# Data Mining
# Classification: Basic Concepts, Decision Trees, and Model Evaluation

## Lecture Notes for Chapter 4

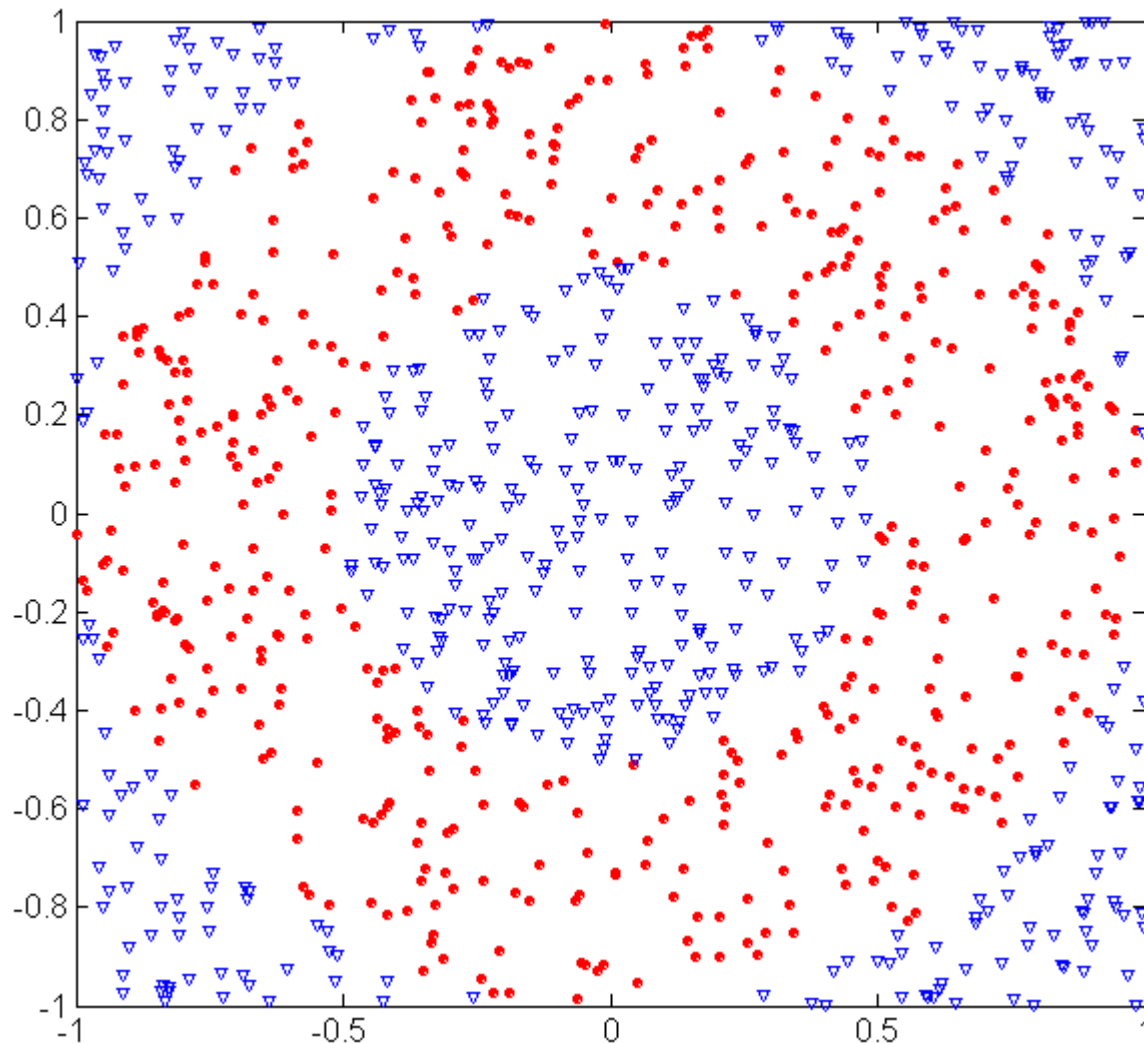## Part III

## Introduction to Data Mining

## by

## Tan, Steinbach, Kumar

Adapted by Qiang Yang (2010)

# Practical Issues of Classification

l  Underfitting and Overfitting

l  Missing Values

l  Costs of Classification

# Underfitting and Overfitting (Example)



500 circular and 500 triangular data points.
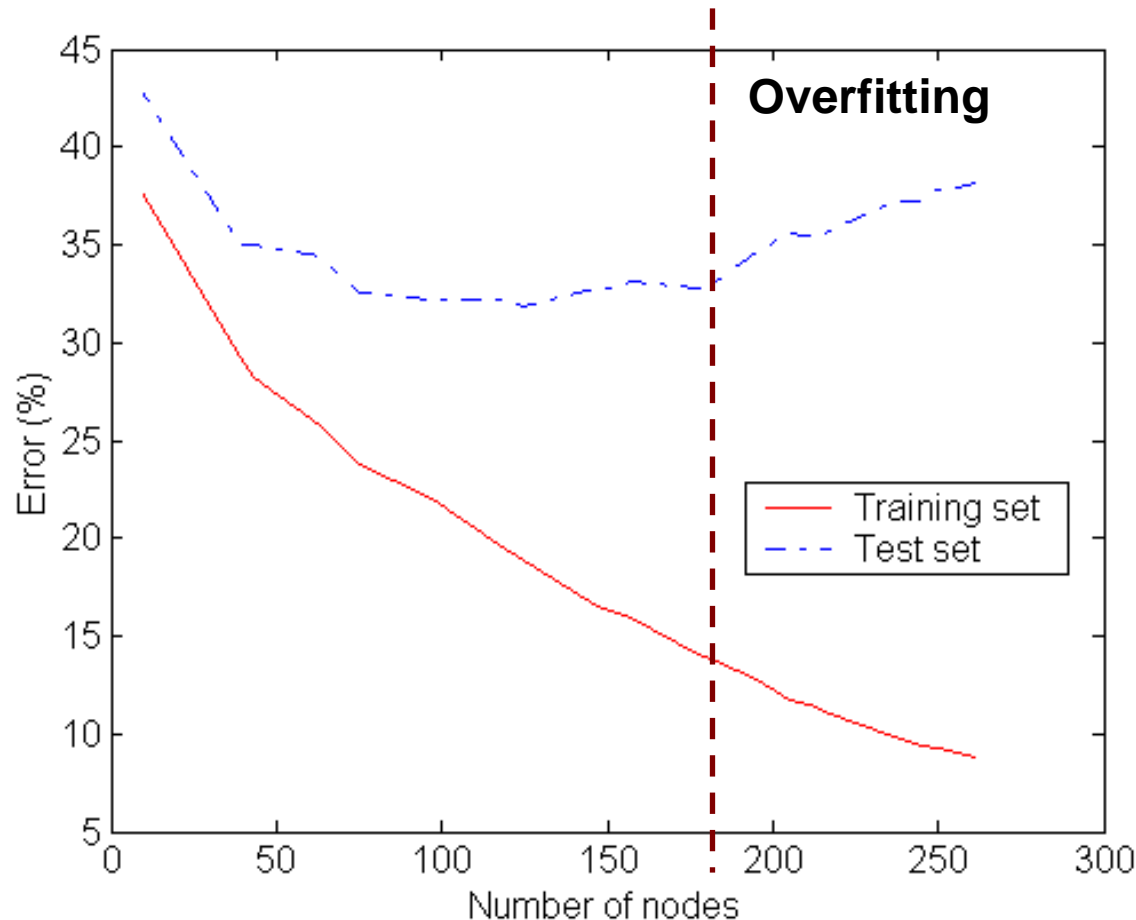
**Circular points:**

$0.5 \leq \text{sqrt}(x_1^2 + x_2^2) \leq 1$

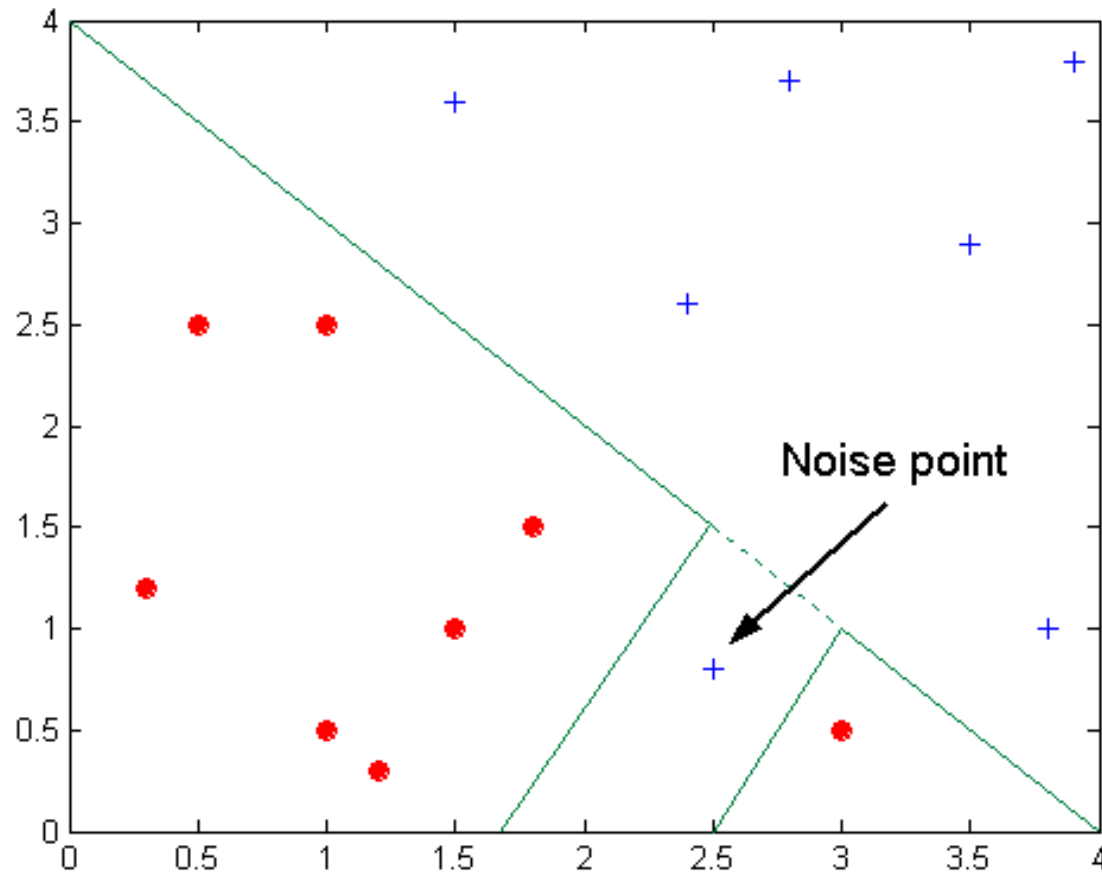**Triangular points:**

$\text{sqrt}(x_1^2 + x_2^2) > 0.5$ or

$\text{sqrt}(x_1^2 + x_2^2) < 1$

# Underfitting and Overfitting



**Underfitting**: when model is too simple, both training and test errors are large

# Overfitting due to Noise



Noise point

**Decision boundary is distorted by noise point**

# Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary

- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records

- Need new ways for estimating errors
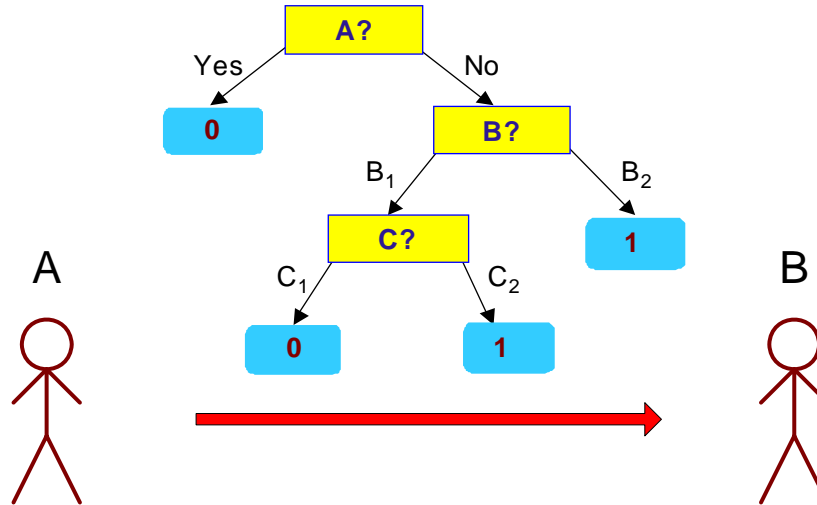
# Estimating Generalization Errors

- Re-substitution errors: error on training ($\Sigma\ e(t)$ )
- Generalization errors: error on testing ($\Sigma\ e'(t)$)

- Methods for estimating generalization errors:
  - Optimistic approach:  $e'(t) = e(t)$
  - Pessimistic approach:
    - For each leaf node: $e'(t) = (e(t)+0.5)$
    - Total error counts: $e'(T) = e(T) + N \times 0.5$ (N: number of leaf nodes)
    - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
      Training error = 10/1000 = 1%
      **Generalization error** = $(10 + 30 \times 0.5)/1000 = 2.5\%$
  - Reduced error pruning (REP):
    - uses validation data set to estimate generalization error

# Occam's Razor

l  Given two models of similar generalization errors, one should prefer the simpler model over the more complex model

   – For complex models, there is a greater chance that it was fitted accidentally by errors in data

   – Therefore, one should include model complexity when evaluating a model

# Minimum Description Length (MDL)



| X | y |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0 |
| $X_3$ | 0 |
| $X_4$ | 1 |
| … | … |
| $X_n$ | 1 |

| X | y |
|---|---|
| $X_1$ | ? |
| $X_2$ | ? |
| $X_3$ | ? |
| $X_4$ | ? |
| … | … |
| $X_n$ | ? |

- **Cost(Model,Data) = Cost(Data|Model) + Cost(Model)**
  - Cost is the number of bits needed for encoding.
  - We should search for the least costly model.
- Cost(Data|Model) encodes the errors on training data.
- Cost(Model) estimates model complexity, or future error…

# How to Address Overfitting in Decision Trees

l **Pre-Pruning (Early Stopping Rule)**

- Stop the algorithm before it becomes a fully-grown tree

- Typical stopping conditions for a node:

    ◆ Stop if all instances belong to the same class

    ◆ Stop if all the attribute values are the same

- More restrictive conditions:

    ◆ Stop if number of instances is less than some user-specified threshold

    ◆ Stop if class distribution of instances are independent of the available features (e.g., using $\chi^2$ test)

    ◆ Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).

# How to Address Overfitting…

l Post-pruning

- Grow decision tree to its entirety

- Trim the nodes of the decision tree in a bottom-up fashion

- **If generalization error improves after trimming, replace sub-tree by a leaf node.**

    - Heuristic: Class label of leaf node is determined from majority class of instances in the sub-tree

    - generalization error count = error count + 0.5*N, where N is the number of leaf nodes,

    - This is a heuristic used in some algorithms, but there are other ways using statistics

# Post-Pruning based on |leaves|

| Class = Yes | 20 |
|---|---|
| Class = No | 10 |
| Error = 10/30 | |

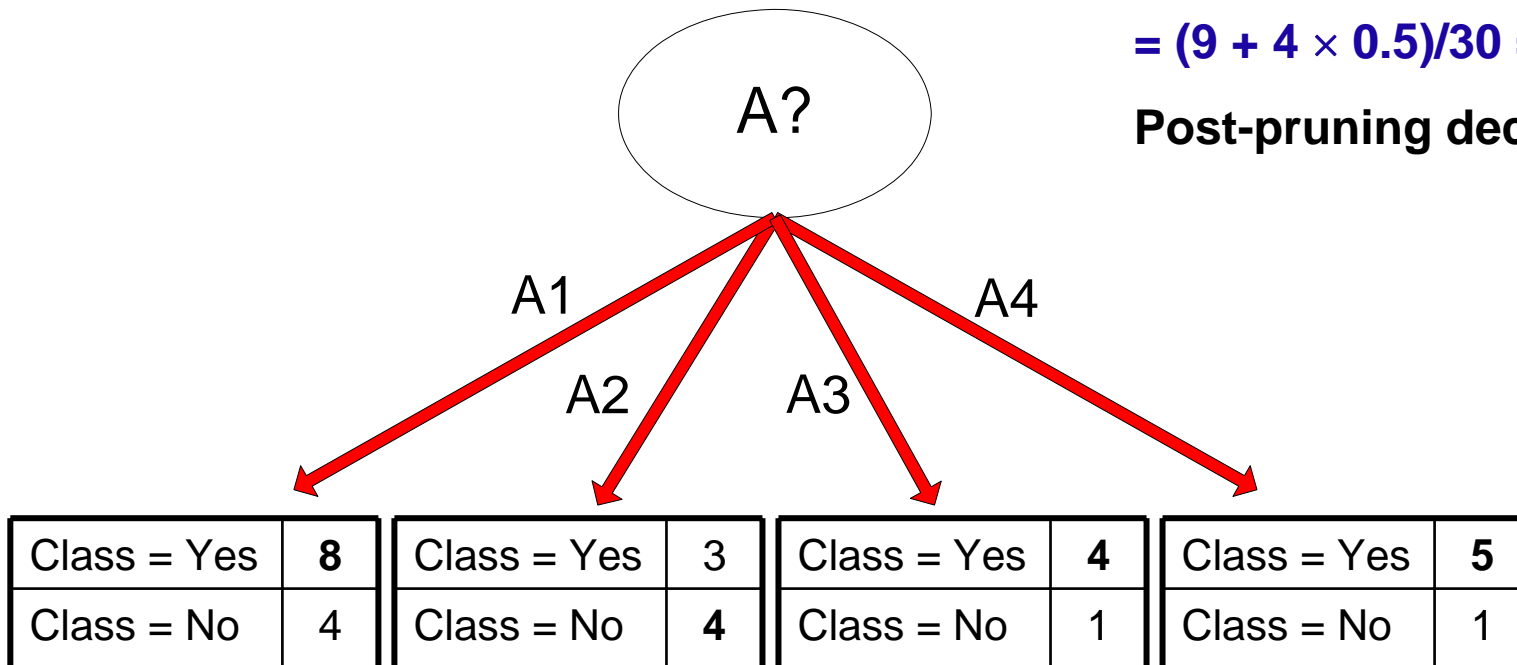**Training Error (Before splitting) = 10/30**

**Pessimistic error (Before splitting) = (10 + 1X 0.5)/30 = 10.5/30**

**Training Error (After splitting) = 9/30**

**Pessimistic error (After splitting)**

$$= (9 + 4 \times 0.5)/30 = 11/30$$

**Post-pruning decision: PRUNE!**



A?

A1    A2    A3    A4

| Class = Yes | **8** | | Class = Yes | 3 | | Class = Yes | **4** | | Class = Yes | **5** |
|---|---|---|---|---|---|---|---|---|---|---|
| Class = No | 4 | | Class = No | **4** | | Class = No | 1 | | Class = No | 1 |

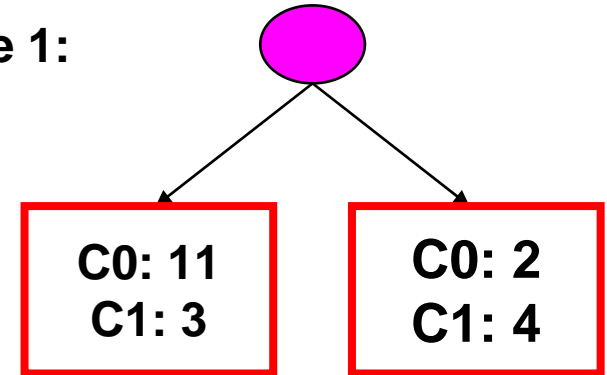# Examples of Post-pruning

– Optimistic error?

  **Don't prune for both cases**

– Pessimistic error?

  **Don't prune case 1, prune case 2**

**Case 1:**

C0: 11
C1: 3

C0: 2
C1: 4

**Case 2:**

C0: 14
C1: 3

C0: 2
C1: 2

# Data Fragmentation

- Number of instances gets smaller as you traverse down the tree

- Number of instances at the leaf nodes could be too small to make any statistically significant decision

- Solution: limit number of instances per leaf node >= a user given value n.

# Decision Trees: Feature Construction



- **Test condition may involve multiple attributes, but hard to automate!**

- **Finding better node test features is a difficult research issue**

# 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?

# Model Evaluation

- <span style="color:red">Metrics for Performance Evaluation</span>
  - 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: count or percentage

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

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

# Metrics for Performance Evaluation…

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

l **Most widely-used metric:**

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

# 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 | | |
|---|---|---|---|
| ACTUAL CLASS | C(i\|j) | **Class=Yes** | **Class=No** |
| | **Class=Yes** | C(Yes\|Yes) | C(No\|Yes) |
| | **Class=No** | C(Yes\|No) | C(No\|No) |

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

- medical diagnosis, customer segmentation

# Computing Cost of Classification

| Cost Matrix | PREDICTED CLASS | | |
|---|---|---|---|
| | C(i\|j) | **+** | **-** |
| ACTUAL CLASS **+** | | -1 | 100 |
| **-** | | 1 | 0 |

**Confusion matrix**

| Model $M_1$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | **+** | **-** |
| ACTUAL CLASS **+** | | 150 | 40 |
| **-** | | 60 | 250 |

| Model $M_2$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | **+** | **-** |
| ACTUAL CLASS **+** | | 250 | 45 |
| **-** | | 5 | 200 |

Accuracy = 80%

Cost = 3910

Accuracy = 90%

Cost = 4255

# Information Retrieval Measures

| | PREDICTED CLASS | | |
|---|---|---|---|
| **ACTUAL CLASS** | | Class=Yes | Class=No |
| | Class=Yes | a | b |
| | Class=No | c | d |

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

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

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

- Let C be cost (can be count in our example)
- 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)

# Model Evaluation

l Metrics for Performance Evaluation
  – How to evaluate the performance of a model?

l Methods for Performance Evaluation
  – How to obtain reliable estimates?

l Methods for Model Comparison
  – How to compare the relative performance among competing models?

# Methods of Estimation

l Holdout

– Reserve 2/3 for training and 1/3 for testing

l 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

# Test of Significance (Sections 4.5,4.6 of TSK Book)

- 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?

# Confidence Interval for Accuracy

l Prediction can be regarded as a Bernoulli trial

   – A Bernoulli trial has 2 possible outcomes

   – Possible outcomes for prediction: correct or wrong

   – Collection of Bernoulli trials has a Binomial distribution:

      ◆ $x \sim Bin(N, p)$      x: number of correct predictions

      ◆ e.g:  Toss a fair coin 50 times, how many heads would turn up?
         Expected number of heads = $N \times p = 50 \times 0.5 = 25$

l Given x (# of correct predictions) or equivalently, acc=x/N, and N =# of test instances,

   – Can we predict p (true accuracy of model)?
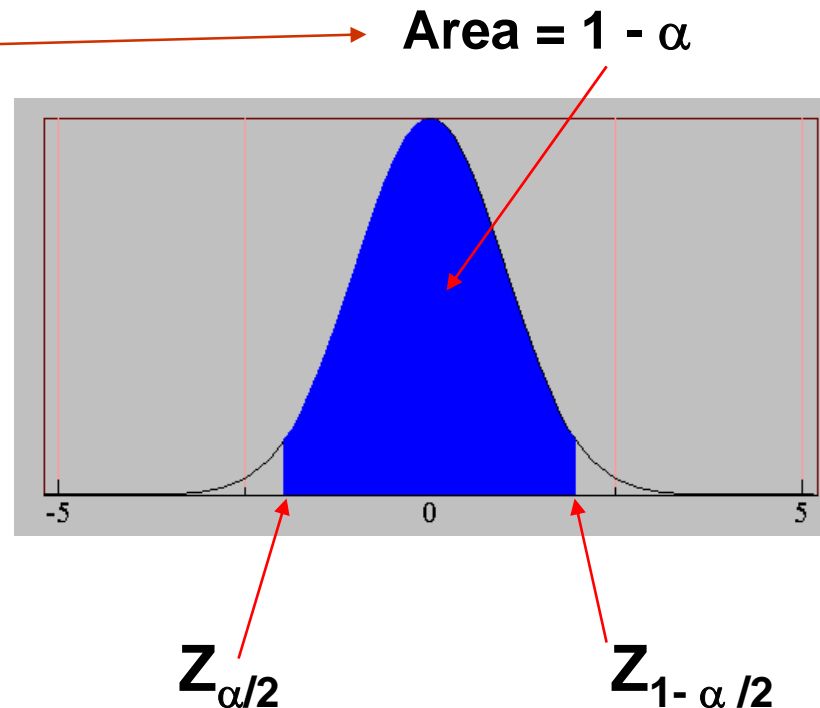
# Confidence Interval for Accuracy

- For large N, let $1-\alpha$ be confidence

  - *acc* has a normal distribution with mean p and variance p(1-p)/N

$$P(Z_{\alpha/2} < \frac{acc - p}{\sqrt{p(1-p)/N}} < Z_{1-\alpha/2})$$

$$= 1 - \alpha$$

**Area = 1 - $\alpha$**

$Z_{\alpha/2}$        $Z_{1-\alpha/2}$

- Confidence Interval for p:

$$p = \frac{2 \times N \times acc + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4 \times N \times acc - 4 \times N \times acc^2}}{2(N + Z_{\alpha/2}^2)}$$

# Confidence Interval for Accuracy

- Consider a model that produces an accuracy of 80% when evaluated on 100 test instances:

  - N=100, acc = 0.8

  - Let $1-\alpha = 0.95$ (95% confidence)

  - From probability table, $Z_{\alpha/2}=1.96$

| $1-\alpha$ | Z |
|---|---|
| 0.99 | 2.58 |
| 0.98 | 2.33 |
| 0.95 | 1.96 |
| 0.90 | 1.65 |

| N | 50 | 100 | 500 | 1000 | 5000 |
|---|---|---|---|---|---|
| p(lower) | 0.670 | 0.711 | 0.763 | 0.774 | 0.789 |
| p(upper) | 0.888 | 0.866 | 0.833 | 0.824 | 0.811 |

# ROC (Receiver Operating Characteristic)

- Page 298 of TSK book.
- Many applications care about ranking (give a queue from the most likely to the least likely)
- Examples…
- Which ranking order is better?
- ROC: 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

# How to Construct an ROC curve

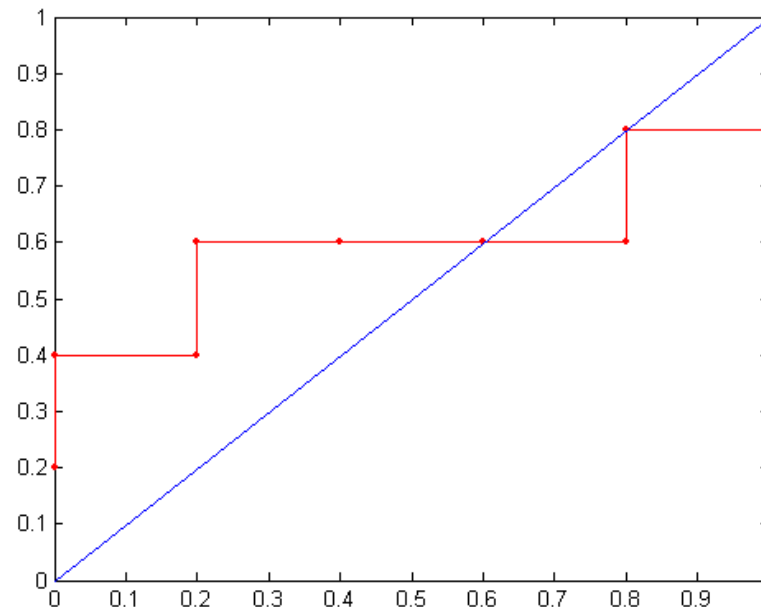| Instance | P(+|A) | True Class |
|----------|--------|------------|
| 1 | 0.95 | + |
| 2 | 0.93 | + |
| 3 | 0.87 | - |
| 4 | 0.85 | - |
| 5 | 0.85 | - |
| 6 | 0.85 | + |
| 7 | 0.76 | - |
| 8 | 0.53 | + |
| 9 | 0.43 | - |
| 10 | 0.25 | + |

**Predicted by classifier**

**This is the ground truth**

- Use classifier that produces posterior probability for each test instance P(+|A) for instance A

- Sort the instances according to P(+|A) in decreasing order

- Apply threshold at each unique value of P(+|A)

- Count the number of TP, FP, TN, FN at each threshold
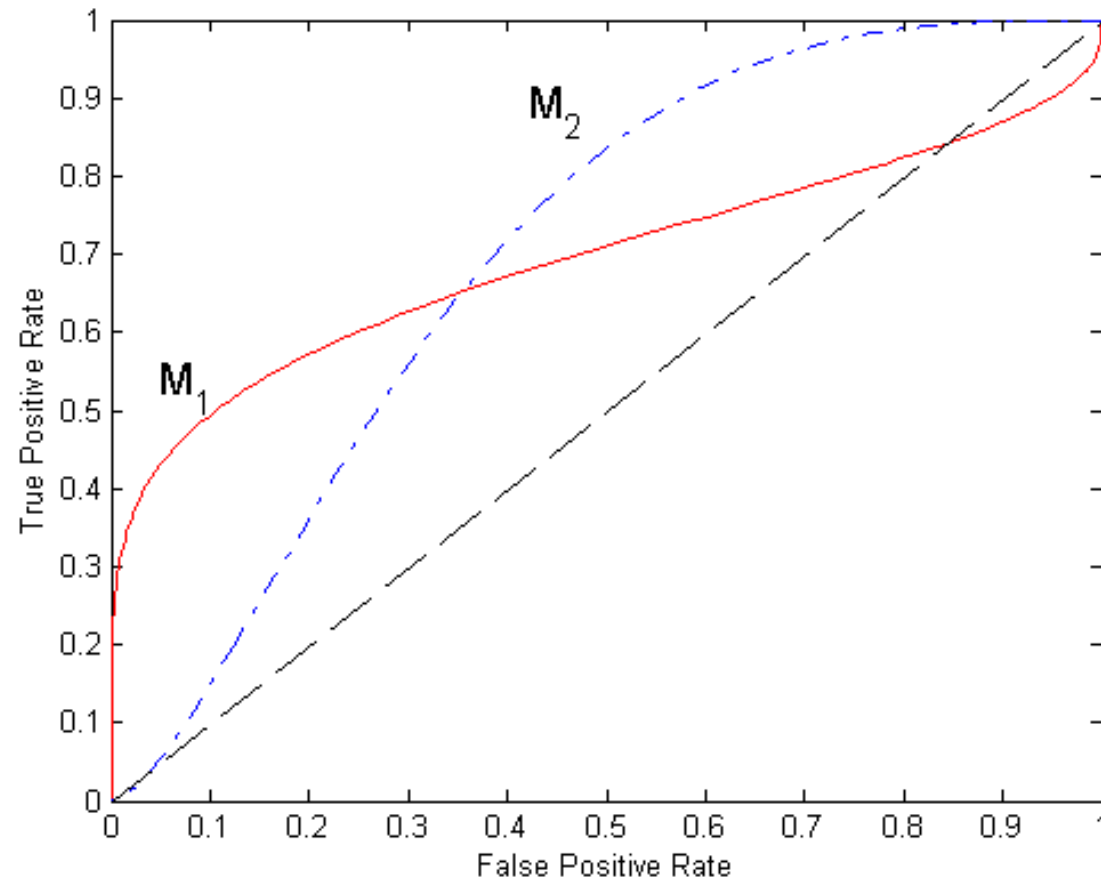
- TP rate, TPR = TP/(TP+FN)

- FP rate, FPR = FP/(FP + TN)

# How to construct an ROC curve

| Class | + | - | + | - | - | - | + | - | + | + | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Threshold >= | 0.25 | 0.43 | 0.53 | 0.76 | 0.85 | 0.85 | 0.85 | 0.87 | 0.93 | 0.95 | 1.00 |
| TP | 5 | 4 | 4 | 3 | 3 | 3 | 3 | 2 | 2 | 1 | 0 |
| FP | 5 | 5 | 4 | 4 | 3 | 2 | 1 | 1 | 0 | 0 | 0 |
| TN | 0 | 0 | 1 | 1 | 2 | 3 | 4 | 4 | 5 | 5 | 5 |
| FN | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 | 3 | 4 | 5 |
| TPR | 1 | 0.8 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.4 | 0.4 | 0.2 | 0 |
| FPR | 1 | 1 | 0.8 | 0.8 | 0.6 | 0.4 | 0.2 | 0.2 | 0 | 0 | 0 |

**ROC Curve:**

# 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: AUC
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5

# ROC Curve

(TP,FP):

l (0,0): declare everything
    to be negative class

l (1,1): declare everything
    to be positive class

l (1,0): ideal


l Diagonal line:

  – Random guessing

  – Below diagonal line:

    ◆ prediction is opposite of
    the true class