# Data Mining Classification: Basic Concepts and Techniques

Lecture Notes for Chapter 3

Introduction to Data Mining, 2<sup>nd</sup> Edition
by

Tan Stainbach Karnatna Kumar

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 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	а	b				
	Class=No	С	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=Yes	a (TP)	b (FN)				
CLASS	Class=No	c (FP)	d (TN)				

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						
	C(i j)	Class=Yes	Class=No				
ACTUAL CLASS	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

# **Computing Cost of Classification**

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

Model M <sub>1</sub>	PREDICTED CLASS					
		+	-			
ACTUAL CLASS	+	150	40			
	-	60	250			

Model M <sub>2</sub>	PRED	ICTED (	CLASS
		+	•
ACTUAL CLASS	+	250	45
	-	5	200

Accuracy = 80%

Cost = 3910

Accuracy = 90%

Cost = 4255

## **Cost-Sensitive Measures**

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

- 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)

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?

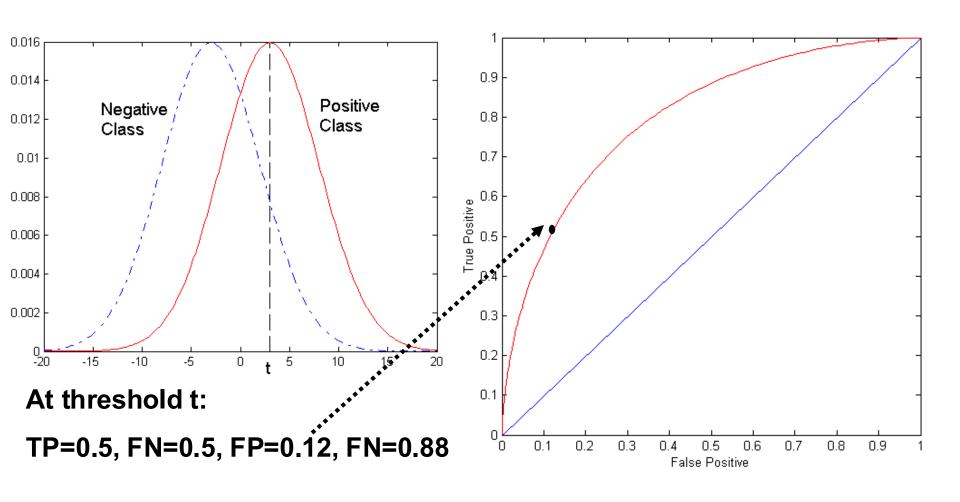
- Methods for Model Comparison
  - How to compare the relative performance among competing models?

# **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**

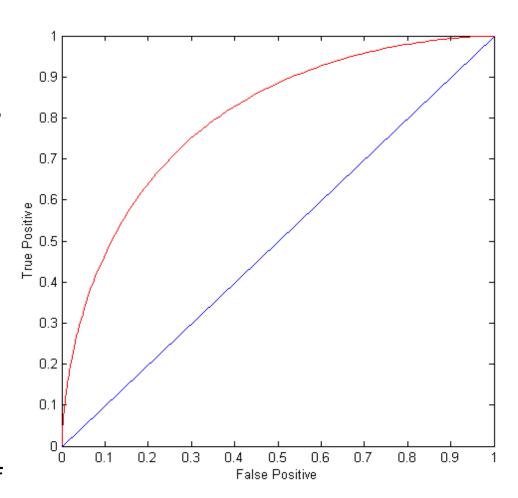
- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at x > t is classified as positive



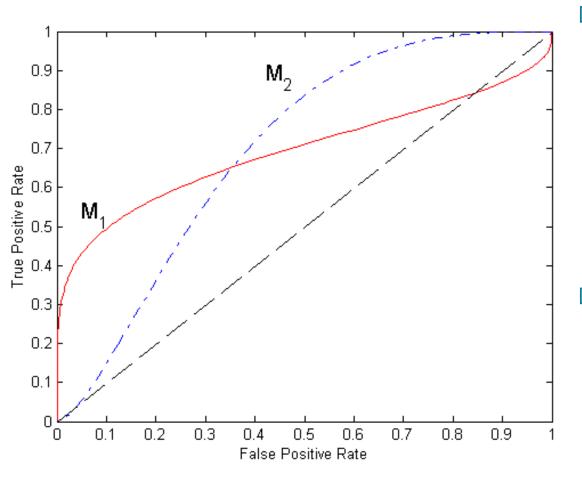
### **ROC Curve**

#### (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<sub>1</sub> is better for small FPR
  - M<sub>2</sub> 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
<b>→</b>	TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0
$\longrightarrow$	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	_

0.9	•							,		+
0.8							,	$\checkmark$		+
0.7										$\frac{1}{2}$
0.6				•—		<del>/</del>		_		$\frac{1}{2}$
0.5										$\frac{1}{2}$
0.4		_	/							$\frac{1}{2}$
0.3										$\frac{1}{2}$
0.2										$\frac{1}{2}$
0.1	. /									$\frac{1}{2}$
ا 0	0.1	0.2	0.3 0	.4	0.5	0.6	0.7	0.8	0.9	_  1

# **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 ~ 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×p = 50 × 0.5 = 25

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

Can we predict p (true accuracy of model)?

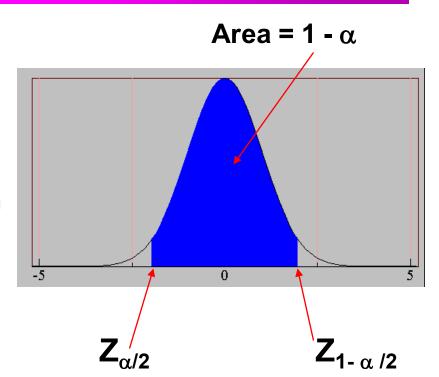
# **Confidence Interval for Accuracy**

## For large 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-α	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=n1), found error rate =  $e_1$
  - M2 is tested on D2 (size=n2), found error rate = e<sub>2</sub>
  - Assume D1 and D2 are independent
  - If n1 and n2 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 D1 and D2 are independent, their variance adds up:

$$\sigma_t^2 = \sigma_1^2 + \sigma_2^2 \cong \hat{\sigma}_1^2 + \hat{\sigma}_2^2$$

$$= \frac{e1(1-e1)}{n1} + \frac{e2(1-e2)}{n2}$$

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

$$d_{t} = d \pm Z_{\alpha/2} \hat{\sigma}_{t}$$

# **An Illustrative Example**

☐ Given: M1: n1 = 30, e1 = 0.15

M2: n2 = 5000, e2 = 0.25

d = |e2 - e1| = 0.1 (2-sided test to check: dt = 0 or dt <> 0)

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

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

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

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