

## 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
	Class=Yes	a                  b
ACTUAL CLASS	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
	Class=Yes	a (TP)                  b (FN)
ACTUAL CLASS	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(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 $M_1$	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 vs Accuracy

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

Accuracy is proportional to cost if

1.  $C(\text{Yes}|\text{No}) = C(\text{No}|\text{Yes}) = q$
2.  $C(\text{Yes}|\text{Yes}) = C(\text{No}|\text{No}) = p$

$$N = a + b + c + d$$

$$\text{Accuracy} = (a + d)/N$$

Cost	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	p	q
	Class=No	q	p

$$\text{Cost} = p(a + d) + q(b + c)$$

$$= p(a + d) + q(N - a - d)$$

$$= qN - (q - p)(a + d)$$

$$= N[q - (q - p) \times \text{Accuracy}]$$



## Cost-Sensitive Measures

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

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

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

- 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?

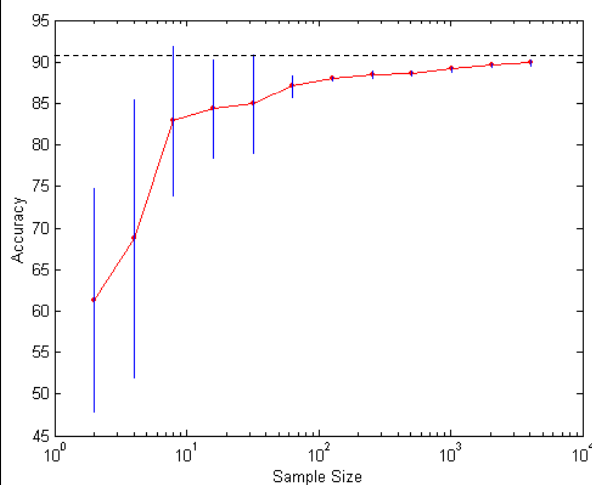


## Methods for Performance Evaluation

- ▶ 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:
  - Arithmetic sampling (Langley, et al)
  - Geometric sampling (Provost et al)

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?

▶

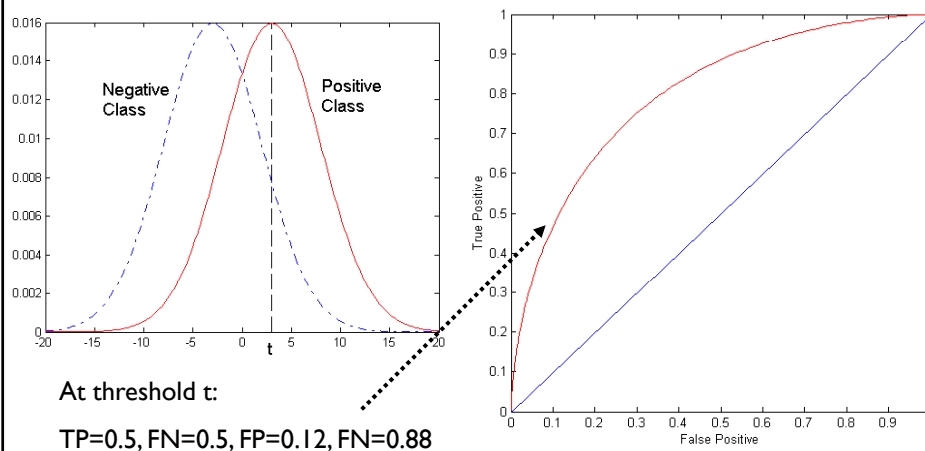
## ROC (Receiver Operating Characteristic)

- ▶ 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

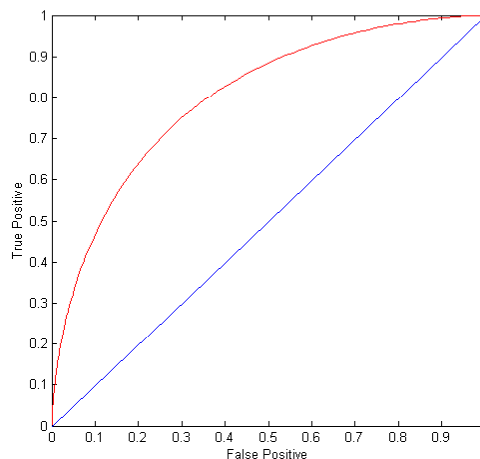




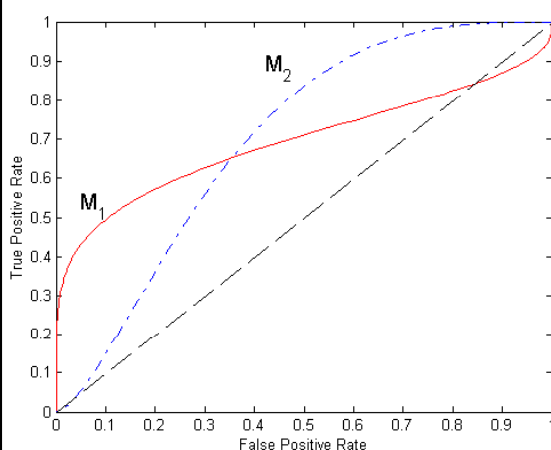
## ROC Curve

(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

## 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

