapping?



# Cross Validation

dividing the dataset into set of subsets  $k$   
 $k-1$  of the sets for training and one (the  $-1$ ) for testing

- K-fold cross-validation avoids overlapping test sets
  - First step: split data into  $k$  subsets of equal size
  - Second step: use each subset in turn for testing, the remainder for training
  - This means the learning algorithm is applied to  $k$  different training sets
- Often the subsets are stratified before the cross-validation is performed to yield stratified k-fold cross-validation
- The error estimates are averaged to yield an overall error estimate; also, standard deviation is often computed
- Alternatively, predictions and actual target values from the  $k$  folds are pooled to compute one estimate
  - Does not yield an estimate of standard deviation



# Cross Validation (Cont.)

- Standard method for evaluation: stratified ten-fold cross-validation
- Why ten?
  - Extensive experiments have shown that this is the best choice to get an accurate estimate
  - There is also some theoretical evidence for this
- Stratification reduces the estimate's variance
- Even better: repeated stratified cross-validation
  - E.g., ten-fold cross-validation is repeated ten times and results are averaged (reduces the variance)



# Leave-One-Out Cross Validation

- Leave-one-out: a particular form of  $k$ -fold cross-validation:
  - Set number of folds to number of training instances
  - I.e., for  $n$  training instances, build classifier  $n$  times
- Makes best use of the data
- Involves no random subsampling
- Very computationally expensive (exception: using lazy classifiers such as the nearest-neighbor classifier)





# Leave-One-Out CV and Stratification

- Disadvantage of Leave-one-out CV: stratification is not possible
  - It guarantees a non-stratified sample because there is only one instance in the test set!
- Extreme example: random dataset split equally into two classes
  - Best inducer predicts majority class
  - 50% accuracy on fresh data
  - Leave-one-out CV estimate gives 100% error!



# The Bootstrap

- CV uses sampling without replacement
  - The same instance, once selected, can not be selected again for a particular training/test set
- The bootstrap uses sampling with replacement to form the training set
  - Sample a dataset of  $n$  instances  $n$  times with replacement to form a new dataset of  $n$  instances
  - Use this data as the training set
  - Use the instances from the original dataset that do not occur in the new training set for testing



# The 0.632 Bootstrap

- A particular instance has a probability of  $1 - 1/n$  of not being picked
- Thus its probability of ending up in the test data is:

$$\left(1 - \frac{1}{n}\right)^n \approx e^{-1} = 0.368$$

- This means the training data will contain approximately 63.2% of the instances



# Estimating Classification Error with The 0.632 Bootstrap

- The error estimate on the test data will be quite pessimistic
- Trained on just  $\sim 63\%$  of the instances
- Idea: combine it with the resubstitution error:

$$E = 0.632E_{test\ instances} + 0.368E_{training\ instances}$$

- The resubstitution error gets less weight than the error on the test data
- Repeat process several times with different samples; average the results



# More on Bootstrapping

- Probably the best way of estimating performance for very small datasets
- However, it has some problems
- Consider the random dataset from above
- A perfect memorizer will achieve  
0% resubstitution error and  $\sim 50\%$  error on test data
- Bootstrap estimate for this classifier:  
$$(0.632 \times 50\% + 0.368 \times 0\%) = 31.6\%$$
- True expected error: 50%



# Comparing Machine Learning Schemes

- Frequent question: which of two learning schemes performs better?
- Note: this is domain dependent!
- Obvious way: compare 10-fold cross-validation estimates
- Generally sufficient in applications (we do not loose if the chosen method is not truly better)
- However, what about machine learning research?
- Need to show convincingly that a particular method works better in a particular domain from which data is taken



# Comparing Machine Learning Schemes (Cont.)

- Want to show that scheme  $A$  is better than scheme  $B$  in a particular domain
  - For a given amount of training data (i.e., data size)
  - On average, across all possible training sets from that domain
- Let's assume we have an infinite amount of data from the domain
- Then, we can simply
  - sample infinitely many dataset of a specified size
  - obtain a cross-validation estimate on each dataset for each scheme
  - check if the mean accuracy for scheme  $A$  is better than the mean accuracy for scheme  $B$



# Counting the Cost

- In practice, different types of classification errors often incur different costs
- Examples:
- Terrorist profiling: “Not a terrorist” correct 99.99...% of the time
  - Loan decisions
  - Oil-slick detection
  - Fault diagnosis
  - Promotional mailing





# Confusion Matrix

min 48 of lab 2

IMPORTANT

- Focus on the predictive capability of a model

Confusion Matrix:

	Predicted Label		
		Class = Y	Class = N
Actual Label	Class = Y	TP	FN
	Class = N	FP	TN

TP = True Positive  
FP = False Positive  
FN = False Negative  
TN = True Negative

- Most widely used metric:

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN}$$



# Accuracy

- Limitation of accuracy metric:
  - The problem of **unbalanced classes**
  - Consider a 2-class problem
    - Number of class 1 examples = 9990
    - Number of class 0 examples = 10
  - If the model predict everything as class 1
    - $Accuracy = \frac{9990}{10000} = 99.9\%$
- Accuracy is misleading because the classifier didn't predict any class 0 examples
- $Weighted\ Accuracy = \frac{w_1 TP + w_4 TN}{w_1 TP + w_2 FN + w_3 FP + w_4 TN}$



# Classification with Cost

- Two cost Matrices:

	Predicted Label					Predicted Label			
		Class = Y	Class = N				a	b	c
(A) Actual Label	Class = Y	0	1		(B) Actual Label	a	0	1	1
	Class = N	1	0			b	1	0	1
						c	1	1	0

- In cost-sensitive evaluation of classification methods, success rate is replaced by average cost per prediction
  - Cost is given by appropriate entry in the cost matrix

# Cost–Sensitive Classification

- Can take costs into account when making predictions
  - Basic idea: only predict high-cost class when very confident about prediction
- Given: predicted class probabilities
  - Normally, we just predict the most likely class
  - Here, we should make the prediction that minimizes the expected cost
    - Expected cost: dot product of vector of class probabilities and appropriate column in cost matrix
    - Choose column (class) that minimizes expected cost
- This is the minimum-expected cost approach to cost-sensitive classification



# Cost–Sensitive Learning

- So far we haven't taken costs into account at training time
- Most learning schemes do not perform cost-sensitive learning
  - They generate the same classifier no matter what costs are assigned to the different classes
  - Example: standard decision tree learner
- Simple methods for cost-sensitive learning:
  - Resampling of instances according to costs
  - Weighting of instances according to costs
- Some schemes can take costs into account by varying a parameter, e.g., naïve Bayes



# Lift Charts

- In practice, costs are rarely known
- Decisions are usually made by comparing possible scenarios
- Example: promotional mailout to 1,000,000 households
  - Mail to all; 0.1% respond (1000)
  - Data mining tool identifies subset of 100,000 most promising, 0.4% of these respond (400)  
40% of responses for 10% of cost may pay off
  - Identify subset of 400,000 most promising, 0.2% respond (800)
- A lift chart allows a visual comparison



# Generating a Lift Chart

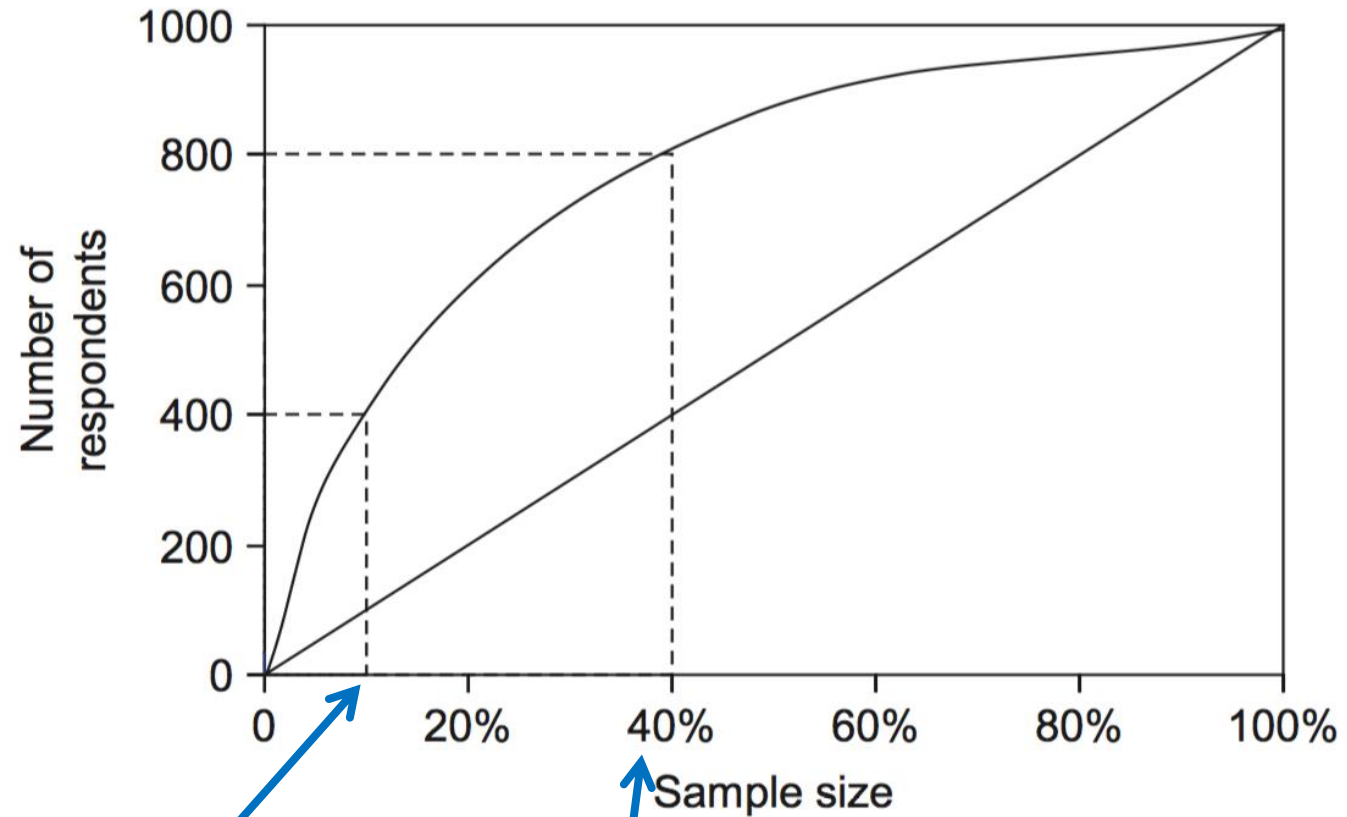
- Sort the instances based on predicted probability of being positive

<i>rank</i>	<i>predicted</i>	<i>actual</i>	<i>rank</i>	<i>predicted</i>	<i>actual</i>
1	0.95	Yes	8	0.80	No
2	0.93	Yes	9	0.80	Yes
3	0.93	No	10	0.79	Yes
4	0.88	Yes	11	0.77	No
5	0.86	Yes	12	0.77	Yes
6	0.85	Yes	13	0.76	Yes
7	0.82	Yes	...	...	...

- $x$ -axis in lift charts is sample size for each probability threshold
- $y$ -axis is the number true positives above the threshold



# A Hypothetical Lift Chart



40% of responses for 10% of cost

80% of responses for 40% of cost





# ROC Curves

- ROC curves are similar to lift charts
  - Stands for “receiver operating characteristic”
  - Used in signal detection to show tradeoff between hit rate and false alarm rate over noisy channel
- Differences to lift chart:
  - $x$ -axis shows percentage of true positives in sample rather than absolute number
  - $x$ -axis shows percentage of false positives in sample rather than sample size



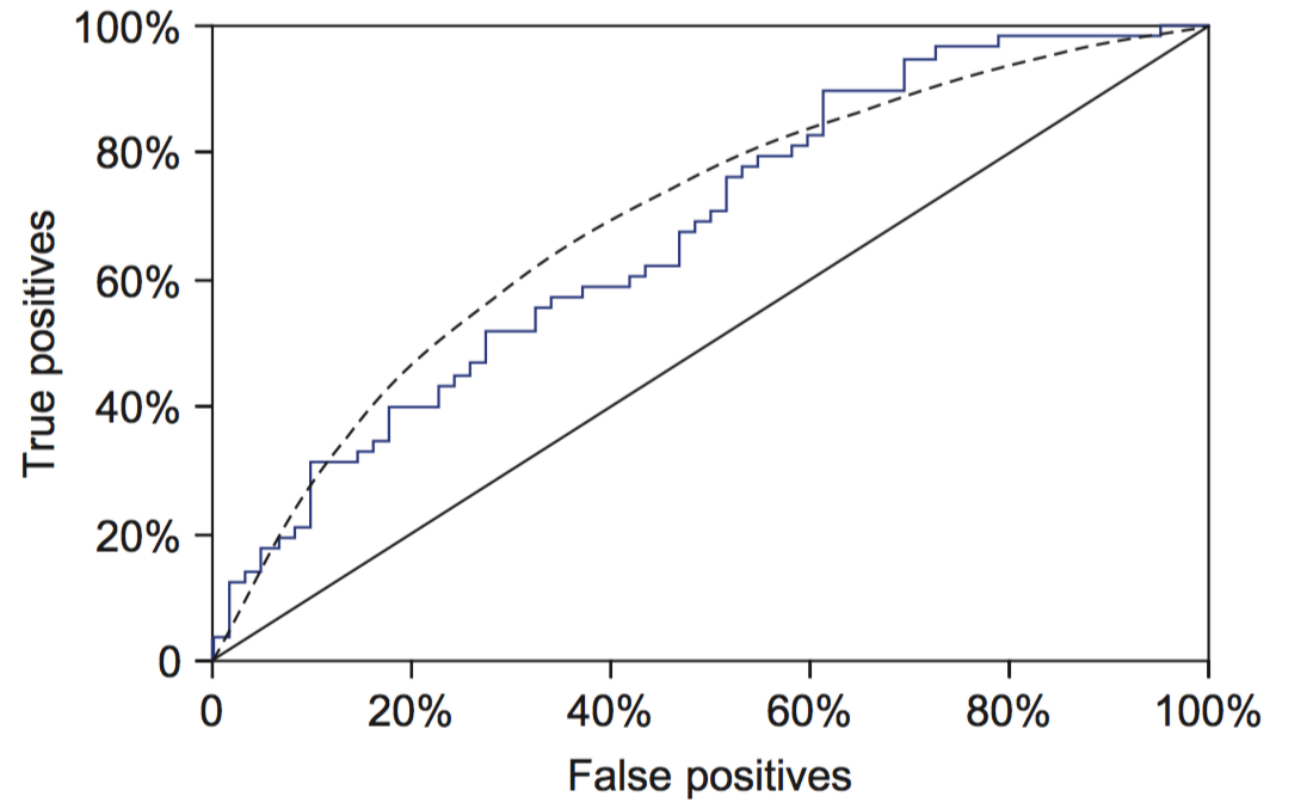
## ROC Curves (Cont.)

- ROC curves: provides visual comparison of classification models
- The area under the ROC 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



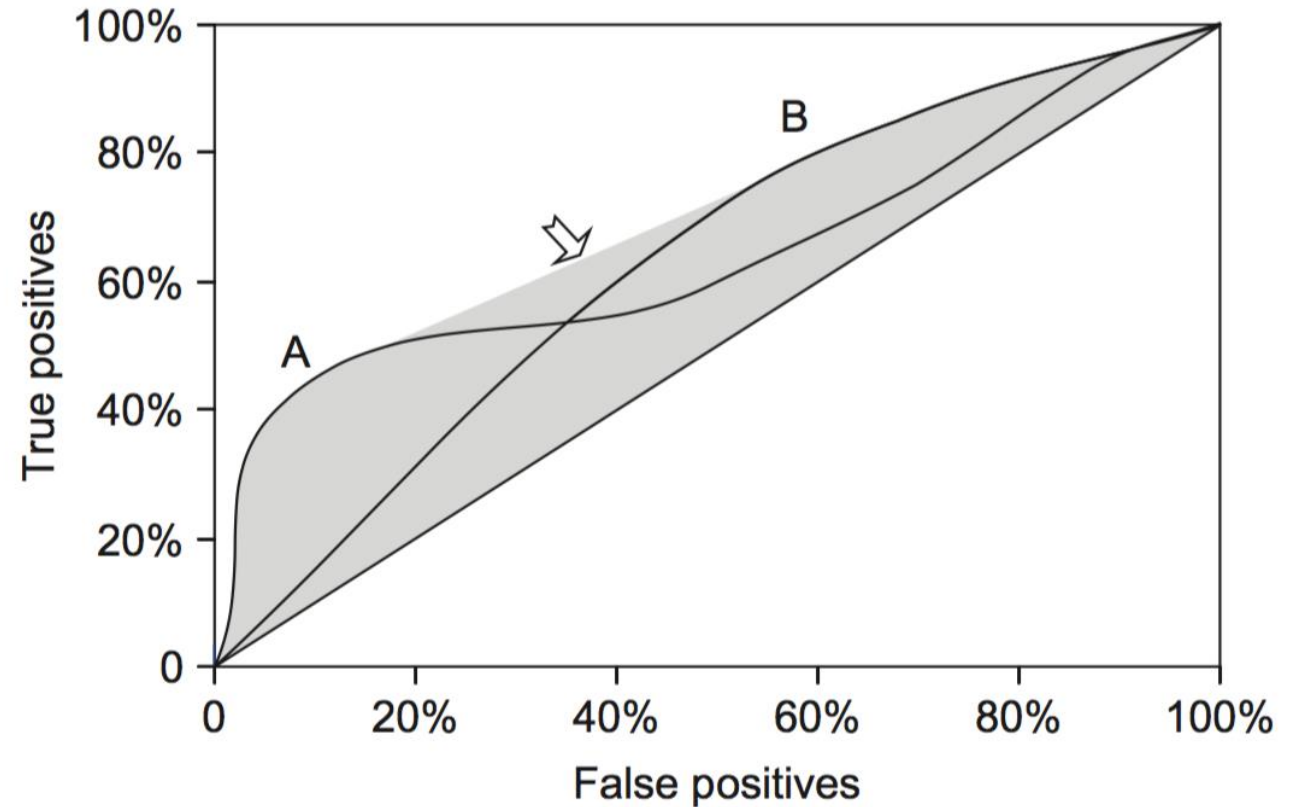
## A Sample ROC Curve

- y-axis represents the true positive rate
- x-axis rep. the false positive rate
- The plot also shows a diagonal line
- A model with perfect accuracy will have an area of 1.0



## ROC Curves for Two Schemes

- For a small, focused sample, use method A
- For a larger one, use method B
- In between, choose between A and B with appropriate probabilities



# More Evaluation Metrics

IMPORTANT

- Percentage of retrieved documents that are relevant:

$$\textit{precision}(P) = TP / (TP + FP)$$

- Percentage of relevant documents that are returned:

$$\textit{recall}(R) = TP / (TP + FN)$$

- Precision/recall curves have hyperbolic shape

- F-measure:  $F_1 = (2 \times R \times P) / (R + P)$

- In general:  $F_\beta = (1 + \beta^2) (R \times P) / (\beta^2 \times P + R)$

- $\textit{sensitivity} \times \textit{specificity} = (TP / (TP + FN)) \times (TN / (FP + TN))$

- Area under the ROC curve (AUC): probability that randomly chosen positive instance is ranked above randomly chosen negative one



# Summary of Evaluation Metrics

	Domain	Plot	Explanation
Lift Chart	Marketing	TP	TP
		Subset Size	$(TP+FP)/(TP+FP+TN+FN)$
ROC Curve	Communication	TP rate	$TP/(TP+FN)$
		FP rate	$FP/(FP+TN)$
Recall-Precision Curve	Information Retrieval	Recall	$TP/(TP+FN)$
		Precision	$TP/(TP+FP)$



# Issues Affecting Model Selection

- Accuracy: classifier accuracy = correctly predicting class label
- Speed
  - Time to construct the model (training time)
  - Time to use the model (classification/prediction time)
- Robustness: handling noise and missing values
- Scalability: efficiency in disk-resident databases
- Interpretability: understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules



# Reading Material & Exercises

- Chapters 8, 9 of the Data Mining: Concepts and Techniques Book
- Chapter 5 of the Data Mining: Practical Machine Learning Tools and Techniques

