

# 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 (TP)	b (FN)
	Class=No	c (FP)	d (TN)

TP: True Positive  
FP: False Positive  
TN: True Negative  
FN: False Negative

- 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

**Unbalanced classes**

- 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

# Other Measures

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

$$\text{False positive rate} = \frac{c}{c + d}$$

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

$$\text{True positive rate} = \frac{a}{a + b} = \text{Sensitivity}$$

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

$$\text{True Negative Rate} = \frac{d}{c + d} = \text{Specificity}$$

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

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

# Model Evaluation

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

# 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

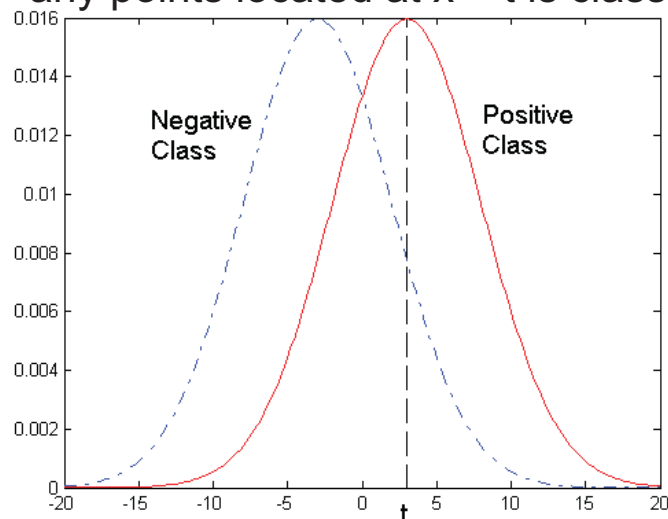
- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?
- Methods for Performance Evaluation
  - How to obtain reliable estimates?
- **Methods for Model Comparison**
  - How to **compare the relative performance** among competing **models**?

# 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 rate (y-axis)** against **FP rate (x-axis)**
- Performance of **each classifier** represented as **a point** on ROC curve  
**changing the threshold** of algorithm, or **sample** distribution changes the **location of the point**

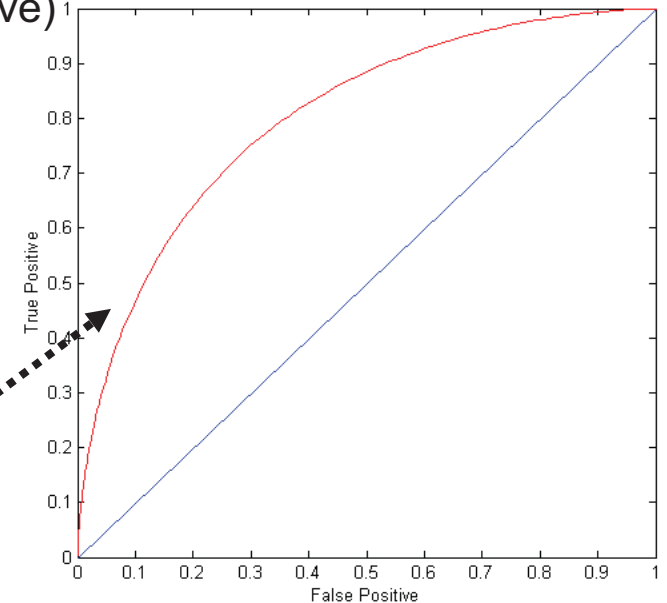
1-dimensional data set containing 2 classes (positive and negative)

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



At threshold  $t$ :

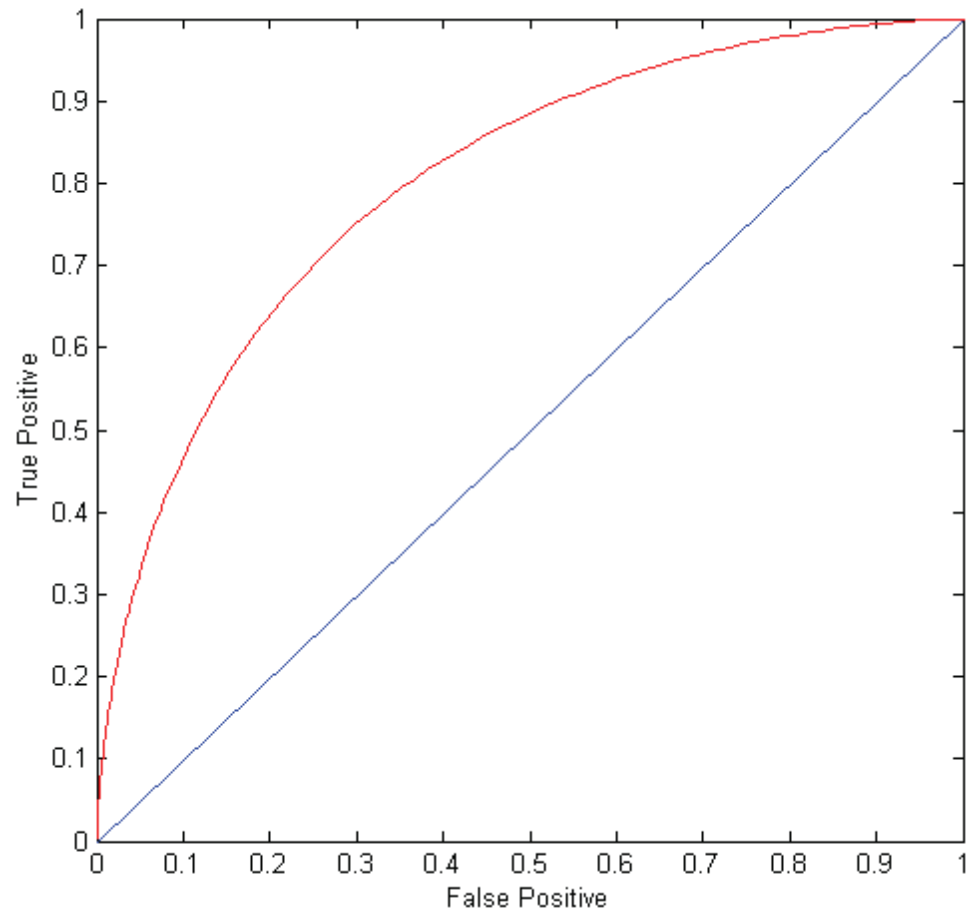
TPR=0.5, FPR=0.12



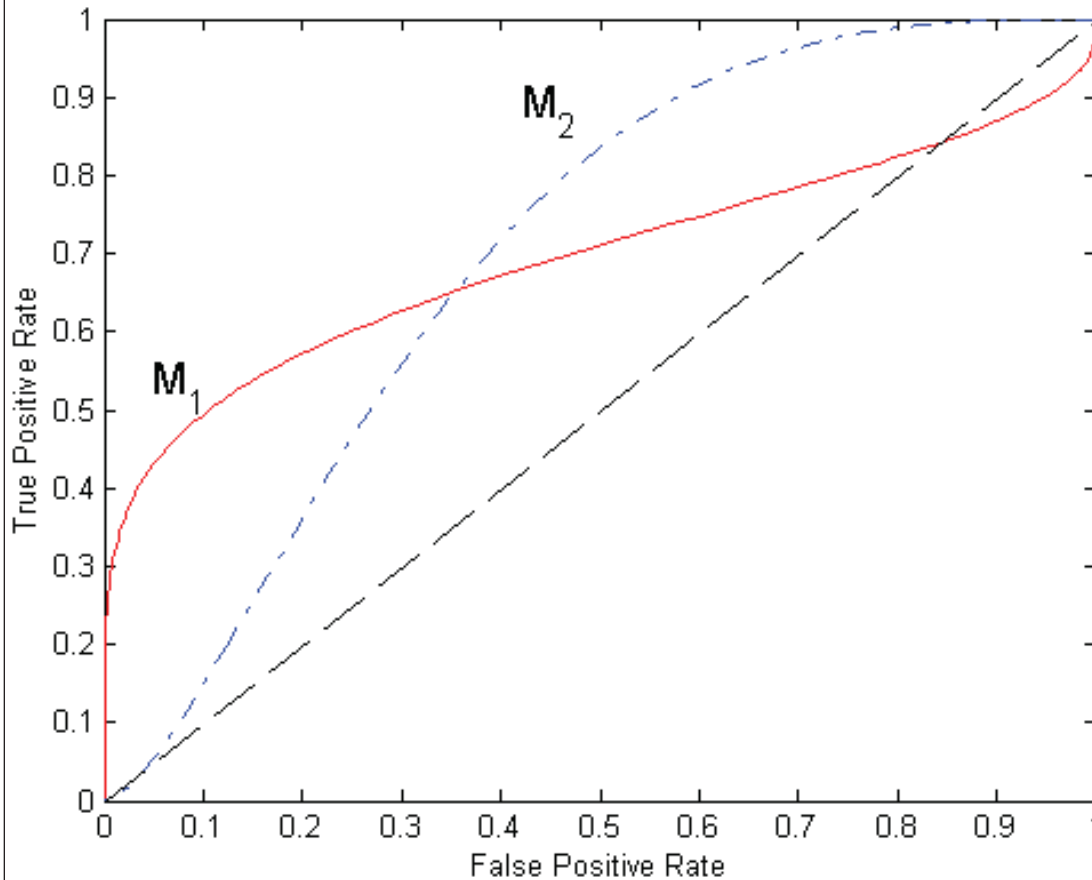
# ROC Curve

(TPR,FPR):

- (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

Posterior probability of test instance  $x$

Threshold:  $t$

# of  $+$   $\geq t$

# of  $-$   $\geq t$

P(+ x)	0.95	0.93	0.87	0.85	0.85	0.85	0.76	0.53	0.43	0.25	
Class	+	-	+	-	-	-	+	-	+	+	
Threshold: $t$	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:

