

Data Mining

Classification: Basic Concepts and Techniques

Lecture Notes for Chapter 3

Introduction to Data Mining, 2nd Edition

by

Tan, Steinbach, Karpatne, Kumar

Classification Model Evaluation

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

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

| Most widely-used metric:

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

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 |
| ACTUAL CLASS | 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

Computing Cost 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 |

Accuracy = 80%

Cost = 3910

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

Accuracy = 90%

Cost = 4255

Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{TP}{TP + FP}$$

$$\text{Recall (r)} = \frac{TP}{TP + FN}$$

$$\text{F-measure (F)} = \frac{2rp}{r + p} = \frac{2TP}{2TP + FN + FP}$$

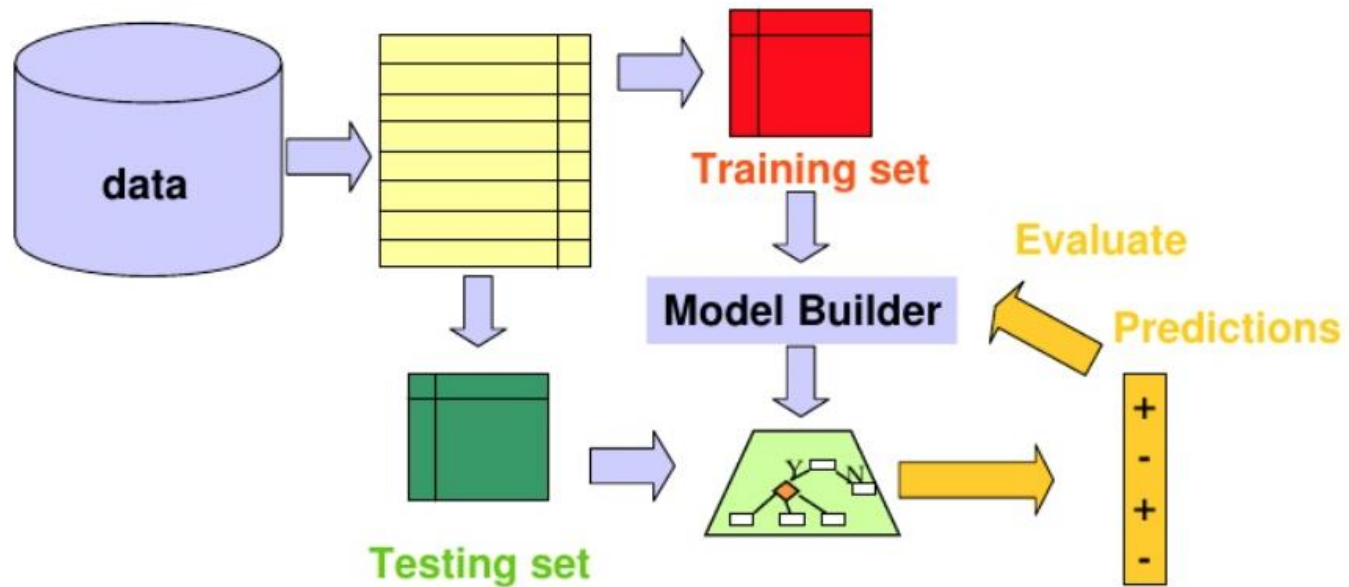
- Precision is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{Yes}|\text{No})$
- Recall is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{No}|\text{Yes})$
- F-measure is biased towards all except $C(\text{No}|\text{No})$

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

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?

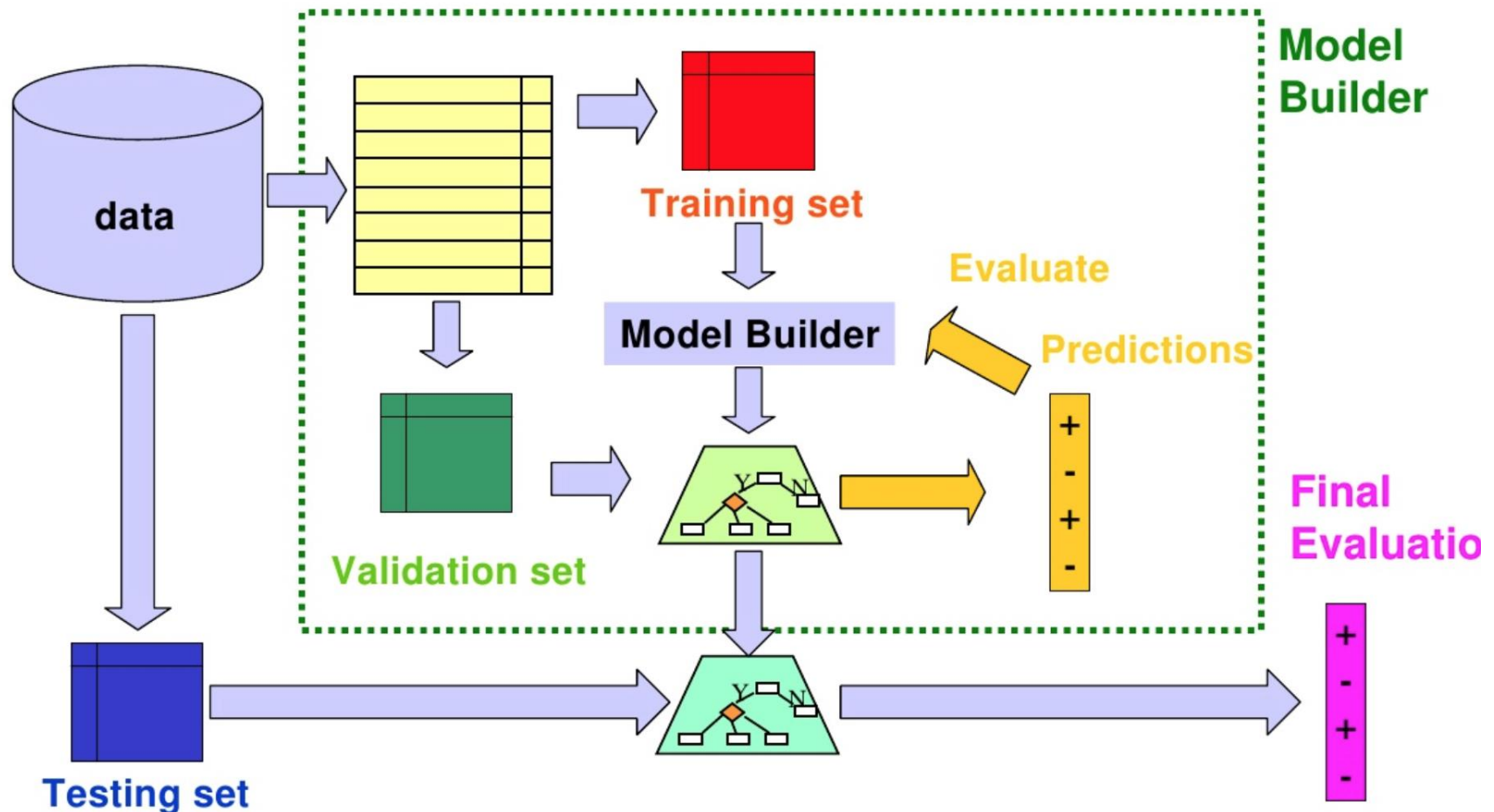
Methods for evaluation



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**: builds the basic structure
 - **Stage 2**: optimizes parameter settings
 - **The test data can't be used for parameter tuning!**
 - Proper procedure uses three sets:
 - ◆ training data,
 - ◆ validation data,
 - ◆ test data
 - **Validation data is used to optimize parameters**
- 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
- The larger the test data the more accurate the error estimate

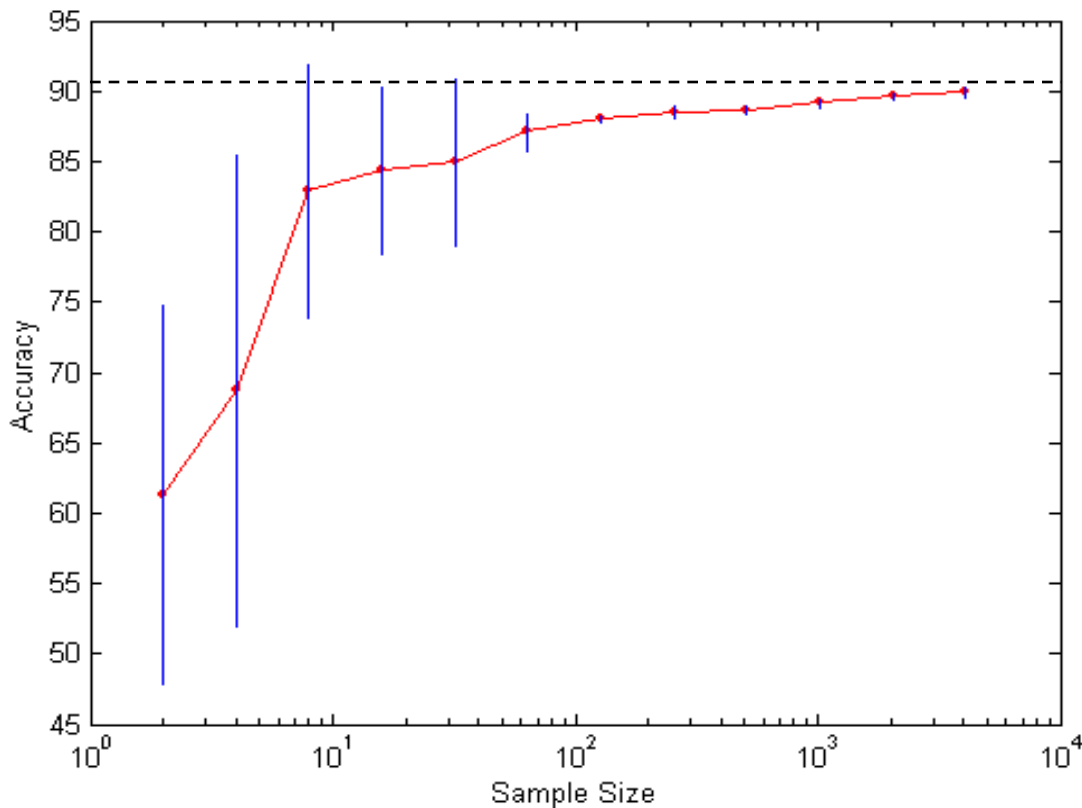
Evaluation: training, validation, test



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
- Requires a **sampling** schedule for creating learning curve

Effect of small sample size:

- Bias in the estimate
- Variance of estimate

1. How much a classification model benefits from adding more training data?
2. Does the model suffer from a variance error or a bias error?

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

Small & Unbalanced Data

- The holdout method reserves a certain amount for **testing** and uses the remainder for **training**
- Usually, **one third for testing**, the rest for training
- For small or “unbalanced” datasets, **samples might not be representative**
 - For instance, few or none instances of some classes
- Stratified sample
 - **Balancing the data**
 - Make sure that each class is represented with approximately equal proportions in both subsets

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

Cross-validation

- Avoids overlapping test sets
 - **First step:** data is split into k subsets of equal size
 - **Second step:** each subset in turn is used for testing and the remainder for training
- This is called **k-fold cross-validation**
- Often the subsets are stratified before cross-validation is performed
- The **error estimates** are **averaged** to yield an overall error estimate
- **Even better:** repeated stratified cross-validation
E.g. ten-fold cross-validation is repeated ten times and results are averaged (reduces the variance)



Model Evaluation

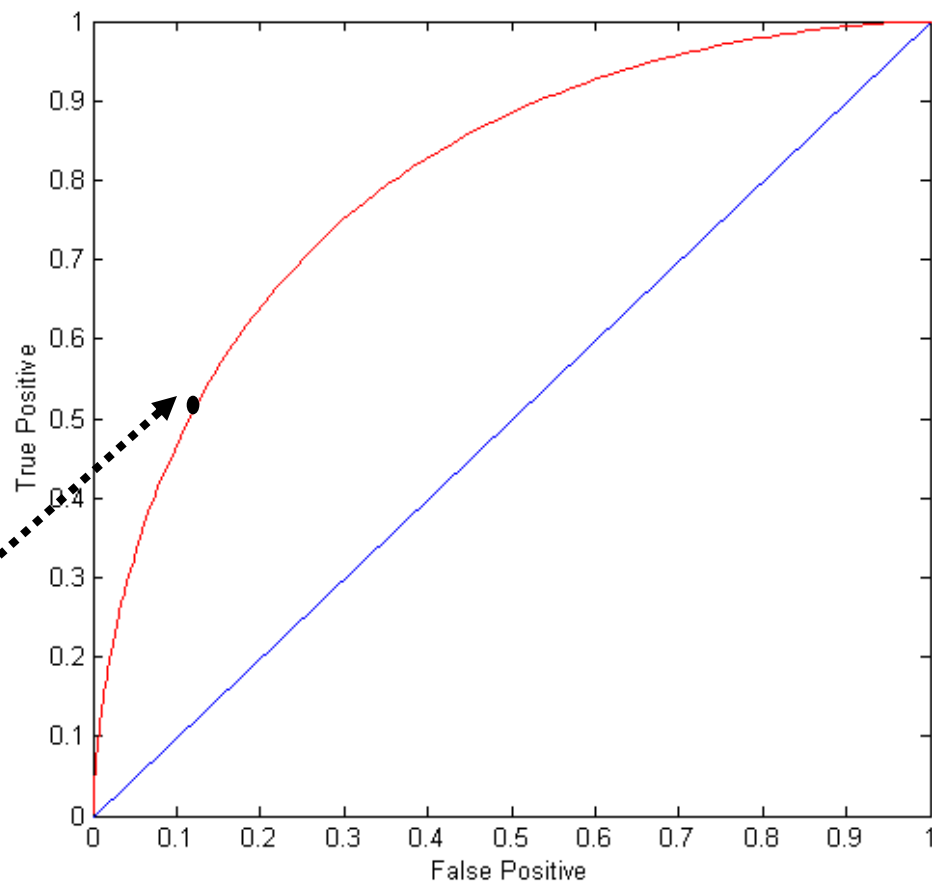
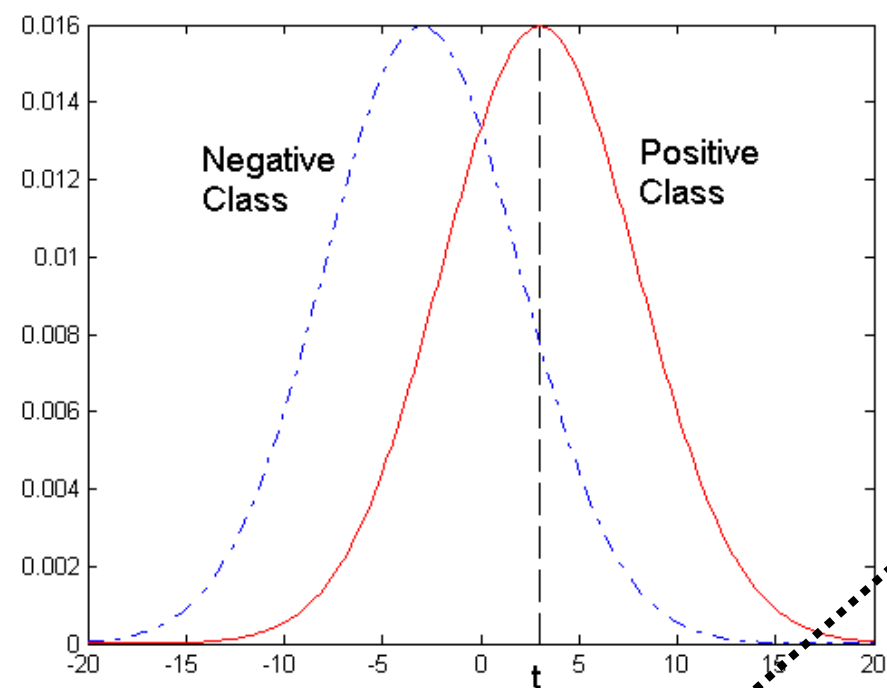
- | Metrics for Performance Evaluation
 - How to evaluate the performance of a model?
- | Methods for Performance Evaluation
 - How to obtain reliable estimates?
- | **Methods for Model Comparison**
 - How to compare the relative performance among competing models?

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



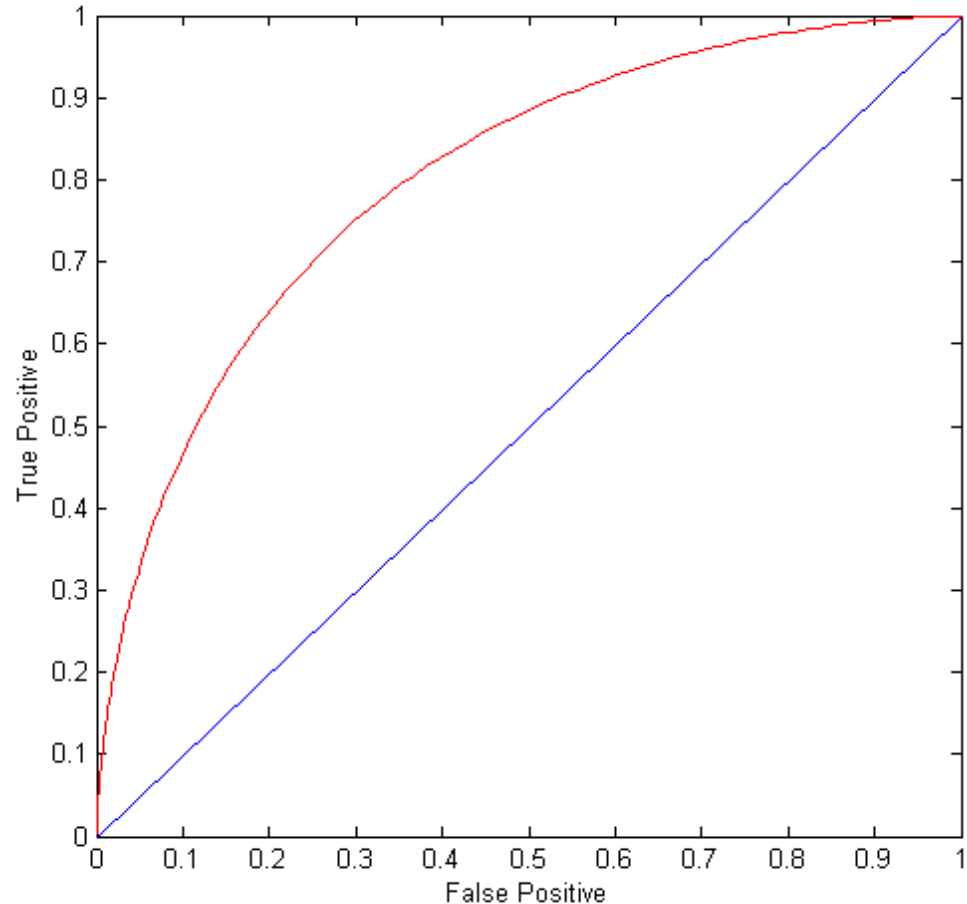
At threshold t :

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

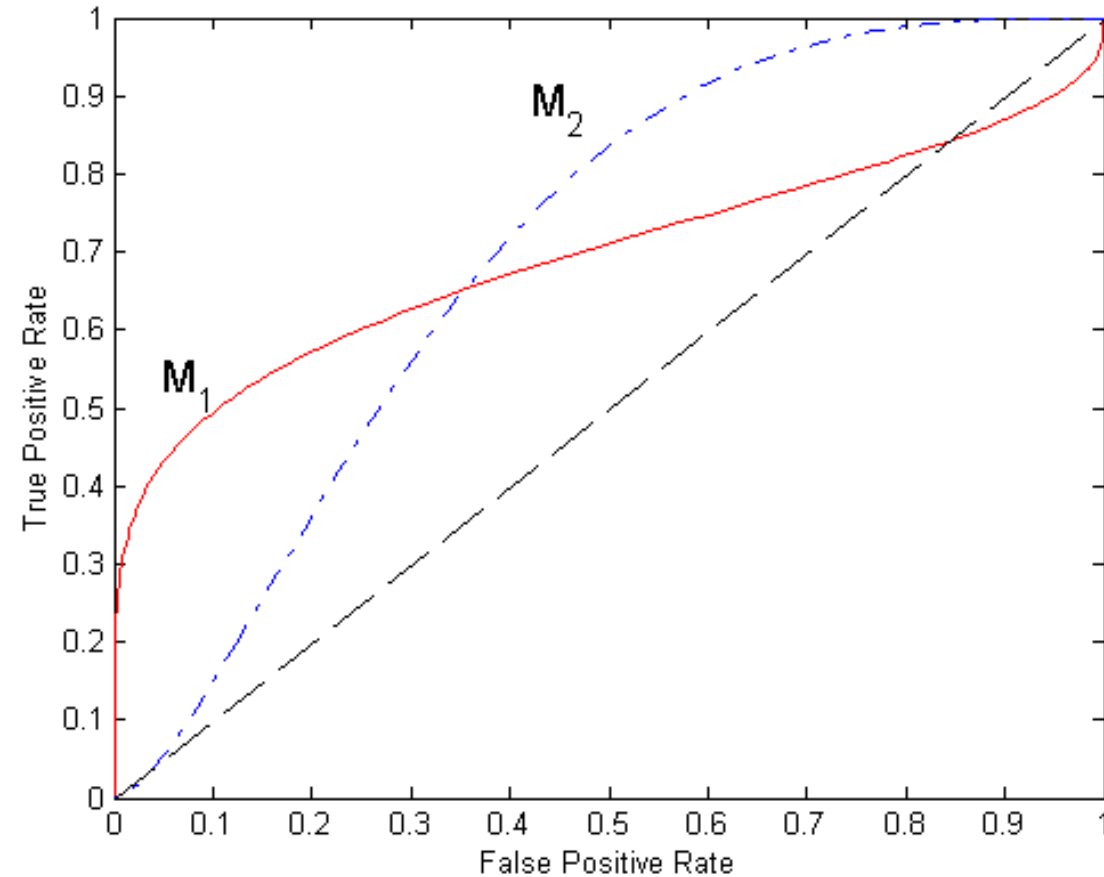
ROC Curve

(TP,FP):

- | (0,0): declare everything to be negative class
- | (1,1): declare everything to be positive class
- | (0,1): 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

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 | + |

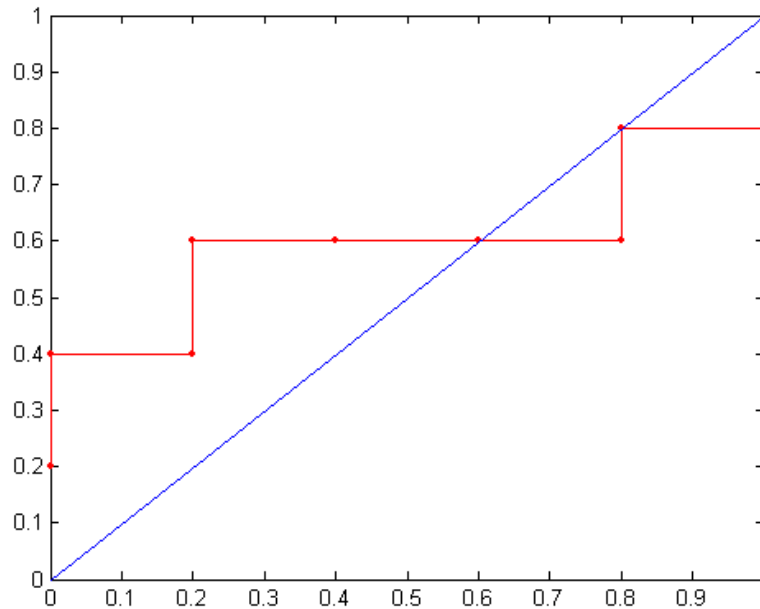
- Use classifier that produces posterior probability for each test instance $P(+|A)$
- Sort the instances according to $P(+|A)$ in decreasing order
- Apply threshold at each unique value of $P(+|A)$
- Count the number of TP, FP, TN, FN at each threshold
- TP rate, $TPR = TP/(TP+FN)$
- FP rate, $FPR = FP/(FP + TN)$

How to construct an ROC curve

| Class | + | - | + | - | - | - | + | - | + | + | |
|------------------|------|------|------|------|------|------|------|------|------|------|------|
| Threshold \geq | 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 |

| Inst. | 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 | + |

ROC Curve:



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?

Confidence Interval for Accuracy

- | Prediction can be regarded as a Bernoulli trial (binomial random experiment)
 - A Bernoulli trial has 2 possible outcomes
 - Possible outcomes for prediction: correct or wrong
 - Probability of success is constant
 - Collection of Bernoulli trials has a Binomial distribution:
 - ◆ $x \sim \text{Bin}(N, p)$ x : # of correct predictions, N trials, p constant prob.
 - ◆ e.g: Toss a fair coin 50 times, how many heads would turn up?
Expected number of heads = $N \times p = 50 \times 0.5 = 25$

Given x (# of correct predictions) or equivalently, $\text{acc} = x/N$, and N (# of test instances)

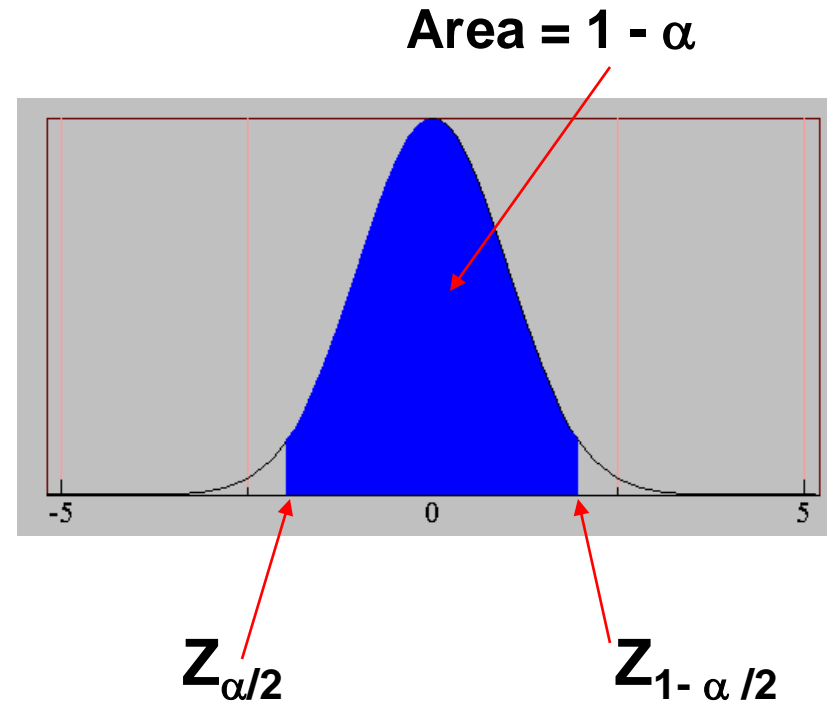
Can we predict p (true accuracy of model)?

Confidence Interval for Accuracy

For large test sets ($N > 30$),

- acc has a normal distribution with mean p and variance $p(1-p)/N$
- the confidence interval for acc can be derived as follows:

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



Confidence Interval for p :

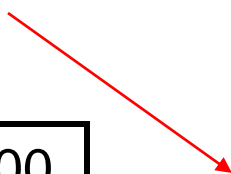
$$p = \frac{2 \times N \times acc + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4 \times N \times acc - 4 \times N \times acc^2}}{2(N + Z_{\alpha/2}^2)}$$

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)
 - Which is the confidence interval?**
 - From probability table, $Z_{\alpha/2}=1.96$

| 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 |

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



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 variance of error rates: $\hat{\sigma}_i = \frac{e_i(1-e_i)}{n_i}$

Comparing Performance of 2 Models

- To test if performance difference is statistically significant: $d = e_1 - e_2$
 - $d \sim N(d_t, \sigma_t)$ where d_t is the true difference
 - Since D_1 and D_2 are independent, their variance adds up:

$$\begin{aligned} S_t^2 &= S_1^2 + S_2^2 @ \hat{S}_1^2 + \hat{S}_2^2 \\ &= \frac{e_1(1 - e_1)}{n_1} + \frac{e_2(1 - e_2)}{n_2} \end{aligned}$$

- It can be shown at $(1-\alpha)$ confidence level,

$$d_t = d \pm Z_{\alpha/2} \hat{\sigma}_t$$

An Illustrative Example

- Given: M1: $n_1 = 30$, $e_1 = 0.15$
M2: $n_2 = 5000$, $e_2 = 0.25$

- $d = |e_2 - e_1| = 0.1$ (2-sided test to check: $d_t = 0$ or $d_t \neq 0$)

$$\hat{S}_d^2 = \frac{0.15(1 - 0.15)}{30} + \frac{0.25(1 - 0.25)}{5000} = 0.0043$$

- At 95% confidence level, $Z_{\alpha/2} = 1.96$

$$d_t = 0.100 \pm 1.96 \times \sqrt{0.0043} = 0.100 \pm 0.128$$

=> Interval contains 0 => difference may not be statistically significant