# Comparing Costs of Classification

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

| 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

# Cost vs Accuracy

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

| Cost | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | p | q |
| Class=No | q | p |

Accuracy is proportional to cost if
1. $C(Yes|No) = C(No|Yes) = q$
2. $C(Yes|Yes) = C(No|No) = p$

$N = a + b + c + d$

$Accuracy = (a + d)/N$

$$Cost = p(a + d) + q(b + c)$$
$$= p(a + d) + q(N - a - d)$$
$$= qN - (q - p)(a + d)$$
$$= N[q - (q-p) \times Accuracy]$$

# Cost-Sensitive Measures

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

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

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

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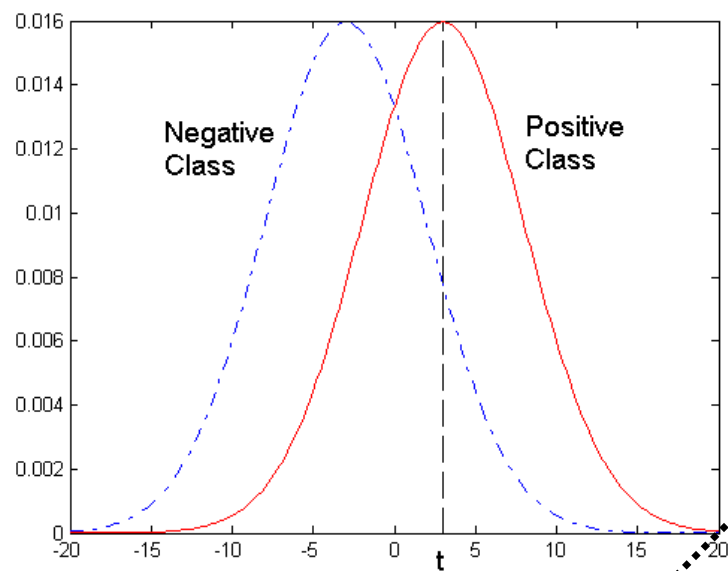
# Model Evaluation

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


- Methods for 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
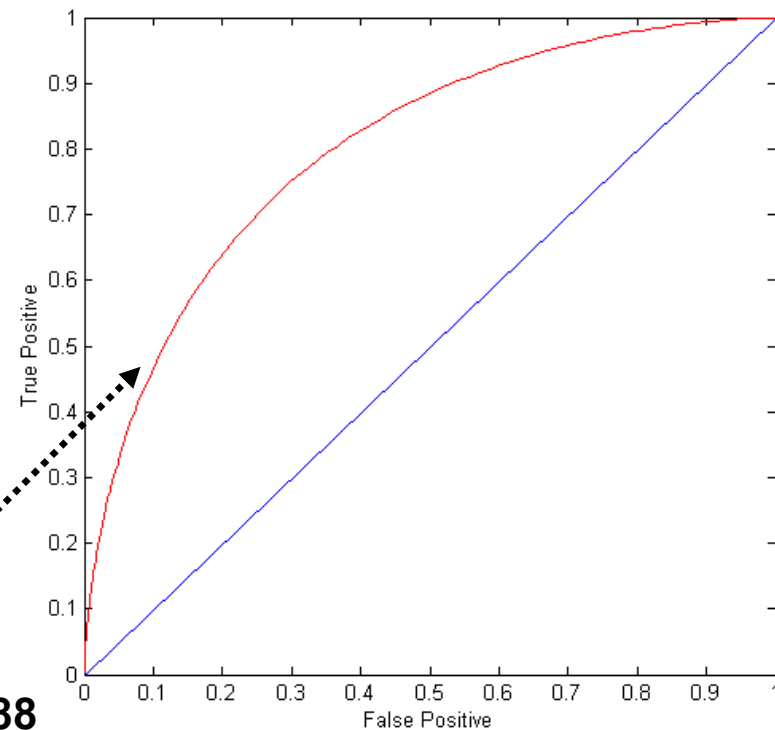
# ROC Curve

- 1-dimensional data set containing 2 classes (positive and negative)

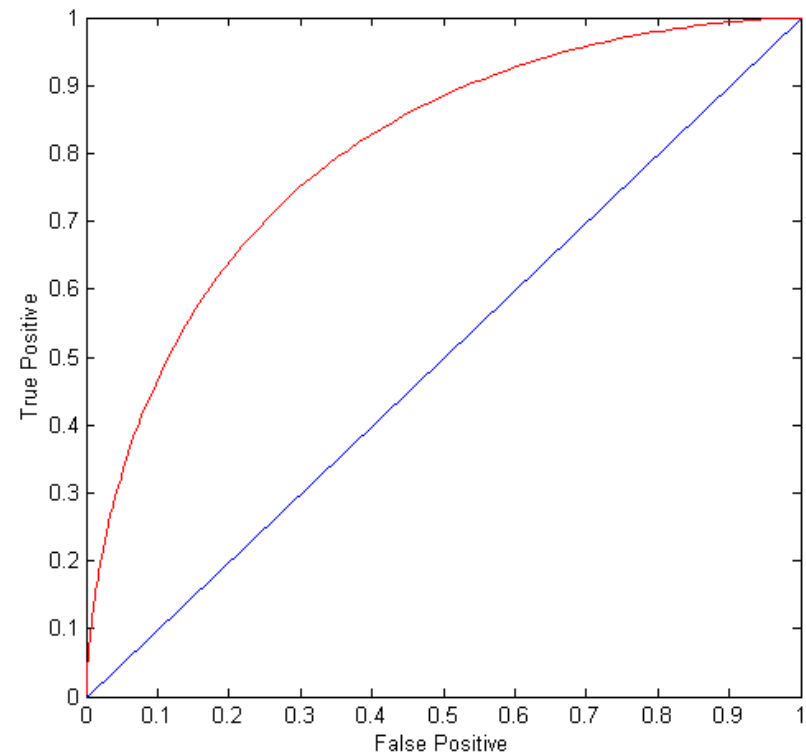- any points located at x > t is classified as positive



**At threshold t:**
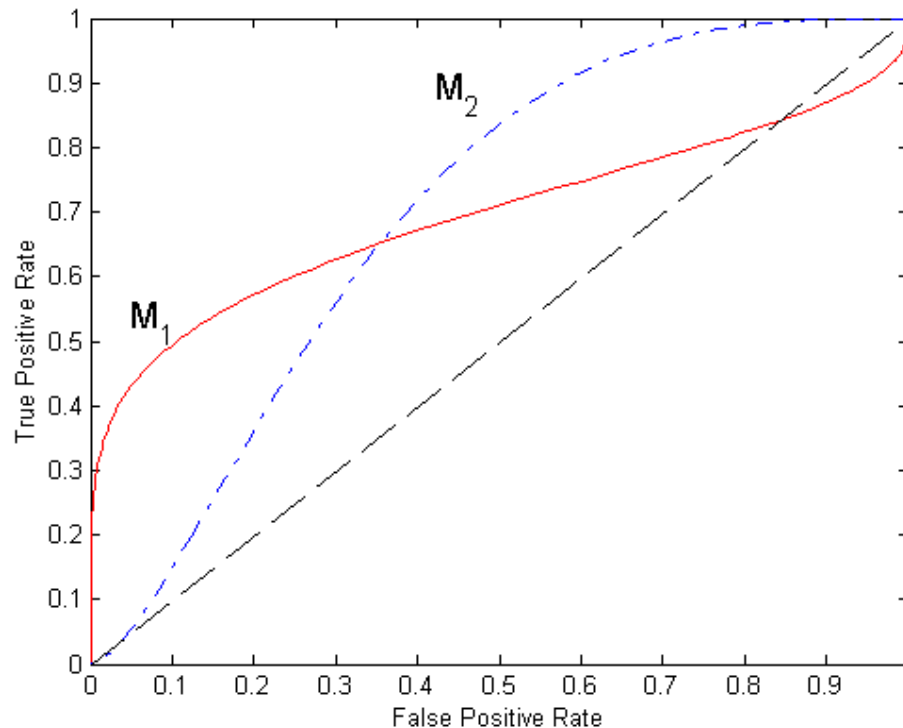
**TP=0.5, FN=0.5, FP=0.12, FN=0.88**

# 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 outperforms 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:
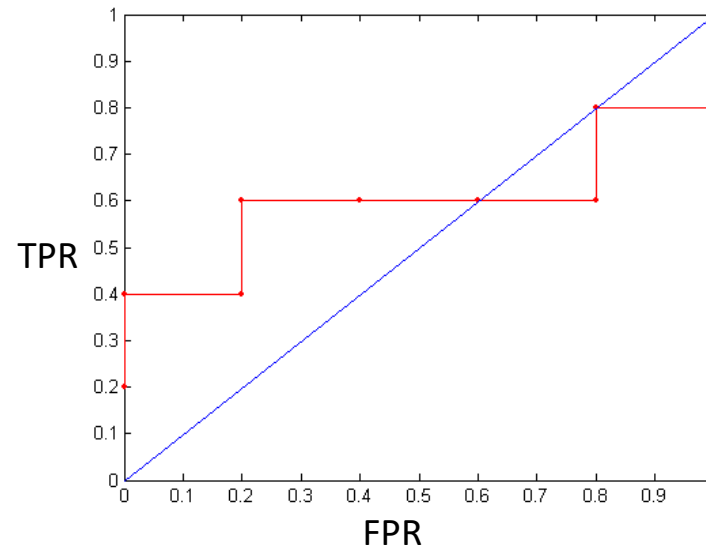    - ▪ Area = 0.5

# 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    | +          |

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

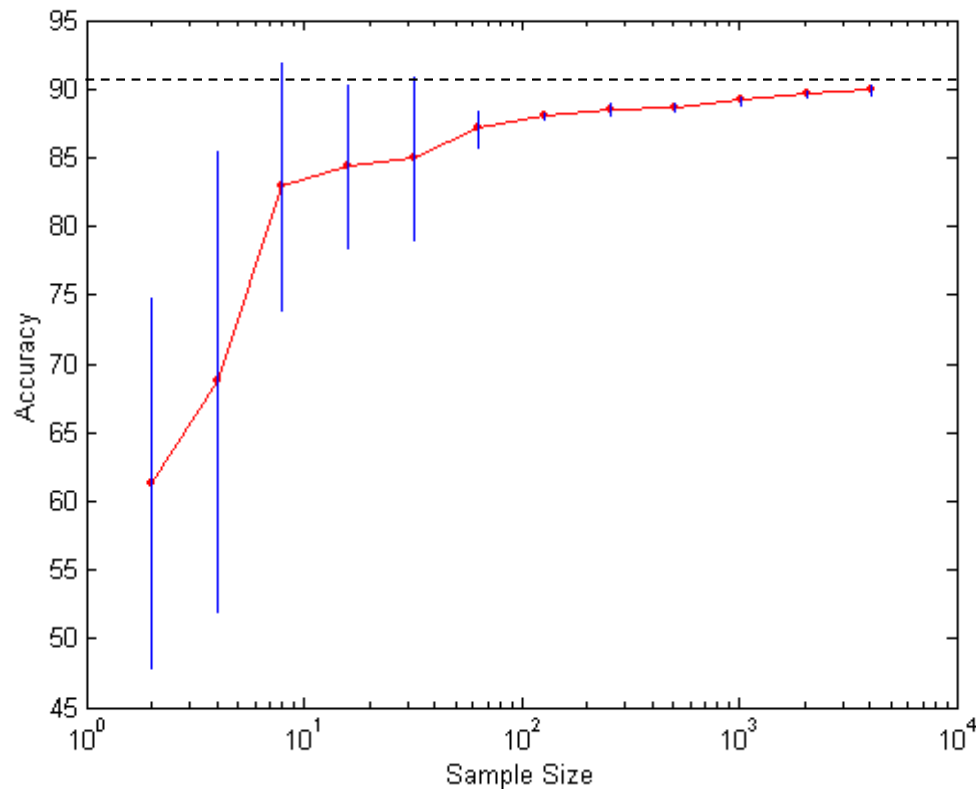| Class | + | - | + | - | - | - | + | - | + | + | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Threshold >= | 0.25 | 0.43 | 0.53 | 0.76 | 0.85 | 0.85 | 0.85 | 0.87 | 0.93 | 0.95 | 1.00 |
| TP | 5 | 4 | 4 | 3 | 3 | 3 | 3 | 2 | 2 | 1 | 0 |
| FP | 5 | 5 | 4 | 4 | 3 | 2 | 1 | 1 | 0 | 0 | 0 |
| TN | 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:**

# 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

# Methods of Estimation

- Holdout
  - Reserve 2/3 for training and 1/3 for testing
- Random subsampling
  - Repeated holdout
- Cross validation
- Stratified sampling
  - Keeps relative frequency of different labels intact
- Bootstrap
  - Sampling with replacement

# Learning Curve



- ☐ Learning curve shows how accuracy changes with varying sample size
- ☐ Requires a sampling schedule for creating learning curve:
  - ☐ Arithmetic sampling (Langley, et al)
  - ☐ Geometric sampling (Provost et al)

Effect of small sample size:
- ‐ Bias in the estimate
- ‐ Variance of estimate