

Evaluation of Classification Models

Limitation of Accuracy

- Consider a 2-class problem:
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10
- If the model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
 - Accuracy is misleading!

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

Most widely-used metric:

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

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 $C(\text{No}|\text{No})$

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

Perils of Overfitting



Data Science Dojo

@DataScienceDojo

Perils of **#overfitting** @kaggle restaurant revenue prediction Pos 1 drops to 2041 in final ranking.



2041	↑7	Cheng Jiang
2042	↓2041	BAYZ, M.D. 
2043	↓81	Alberto



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?

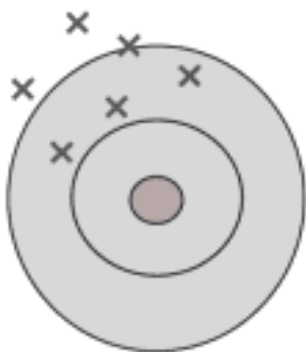
Methods for Performance Evaluation

- 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

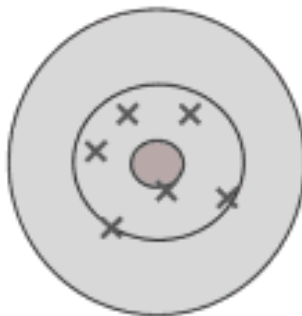
Methods of Estimation

- Holdout
 - Reserve 2/3 for training and 1/3 for testing
- Cross validation
 - Partition data into k disjoint subsets
 - k -fold: train on $k-1$ partitions, test on the remaining one
 - Leave-one-out: $k = n$
- Random subsampling
 - Repeated holdout
- Stratified sampling
 - Oversampling vs undersampling
- Bootstrap
 - Sampling with replacement

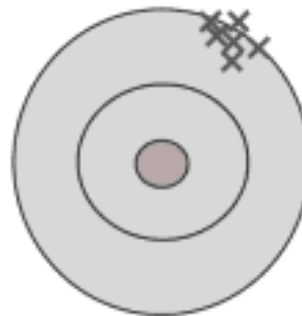
Understanding Bias And Variance In Estimate



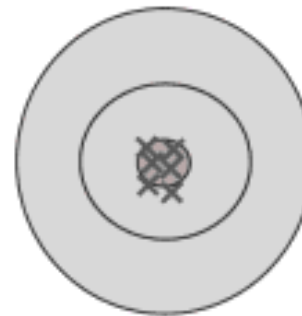
High bias
High variance



Low bias
High variance



High bias
Low variance



Low bias
Low 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

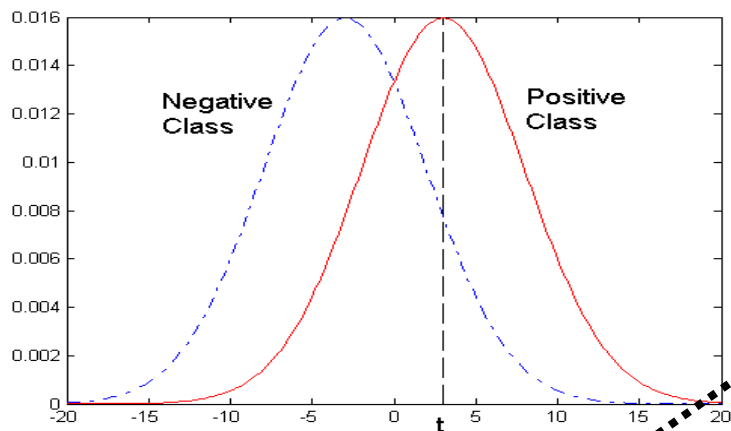
How to compare the relative performance among competing models?

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 the algorithm, sample distribution, or cost matrix changes the location of the point

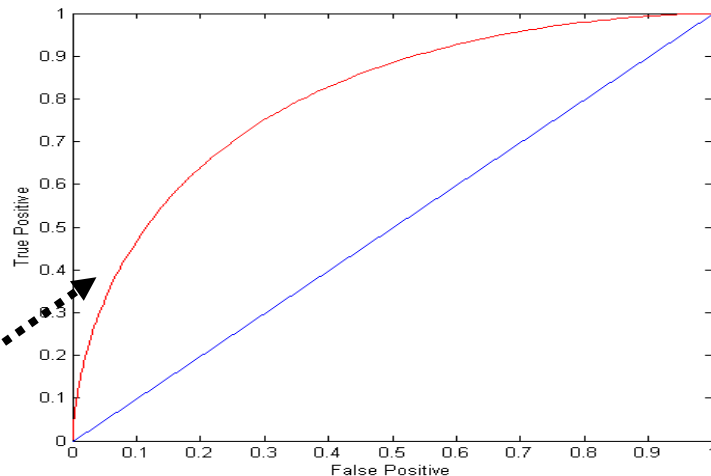
ROC Curve

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



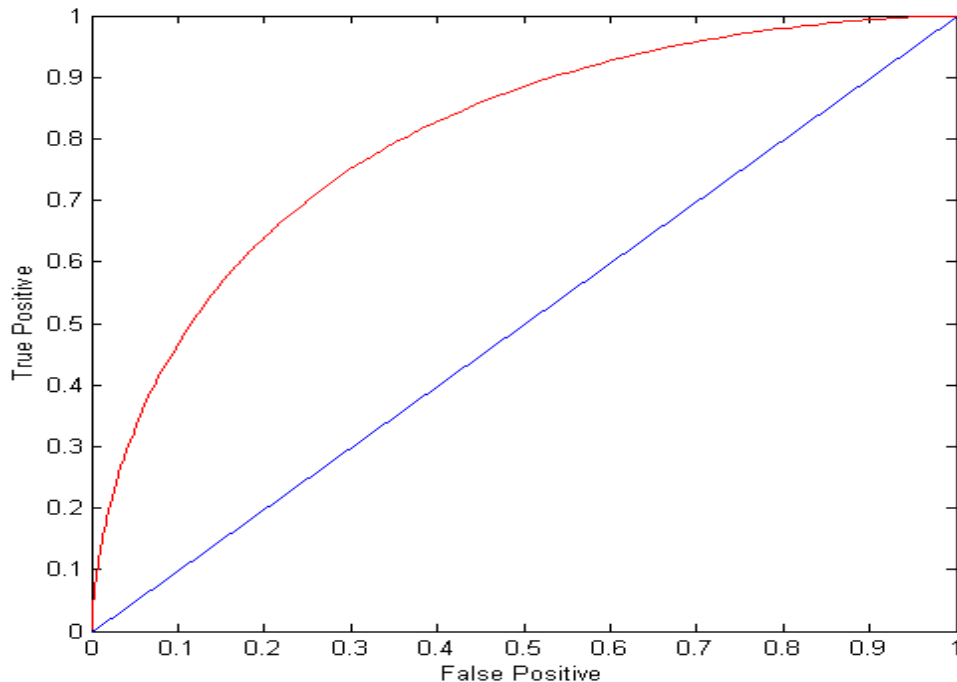
At threshold t:

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

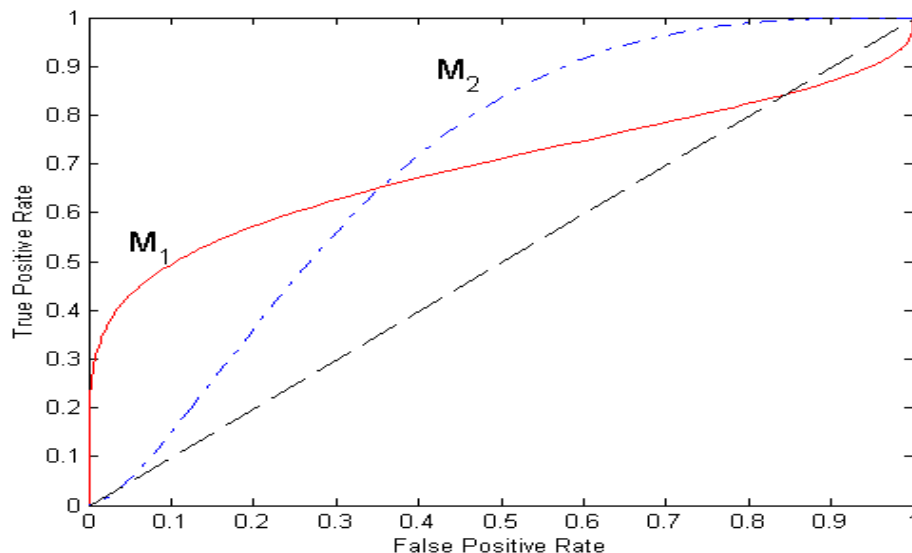


ROC Curve

- (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



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

QUESTIONS