



IIT KHARAGPUR



NPTEL ONLINE  
CERTIFICATION COURSES

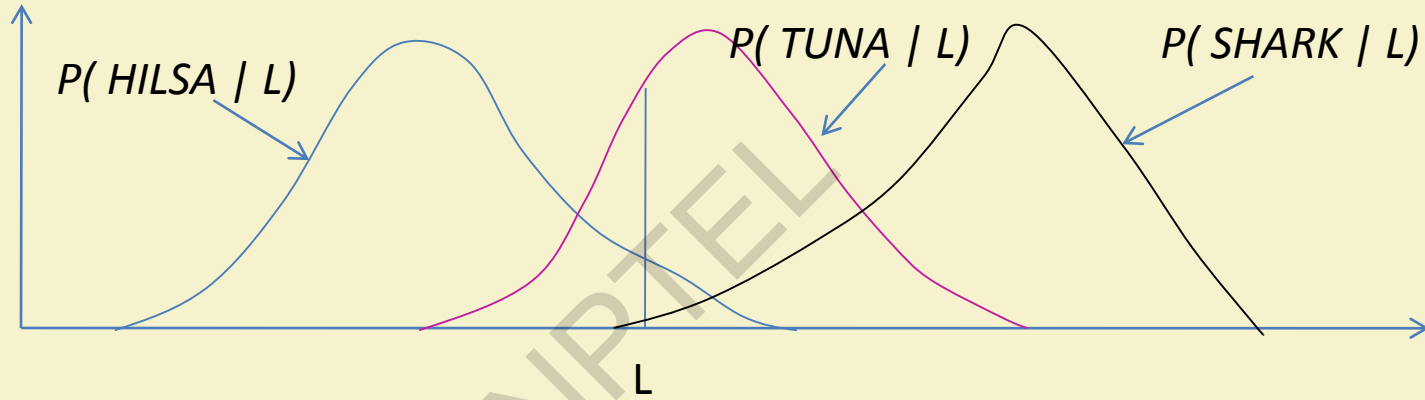
# Data Mining

## Week 4: K-Nearest Neighbor, Classifier Evaluation

Pabitra Mitra

Computer Science and Engineering, IIT Kharagpur

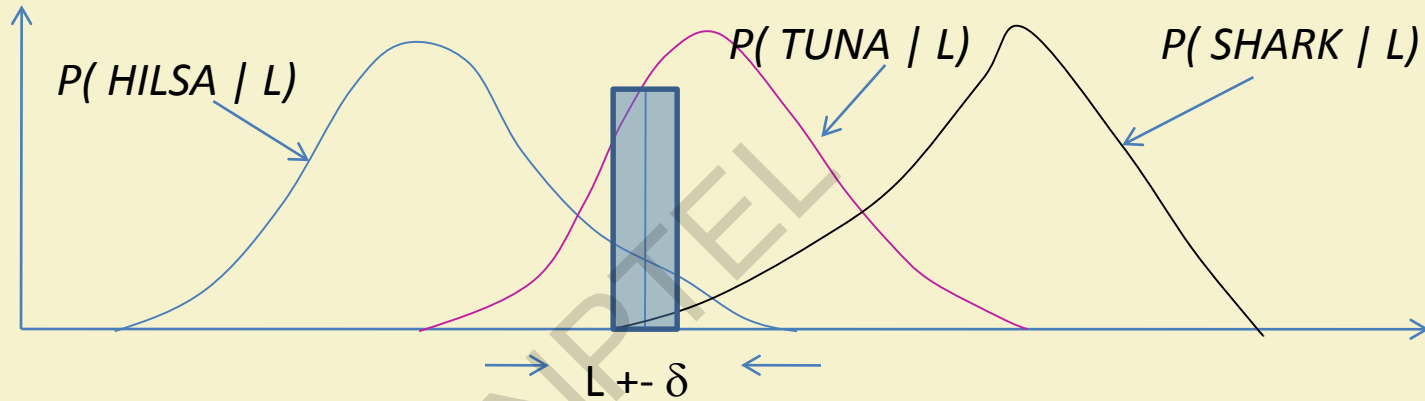
# Bayes Classifier: Recap



Maximum A posteriori (MAP) Rule

Distributions assumed to be of particular family (e.g., Gaussian), and parameters estimated from training data.

# Bayes Classifier: Recap

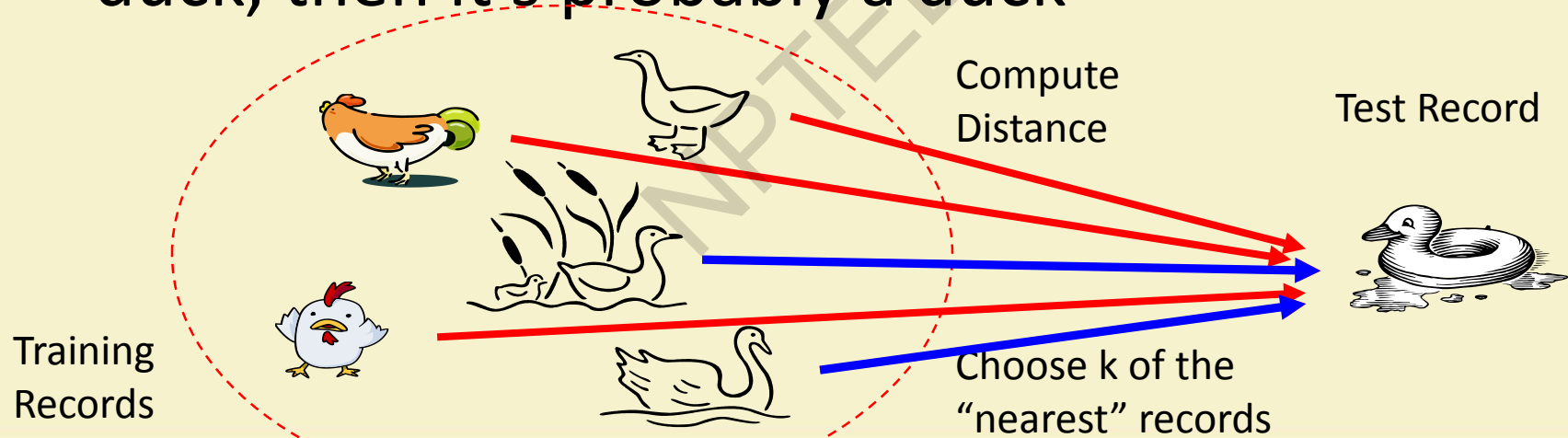


Approximate Maximum A Posteriori (MAP) Rule

*Non-parametric (data driven) approach:* consider a small window around  $L$ , Find which class is most populous in that window.

# Nearest Neighbor Classifiers

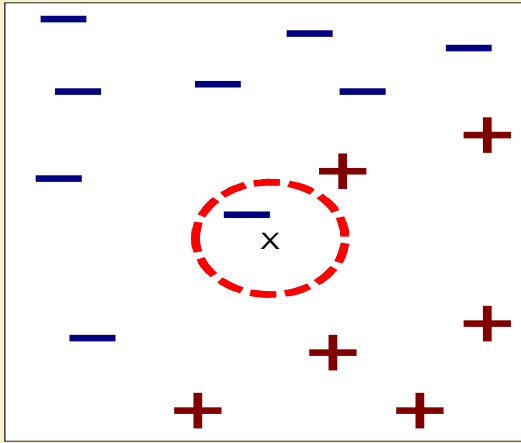
- Basic idea: If it walks like a duck, quacks like a duck, then it's probably a duck



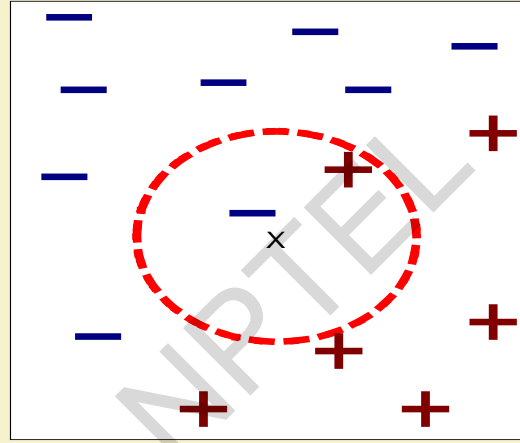
# Basic Idea

- $k$ -NN classification rule is to assign to a test sample the majority category label of its  $k$  nearest training samples
- In practice,  $k$  is usually chosen to be odd, so as to avoid ties
- The  $k = 1$  rule is generally called the nearest-neighbor classification rule

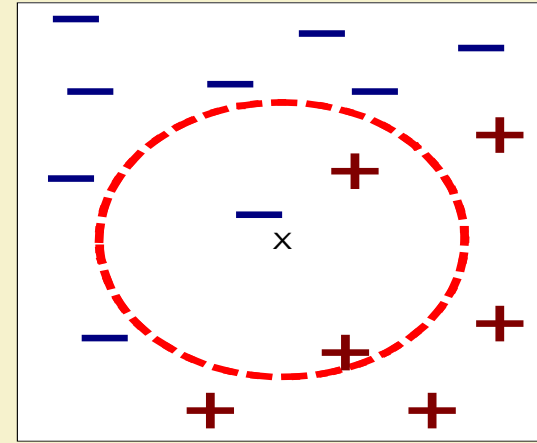
# Definition of Nearest Neighbor



(a) 1-nearest neighbor



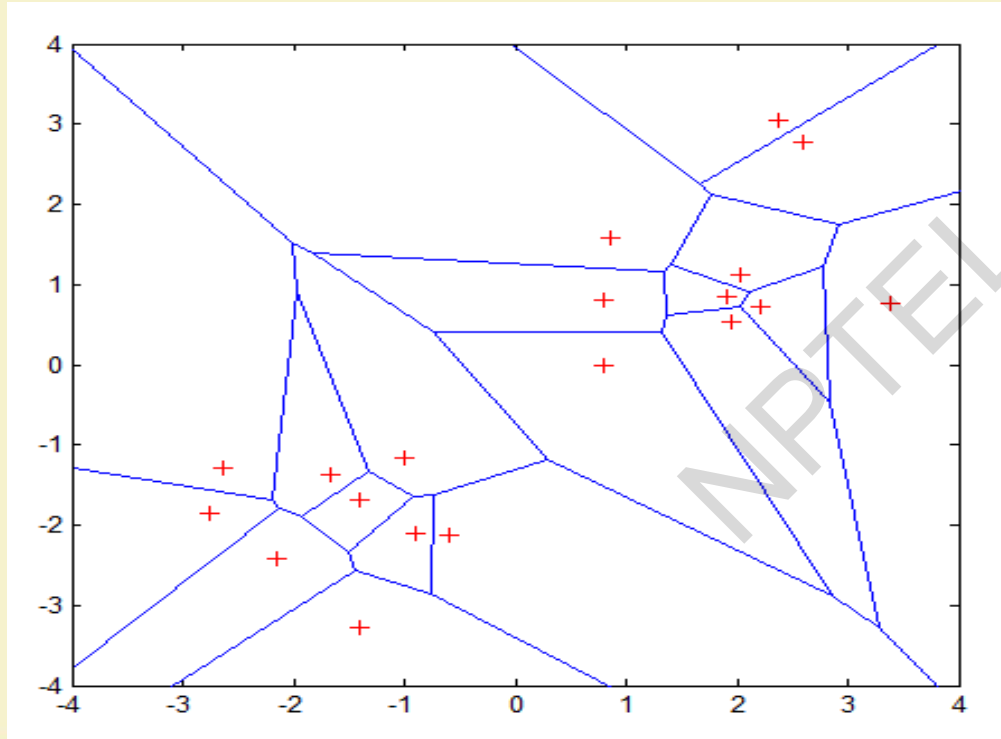
(b) 2-nearest neighbor



(c) 3-nearest neighbor

K-nearest neighbors of a record  $x$  - data points that have the  $k$  smallest distance to  $x$

# Nearest Neighbor: Voronoi Diagram



Properties:

- 1) All possible points within a sample's Voronoi cell are the nearest neighboring points for that sample
- 2) For any sample, the nearest sample is determined by the closest Voronoi cell edge

# Distance-weighted $k$ -NN

- Replace  $\hat{f}(q) = \arg \max_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$  by:

$$\hat{f}(q) = \arg \max_{v \in V} \sum_{i=1}^k \frac{1}{d(x_i, x_q)^2} \delta(v, f(x_i))$$

General Kernel functions like Parzen Windows may be considered Instead of inverse distance.



# Predicting Continuous Values

- Replace  $\hat{f}(q) = \arg \max_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i))$  by:

- Note: unweighted corresponds to  $w_i=1$  for all  $i$   
$$\hat{f}(q) = \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}$$

# Nearest-Neighbor Classifiers: Issues

- The value of  $k$ , the number of nearest neighbors to retrieve
- Choice of Distance Metric to compute distance between records
- Computational complexity
  - Size of training set
  - Dimension of data

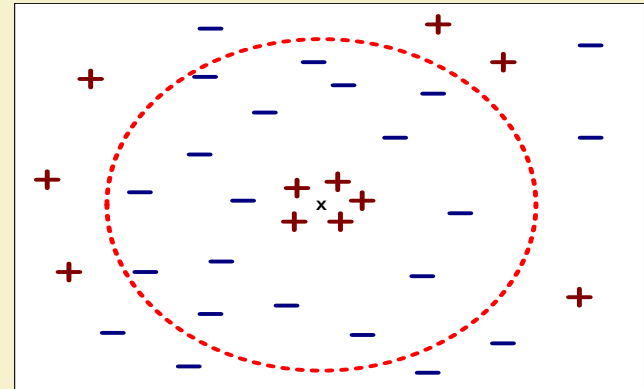
# Value of K

- Choosing the value of k:
  - If k is too small, sensitive to noise points
  - If k is too large, neighborhood may include points from other classes

Rule of thumb:

$K = \sqrt{N}$

N: number of training points



# Distance Metrics

**Minkowsky:**

$$D(\mathbf{x}, \mathbf{y}) = \left( \sum_{i=1}^m |x_i - y_i|^r \right)^{1/r}$$

**Euclidean:**

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

**Manhattan / city-block:**

$$D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m |x_i - y_i|$$

**Camberra:**

$$D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|}$$

**Chebychev:**

$$D(\mathbf{x}, \mathbf{y}) = \max_{i=1}^m |x_i - y_i|$$

**Quadratic:**

$$D(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{Q} (\mathbf{x} - \mathbf{y}) = \sum_{j=1}^m \left( \sum_{i=1}^m (x_i - y_i) q_{ji} \right) (x_j - y_j)$$

$\mathbf{Q}$  is a problem-specific positive definite  $m \times m$  weight matrix

**Mahalanobis:**

$$D(\mathbf{x}, \mathbf{y}) = [\det \mathbf{V}]^{1/m} (\mathbf{x} - \mathbf{y})^T \mathbf{V}^{-1} (\mathbf{x} - \mathbf{y})$$

$\mathbf{V}$  is the covariance matrix of  $A_1..A_m$ , and  $A_j$  is the vector of values for attribute  $j$  occurring in the training set instances  $1..n$ .

**Correlation:**

$$D(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^m (x_i - \bar{x}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^m (x_i - \bar{x}_i)^2 \sum_{i=1}^m (y_i - \bar{y}_i)^2}}$$

$\bar{x}_i = \bar{y}_i$  and is the average value for attribute  $i$  occurring in the training set.

**Chi-square:**

$$D(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \frac{1}{sum_i} \left( \frac{x_i}{size_x} - \frac{y_i}{size_y} \right)^2$$

$sum_i$  is the sum of all values for attribute  $i$  occurring in the training set, and  $size_x$  is the sum of all values in the vector  $\mathbf{x}$ .

**Kendall's Rank Correlation:**

$$D(\mathbf{x}, \mathbf{y}) = 1 - \frac{2}{n(n-1)} \sum_{i=1}^m \sum_{j=1}^{i-1} \text{sign}(x_i - x_j) \text{sign}(y_i - y_j)$$

$\text{sign}(x) = -1, 0$  or  $1$  if  $x < 0$ ,  $x = 0$ , or  $x > 0$ , respectively.

Figure 1. Equations of selected distance functions.  
( $\mathbf{x}$  and  $\mathbf{y}$  are vectors of  $m$  attribute values).

# Distance Measure: Scale Effects

- Different features may have different measurement scales
  - E.g., patient weight in kg (range [50,200]) vs. blood protein values in ng/dL (range [-3,3])
- Consequences
  - Patient weight will have a much greater influence on the distance between samples
  - May bias the performance of the classifier

# Standardization

- Transform raw feature values into z-scores

$$z_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

$x_{ij}$  is the value for the  $i^{th}$  sample and  $j^{th}$  feature

$\mu_j$  is the average of all  $x_{ij}$  for feature  $j$

$\sigma_j$  is the standard deviation of all  $x_{ij}$  over all input samples

- Range and scale of z-scores should be similar (providing distributions of raw feature values are alike)

## Nearest Neighbor : Dimensionality

- Problem with Euclidean measure:
  - High dimensional data
    - **curse of dimensionality**
  - Can produce counter-intuitive results
  - Shrinking density – sparsification effect

1	1	1	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	1	1

$d = 1.4142$

VS

1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	1

$d = 1.4142$

# Distance for Nominal Attributes

## Value Difference Metric (VDM)

[Stanfill & Waltz, 1986]

Providing appropriate distance measurements for nominal attributes.

$$vdm_a(x, y) = \sum_{c=1}^C \left( \frac{N_{a,x,c}}{N_{a,x}} - \frac{N_{a,y,c}}{N_{a,y}} \right)^2$$

$N_{a,x}$  = # times attribute  $a$  had value  $x$

$N_{a,x,c}$  = # times attribute  $a$  had value  $x$  and class was  $c$

$C$  = # output classes

Two values are considered closer if they have more similar classifications, i.e., if they have more similar correlations with the output classes.



# Distance for Heterogeneous Data

In this section, we define a heterogeneous distance function *HVDM* that returns the distance between two input vectors  $\mathbf{x}$  and  $\mathbf{y}$ . It is defined as follows:

$$HVDM(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{a=1}^m d_a^2(x_a, y_a)} \quad (11)$$

where  $m$  is the number of attributes. The function  $d_a(x, y)$  returns a distance between the two values  $x$  and  $y$  for attribute  $a$  and is defined as:

$$d_a(x, y) = \begin{cases} 1, & \text{if } x \text{ or } y \text{ is unknown; otherwise...} \\ \text{normalized\_vdm}_a(x, y), & \text{if } a \text{ is nominal} \\ \text{normalized\_diff}_a(x, y), & \text{if } a \text{ is linear} \end{cases} \quad (12)$$

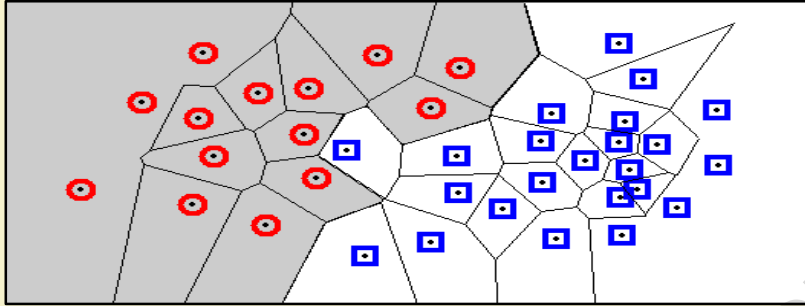
# Nearest Neighbour : Computational Complexity

- Expensive
  - To determine the nearest neighbour of a query point  $q$ , must compute the distance to all  $N$  training examples
    - + Pre-sort training examples into fast data structures (kd-trees)
    - + Compute only an approximate distance (LSH)
    - + Remove redundant data (condensing)
- Storage Requirements
  - Must store all training data  $P$ 
    - + Remove redundant data (condensing)
    - Pre-sorting often increases the storage requirements
- High Dimensional Data
  - “Curse of Dimensionality”
    - Required amount of training data increases exponentially with dimension
    - Computational cost also increases dramatically
    - Partitioning techniques degrade to linear search in high dimension

## Reduction in Computational Complexity

- Reduce size of training set
  - Condensation, editing
- Use geometric data structure for high dimensional search

# Condensation: Decision Regions



Each cell contains one sample, and every location within the cell is closer to that sample than to any other sample.

A Voronoi diagram divides the space into such cells.

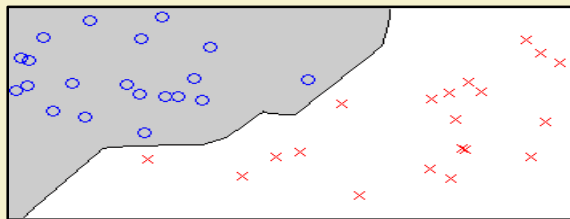
Every query point will be assigned the classification of the sample within that cell. The *decision boundary* separates the class regions based on the 1-NN decision rule.

Knowledge of this boundary is sufficient to classify new points.

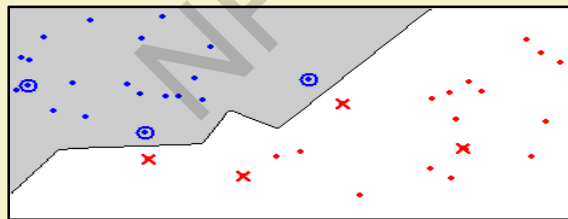
The boundary itself is rarely computed; many algorithms seek to retain only those points necessary to generate an identical boundary.

# Condensing

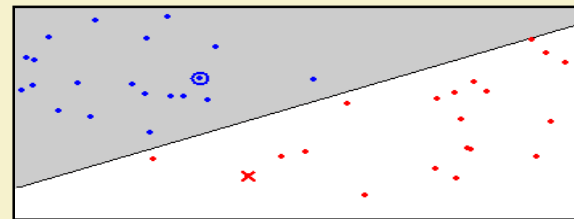
- Aim is to reduce the number of training samples
- Retain only the samples that are needed to define the decision boundary
- Decision Boundary Consistent – a subset whose nearest neighbour decision boundary is identical to the boundary of the entire training set
- Minimum Consistent Set – the smallest subset of the training data that correctly classifies all of the original training data



Original data



Condensed data



Minimum Consistent Set

# Condensed Nearest Neighbor

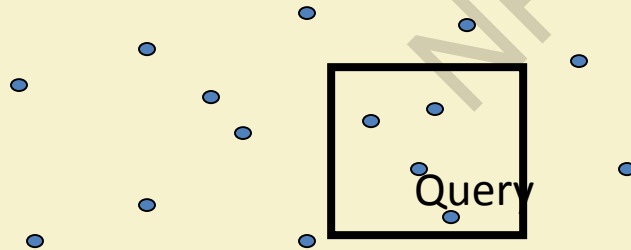
- Condensed Nearest Neighbour (CNN)

1. Initialize subset with a single (or K) training example
2. Classify all remaining samples using the subset, and transfer any incorrectly classified samples to the subset
3. Return to 2 until no transfers occurred or the subset is full

- Incremental
- Order dependent
- Neither minimal nor decision boundary consistent
- $O(n^3)$  for brute-force method

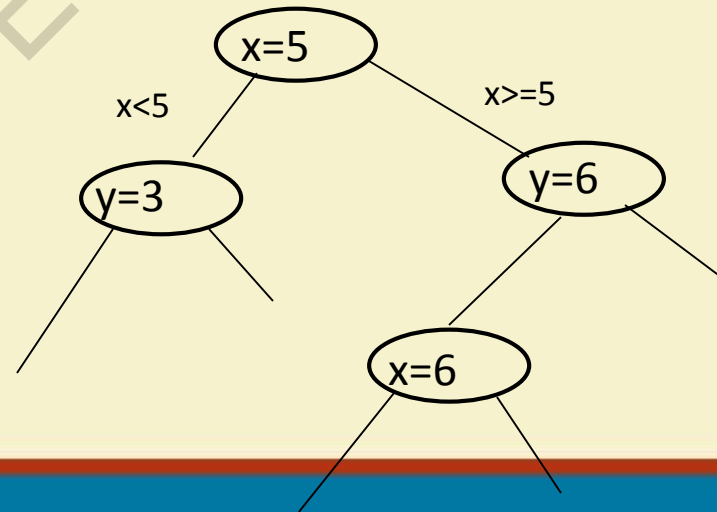
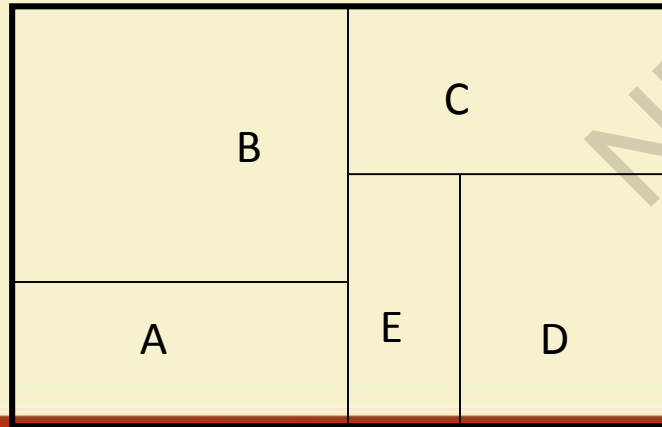
# High dimensional search

- Given a point set and a nearest neighbor query point
- Find the points enclosed in a rectangle (range) around the query
- Perform linear search for nearest neighbor only in the rectangle



# kd-tree: data structure for range search

- Index data into a tree
- Search on the tree
- Tree construction: At each level we use a different dimension to split





# KNN: Alternate Terminologies

- Instance Based Learning
- Lazy Learning
- Case Based Reasoning
- Exemplar Based Learning

# Text Search: Documents as vectors

- We have a  $|V|$ -dimensional vector space
- Terms are axes of the space
- Documents are points or vectors in this space
- Very high-dimensional: tens of millions of dimensions when you apply this to a web search engine
- These are very sparse vectors - most entries are zero.

# Queries as vectors

- [Key idea 1:](#) Do the same for queries: represent them as vectors in the space
- [Key idea 2:](#) Rank documents according to their proximity to the query in this space
- proximity = similarity of vectors
- proximity  $\approx$  inverse of distance
- **Recall: We do this because we want to get away from the you're-either-in-or-out Boolean model.**
- Instead: rank more relevant documents higher than less relevant documents

# Formalizing vector space proximity

- First cut: distance between two points
  - ( = distance between the end points of the two vectors)
- Euclidean distance?
- Euclidean distance is a bad idea . . .
- . . . because Euclidean distance is large for vectors of different lengths.

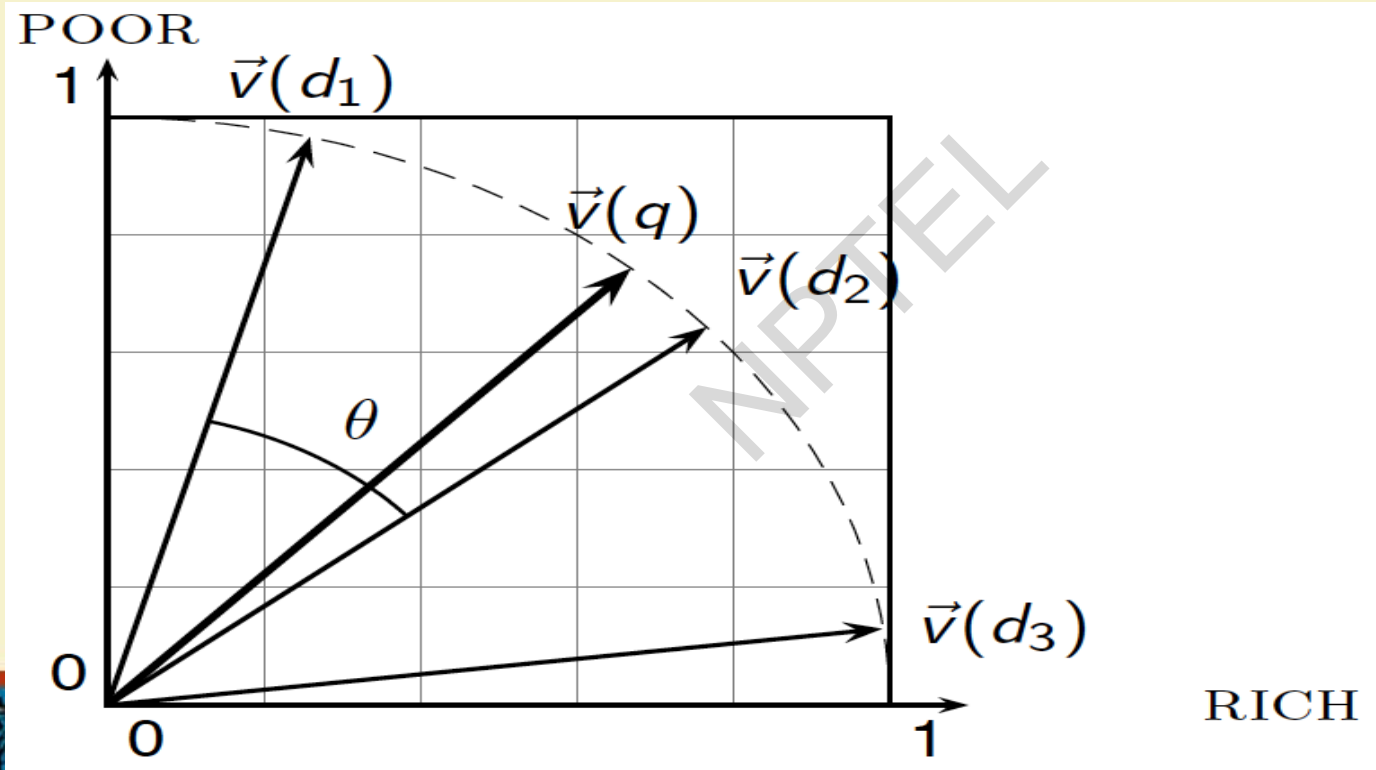
# Use angle instead of distance

- Thought experiment: take a document  $d$  and append it to itself. Call this document  $d'$ .
- “Semantically”  $d$  and  $d'$  have the same content
- The Euclidean distance between the two documents can be quite large
- The angle between the two documents is 0, corresponding to maximal similarity.
- Key idea: Rank documents according to angle with query.

# From angles to cosines

- The following two notions are equivalent.
  - Rank documents in decreasing order of the angle between query and document
  - Rank documents in increasing order of  $\cos(\text{angle}(\text{query}, \text{document}))$
- Cosine is a monotonically decreasing function for the interval  $[0^\circ, 180^\circ]$

# Cosine similarity illustrated



# End of K-Nearest Neighbor



# Classifier Evaluation

# Classifier Evaluation

- **Metrics for Performance Evaluation**
  - How to evaluate the performance of a model
- **Methods for Performance Evaluation**
  - How to obtain reliable estimates?
- **Methods for Model Comparison**
  - How to compare the relative performance among competing models?

# Metrics for Performance Evaluation

- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix:

	PREDICTED CLASS		
		Class=Yes	Class=No
	ACTUAL CLASS	Class=Yes	Class=No
		Class=No	Class=No
	Class=Yes	a	b
	Class=No	c	d

- a: TP (true positive)  
b: FN (false negative)  
c: FP (false positive)  
d: TN (true negative)

# Accuracy

ACTUAL CLASS	PREDICTED CLASS		
		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

# Limitation of Accuracy

- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is  $9990/10000 = 99.9\%$ 
  - Accuracy is misleading because model does not detect any class 1 example

# Cost Matrix

	PREDICTED CLASS		
	$C(i j)$	Class=Yes	Class=No
	Class=Yes	$C(\text{Yes} \text{Yes})$	$C(\text{No} \text{Yes})$
	Class=No	$C(\text{Yes} \text{No})$	$C(\text{No} \text{No})$

$C(i|j)$ : Cost of misclassifying class  $j$  example as class  $i$

## Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

- Precision is biased towards C(Yes|Yes) & C(Yes|No)
- Recall is biased towards C(Yes|Yes) & C(No|Yes)
- F-measure is biased towards all except C(No|No)

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

# Model Evaluation

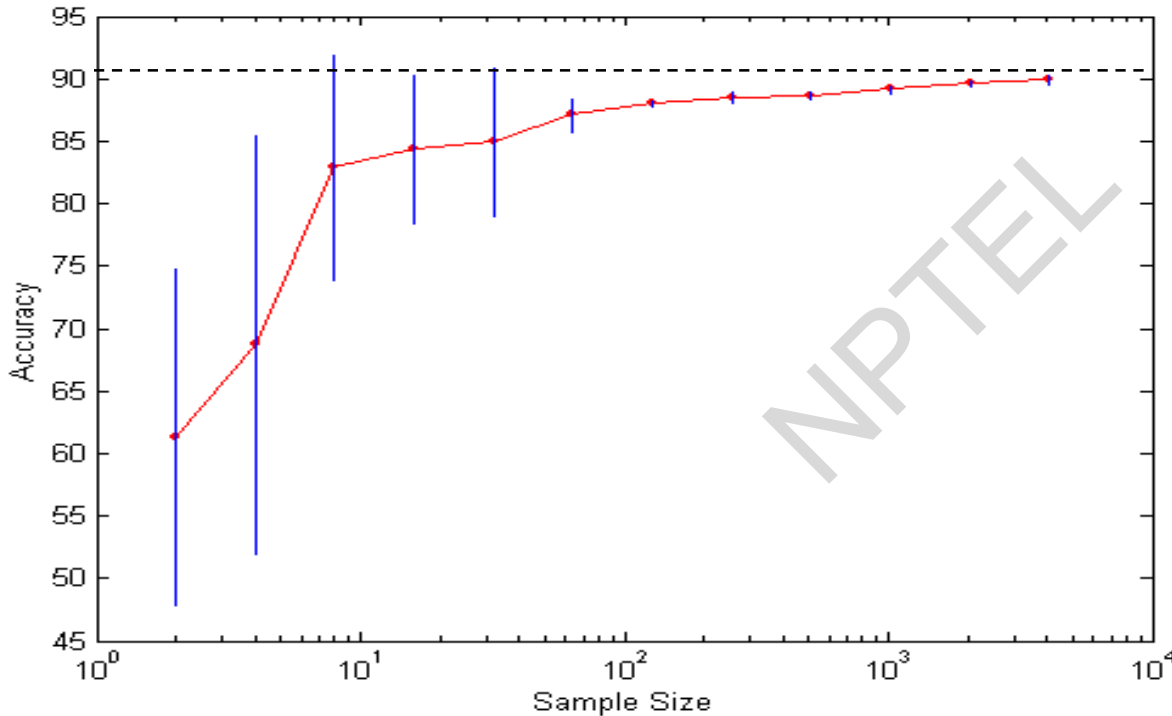
- Metrics for Performance Evaluation
  - How to evaluate the performance of a model
- **Methods for Performance Evaluation**
  - How to obtain reliable estimates?
- Methods for Model Comparison
  - How to compare the relative performance among competing models?



# Methods for Performance Evaluation

- How to obtain a reliable estimate of performance?
- Performance of a model may depend on other factors besides the learning algorithm:
  - Class distribution
  - Cost of misclassification
  - Size of training and test sets

# Learning Curve



- Learning curve shows how accuracy changes with varying sample size

Effect of small sample size:

- Bias in the estimate
- Variance of estimate

# Methods of Estimation

- Holdout
  - Reserve  $2/3$  for training and  $1/3$  for testing
- Random subsampling
  - Repeated holdout
- Cross validation
  - Partition data into  $k$  disjoint subsets
  - $k$ -fold: train on  $k-1$  partitions, test on the remaining one
  - Leave-one-out:  $k=n$
- Stratified sampling
  - oversampling vs undersampling
- Bootstrap
  - Sampling with replacement

# Model Evaluation

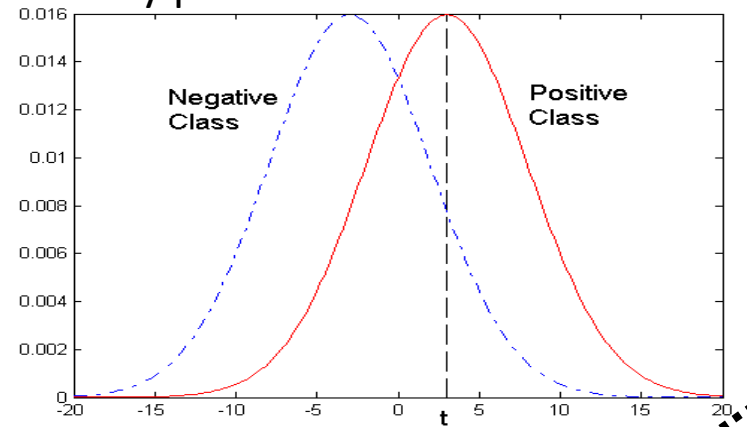
- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?
- Methods for Performance Evaluation
  - How to obtain reliable estimates?
- **Methods for Model Comparison**
  - How to compare the relative performance among competing models?

# ROC (Receiver Operating Characteristic)

- Developed in 1950s for signal detection theory to analyze noisy signals
  - Characterize the trade-off between positive hits and false alarms
- ROC curve plots TP (on the y-axis) against FP (on the x-axis)
- Performance of each classifier represented as a point on the ROC curve
  - changing the threshold of algorithm, sample distribution or cost matrix changes the location of the point

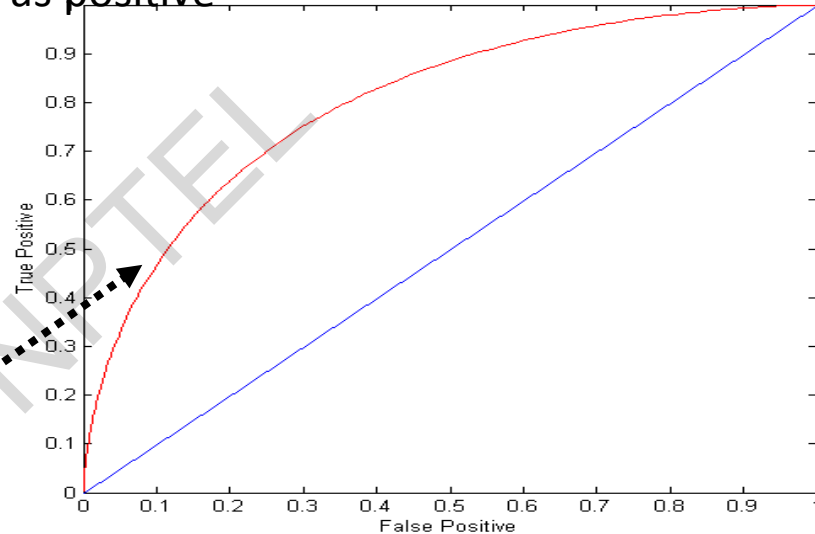
# ROC Curve

- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at  $x > t$  is classified as positive



At threshold  $t$ :

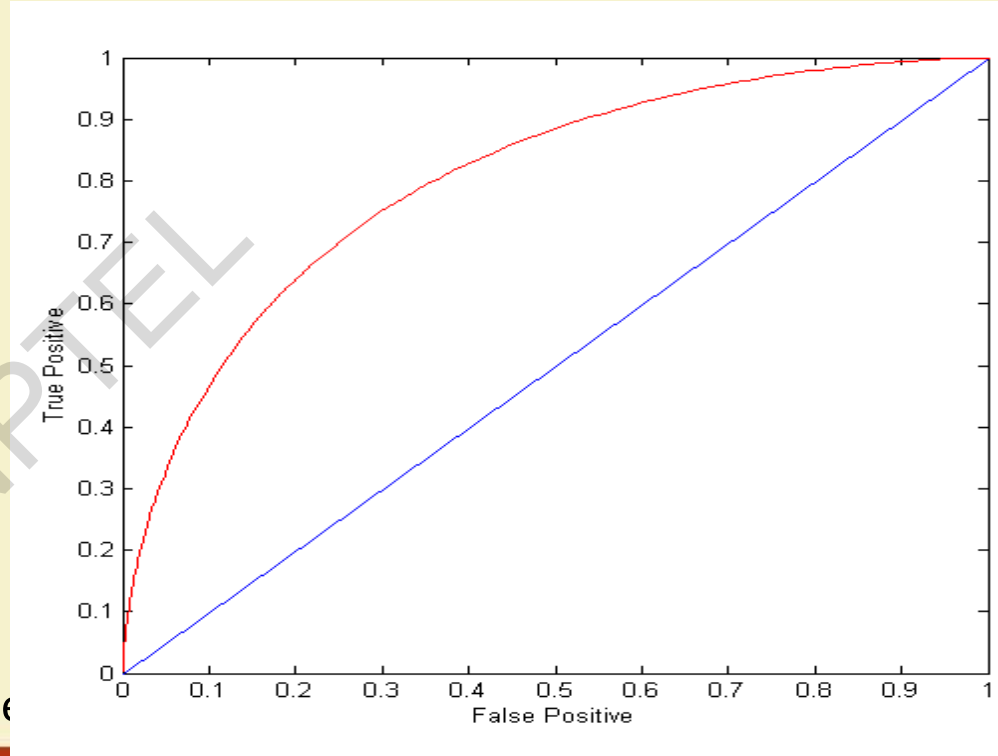
TP=0.5, FN=0.5, FP=0.12, FN=0.88



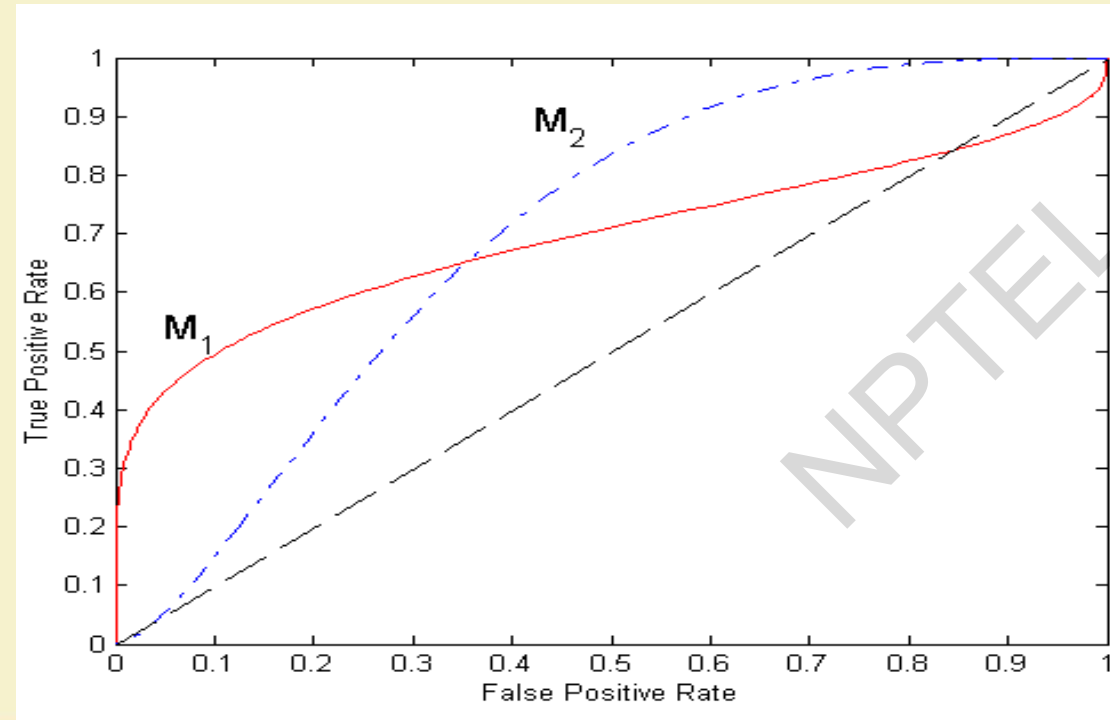
# ROC Curve

(TP,FP):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal
- Diagonal line:
  - Random guessing
  - Below diagonal line:
    - prediction is opposite of the true class



# Using ROC for Model Comparison



- No model consistently outperform the other
  - $M_1$  is better for small FPR
  - $M_2$  is better for large FPR
- Area Under the ROC curve
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5



# Test of Significance

- Given two models:
  - Model M1: accuracy = 85%, tested on 30 instances
  - Model M2: accuracy = 75%, tested on 5000 instances
- Can we say M1 is better than M2?
  - How much confidence can we place on accuracy of M1 and M2?
  - Can the difference in performance measure be explained as a result of random fluctuations in the test set?

# Comparing Performance of 2 Models

- Given two models, say M1 and M2, which is better?
  - M1 is tested on D1 (size= $n_1$ ), found error rate =  $e_1$
  - M2 is tested on D2 (size= $n_2$ ), found error rate =  $e_2$
  - Assume D1 and D2 are independent
  - If  $n_1$  and  $n_2$  are sufficiently large, then

$$e_1 \sim N(\mu_1, \sigma_1)$$

$$e_2 \sim N(\mu_2, \sigma_2)$$

- Approximate: 
$$\hat{\sigma}_i = \frac{e_i(1-e_i)}{n_i}$$

# Comparing Performance of 2 Models

- To test if performance difference is statistically significant:  $d = e1 - e2$ 
  - $d \sim N(d_t, \sigma_t)$  where  $d_t$  is the true difference
  - Since D1 and D2 are independent, their variance adds up:

$$\begin{aligned}\sigma_t^2 &= \sigma_1^2 + \sigma_2^2 \cong \hat{\sigma}_1^2 + \hat{\sigma}_2^2 \\ &= \frac{e1(1-e1)}{n1} + \frac{e2(1-e2)}{n2}\end{aligned}$$

- At  $(1-\alpha)$  confidence level,  $d_t = d \pm Z_{\alpha/2} \hat{\sigma}_t$

# End of Classifier Evaluation