# Classifier Performance and Model Selection

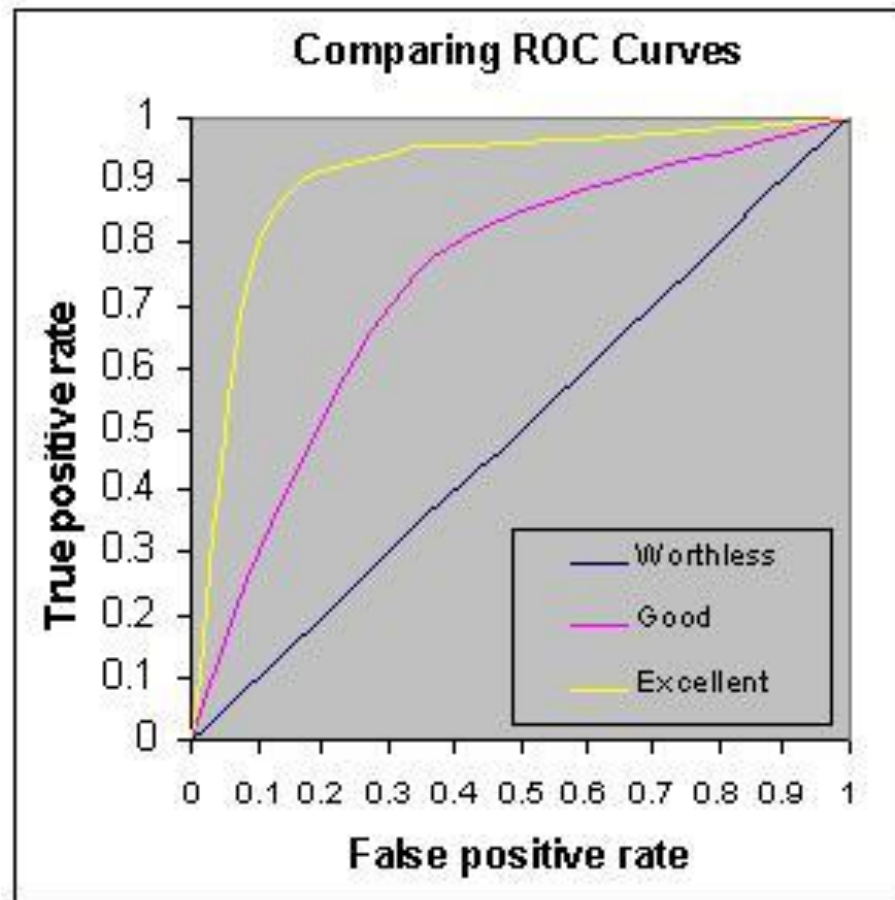# Overview

- Evaluating the performance of a classifier
  - Confusion Matrix
  - **ROC** curve
- Selecting a model for the problem
  - Bias / Variance tradeoff
  - Cross Validation

# Receiver Operating Characteristic Methodology

# Introduction to ROC curves

- *ROC* = *Receiver Operating Characteristic*

- Started in electronic signal detection theory (1940s - 1950s)

- Has become very popular in biomedical applications, particularly radiology and imaging

- Also used in **machine learning applications** to assess classifiers

- Can be used to compare tests/procedures

# ROC curves: simplest case

- Consider diagnostic **test** for a disease

- Test has 2 possible outcomes:

  - '**positive**' = suggesting presence of disease

  - '**negative**'

- An individual can test either positive or negative for the disease

# Hypothesis testing refresher

- 2 'competing theories' regarding a population parameter:
  - *NULL* hypothesis *H* ('straw man')
  - ALTERNATIVE hypothesis *A* ('claim', or theory you wish to test)
- *H:* NO DIFFERENCE
  - any observed deviation from what we expect to see is due to *chance variability*
- *A:* THE DIFFERENCE IS *REAL*

# Test statistic

- Measure how far the observed data are from what is expected *assuming the NULL  H*  by computing the value of a *test statistic*  (TS) from the data

- The particular TS computed depends on the parameter

- For example, to test the population mean $\mu$, the TS is the sample mean (or standardized sample mean)

- The NULL is rejected if the TS falls in a user-specified 'rejection region'

# Types of errors

True class

|  | p | n |
|---|---|---|
| **Y** | True Positives | False Positives |
| **N** | False Negatives | True Negatives |

Hypothesized class

Column totals: **P**     **N**

$$\text{FP rate} = \frac{FP}{N} \qquad \text{TP rate} = \frac{TP}{P} = \text{Recall}$$

$$\text{Precision} = \frac{TP}{TP + FP} \qquad \text{F-score} = \text{Precision} \times \text{Recall}$$

$$\text{Accuracy} = \frac{TP + TN}{P + N}$$

Two parts to each: whether you got it correct or not, and what you guessed.

True Positive

Did we get it correct?
True, we did get it correct.

What did we say?
We said 'positive'

False Negative

Did we get it correct?
False, we did not get it correct.

What did we say?
We said 'negative

# Sensitivity and Specificity

Count up the total number of each label (TP, FP, TN, FN) over a large dataset. In ROC analysis, we use two statistics:

**Sensitivity** $= \dfrac{TP}{TP+FN}$    Can be thought of as the likelihood of spotting a positive case when presented with one.
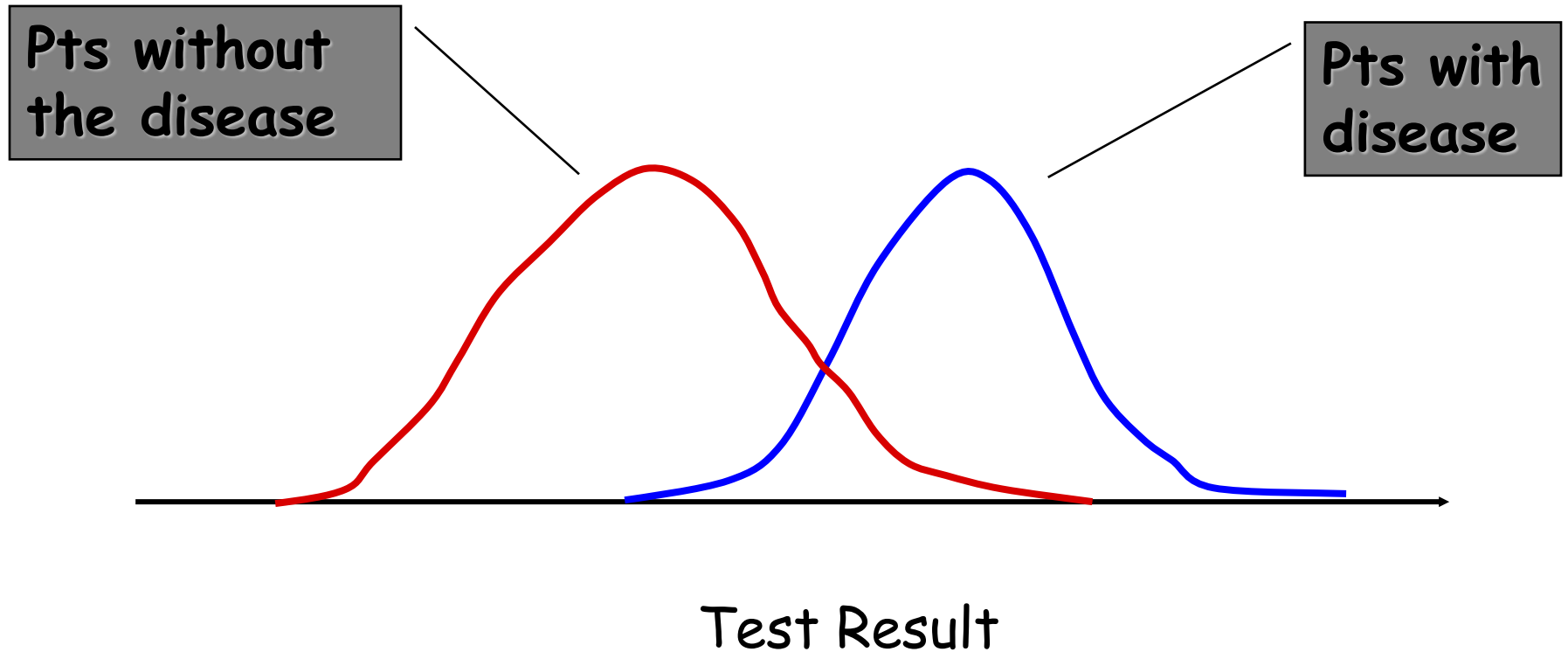
**Specificity** $= \dfrac{TN}{TN+FP}$    Can be thought of as the likelihood of spotting a negative case when presented with one.
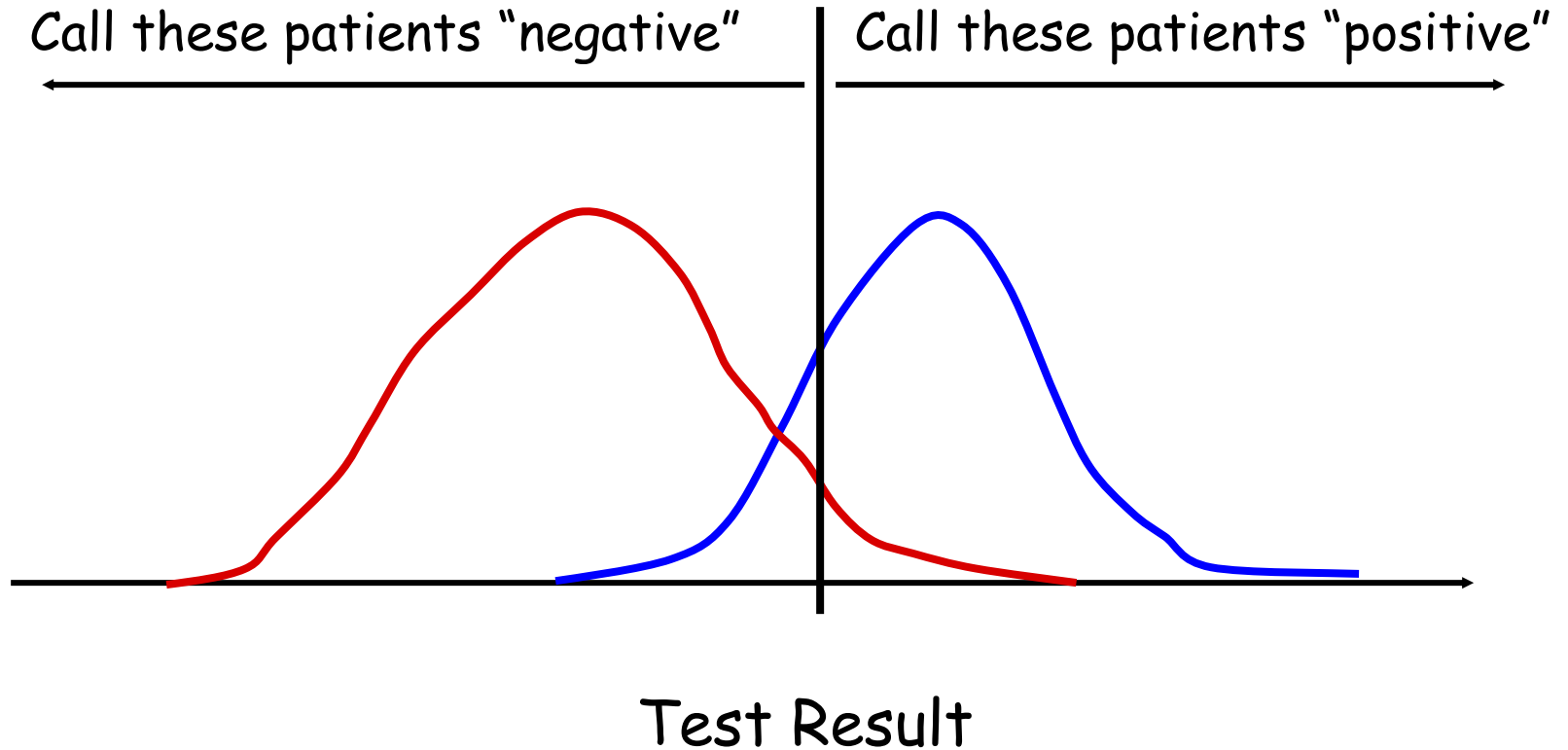
**Ground Truth**

|  | **1** | **0** |
|---|---|---|
| **1** | TRUE POS<br>60 | FALSE POS<br>80 |
| **0** | FALSE NEG<br>30 | TRUE NEG<br>20 |

**Prediction**

60+30 = 90 cases in the dataset were class 1

80+20 = 100 cases in the dataset were class 0

90+100 = 190 examples in the data overall

# Specific Example



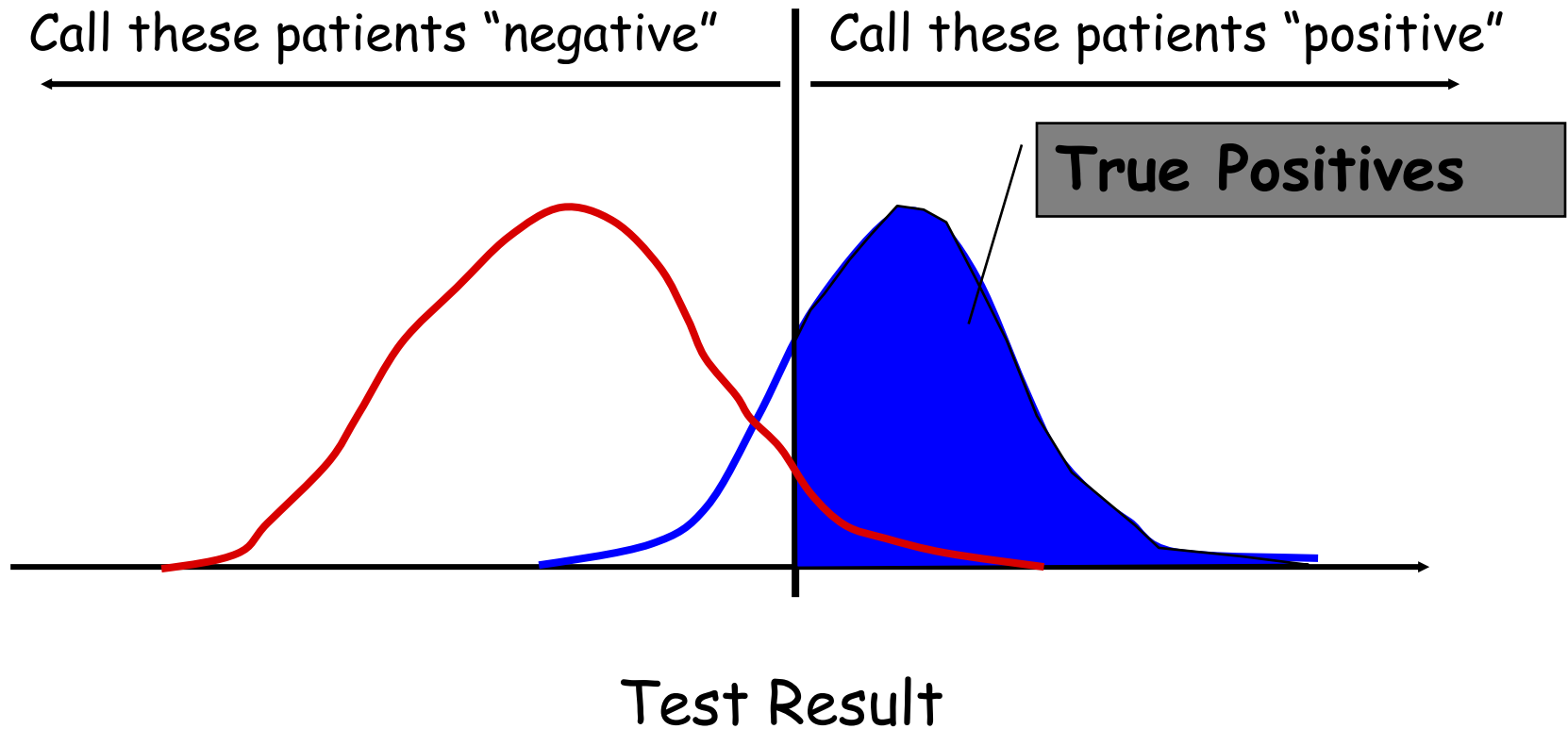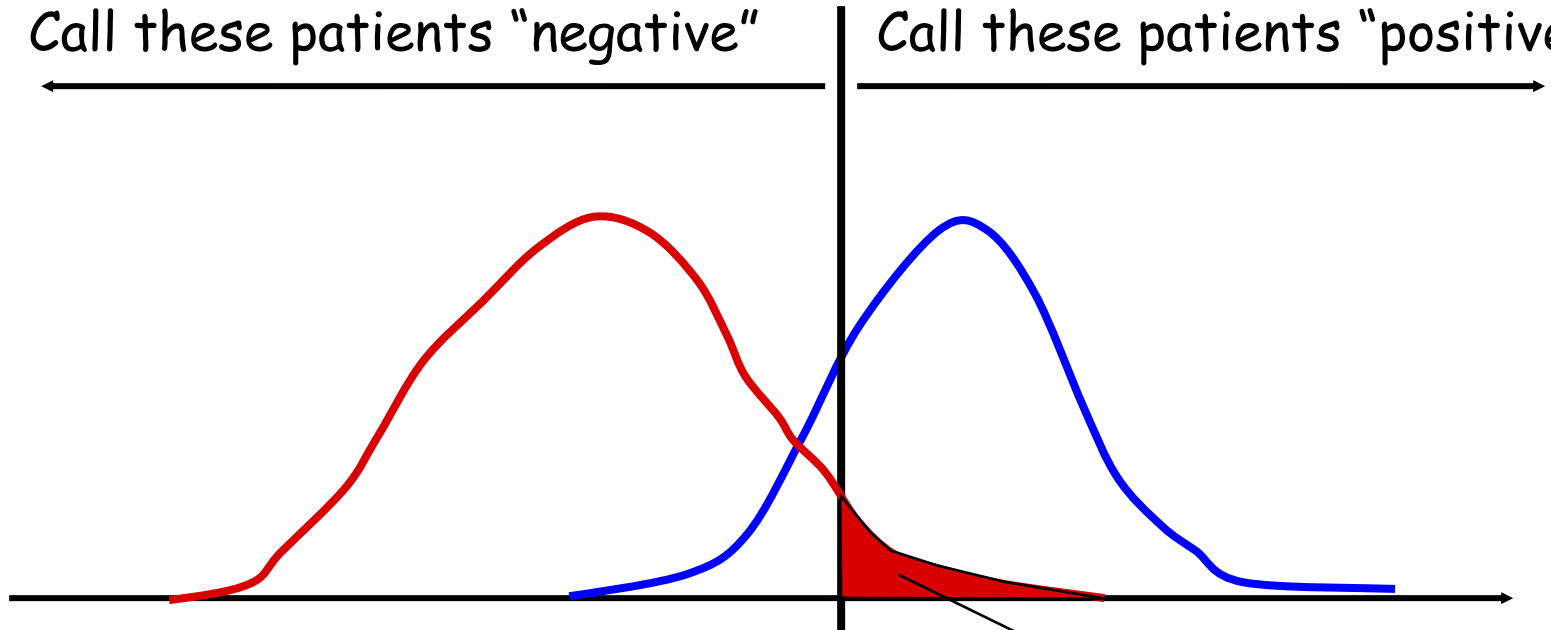Pts without the disease

Pts with disease

Test Result

# Threshold



Call these patients "negative" ← → Call these patients "positive"

Test Result

# Some definitions ...
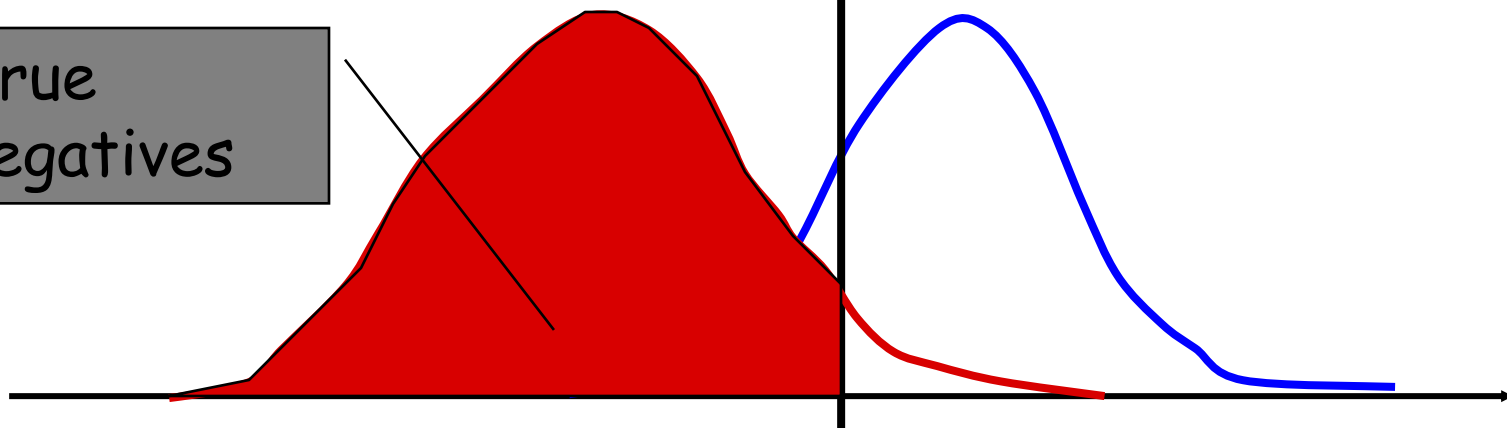


Call these patients "negative"    Call these patients "positive"

True Positives

Test Result

<span style="color:red">~~without the disease~~</span>
<span style="color:blue">**with the disease**</span>

Call these patients "negative"   Call these patients "positive"

Test Result

without the disease
with the disease

False Positives

Call these patients "negative"                    Call these patients "positive"

True
negatives

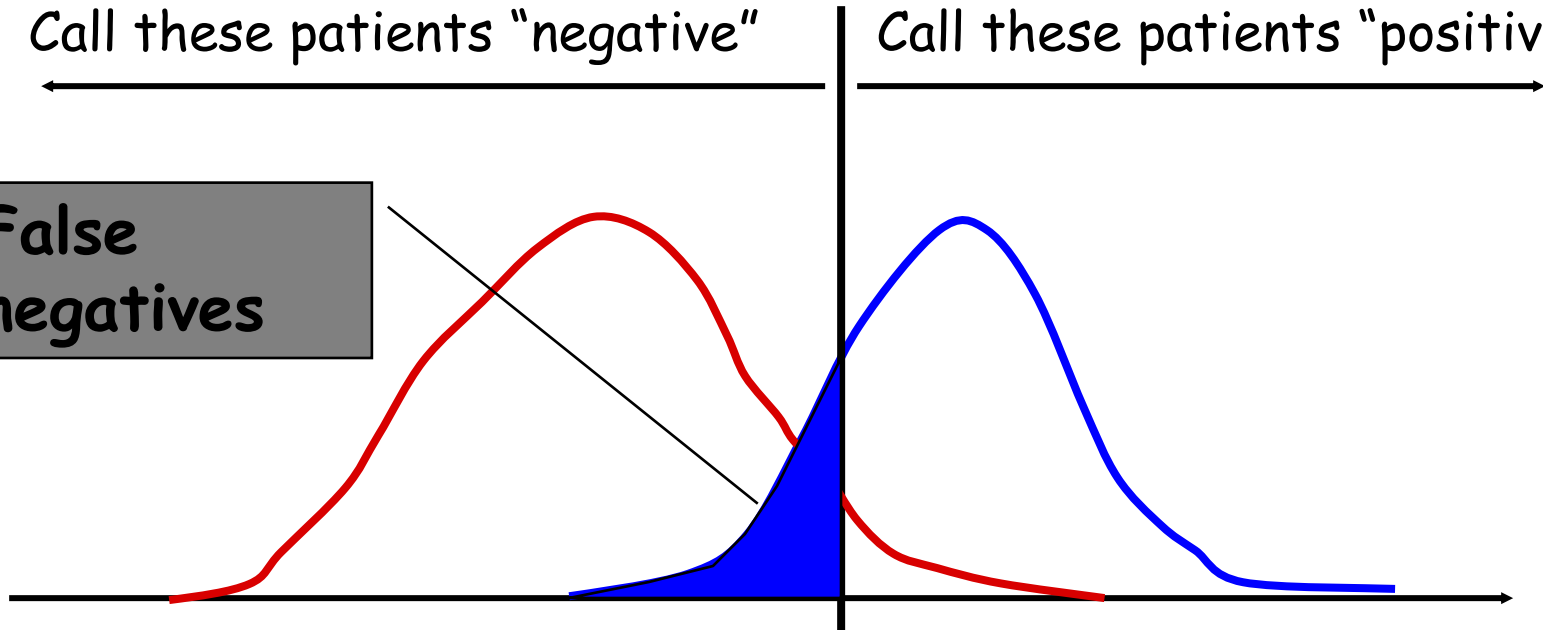Test Result

**without the disease**
**with the disease**

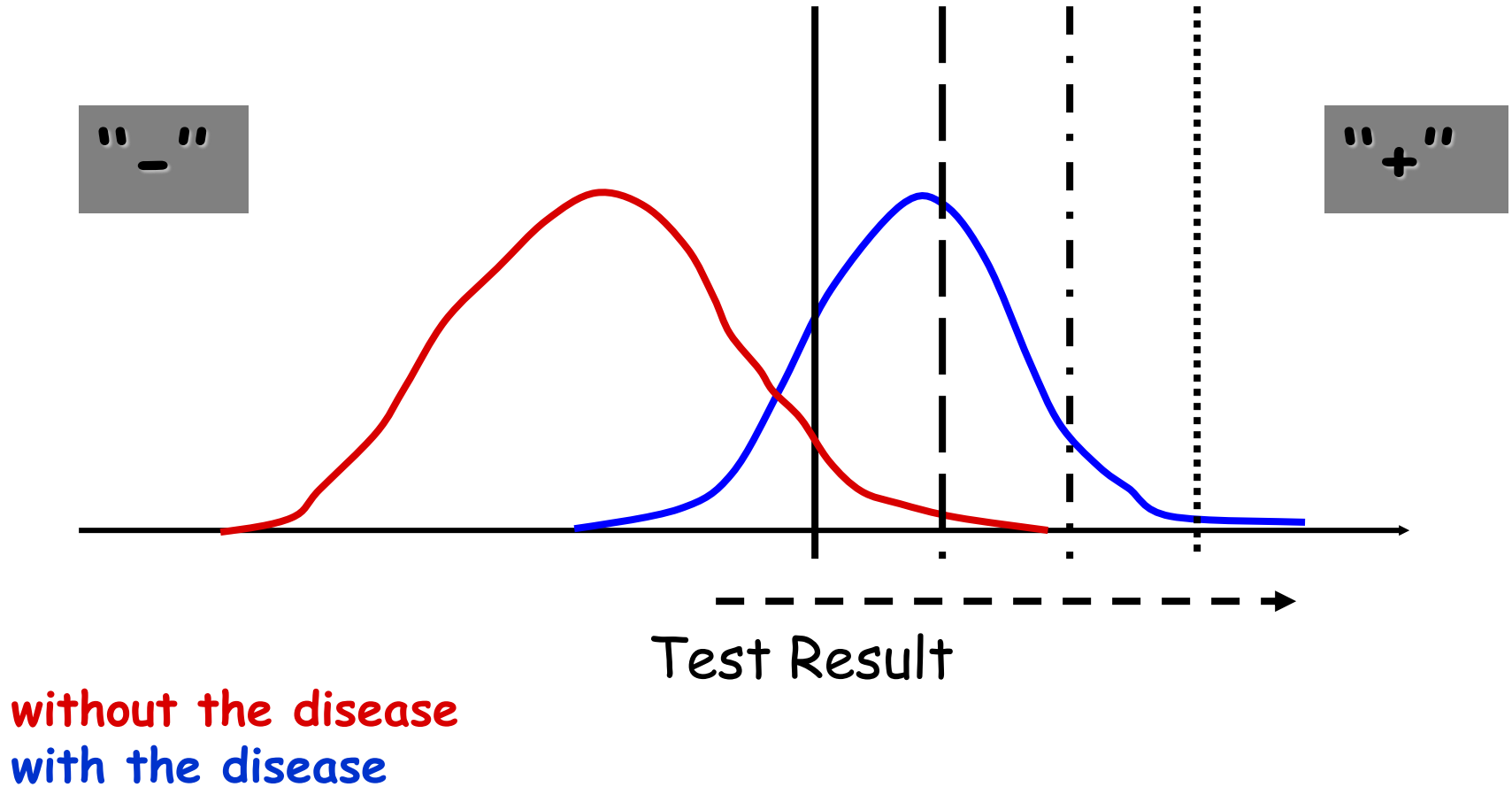Call these patients "negative"    Call these patients "positive"
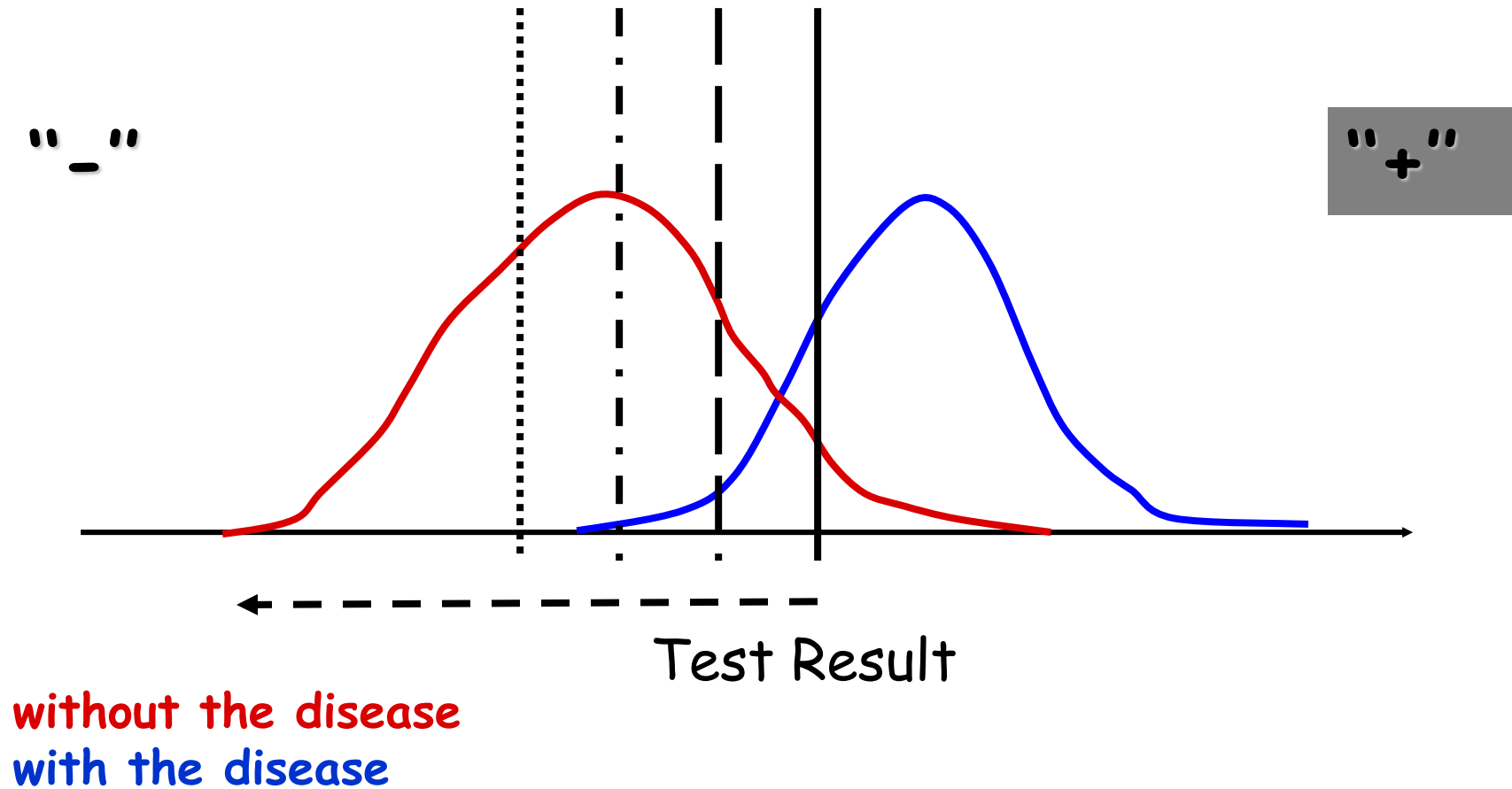
False negatives

Test Result

without the disease
with the disease
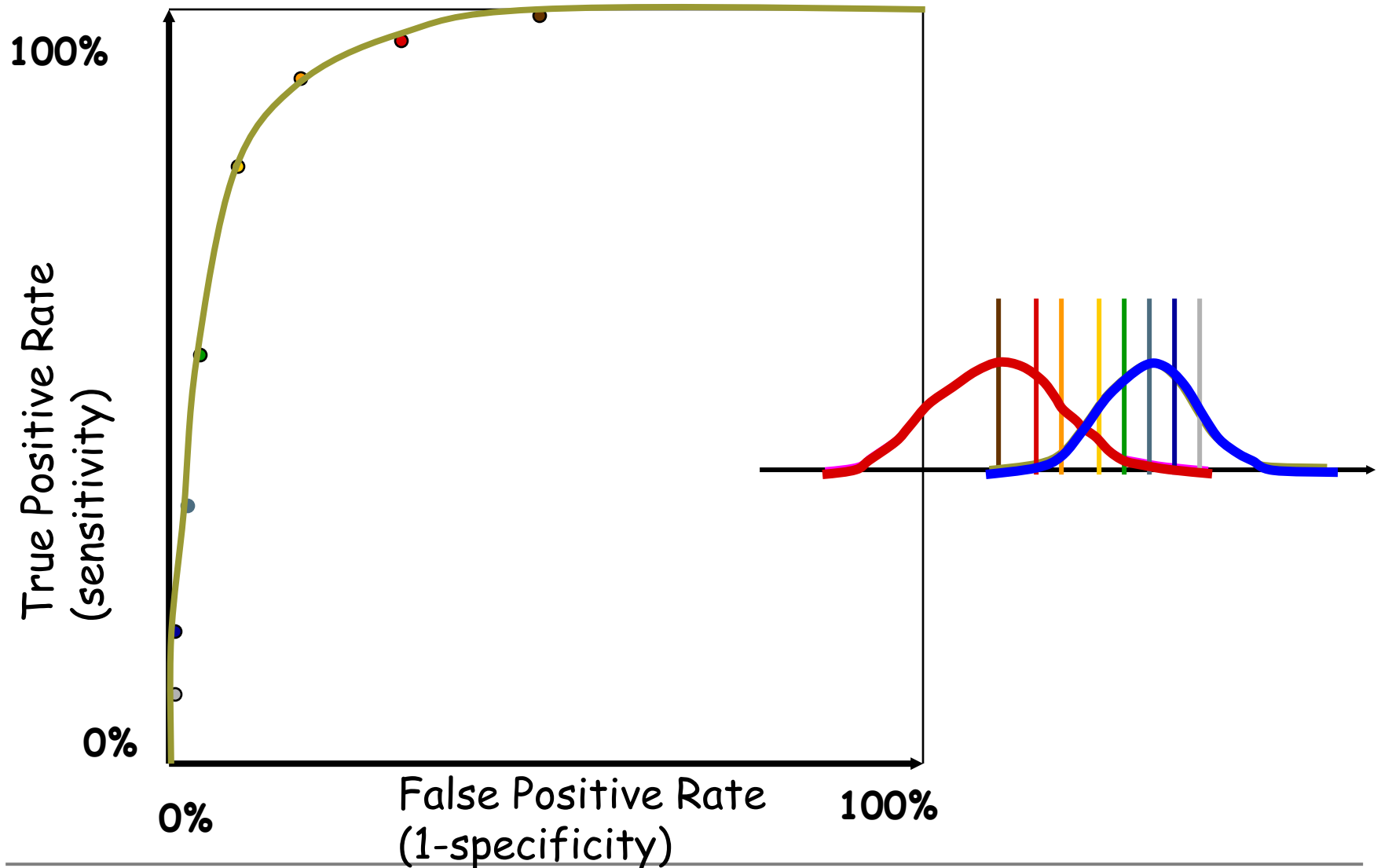
# Moving the Threshold: right



"–"

"+"

Test Result

without the disease
with the disease

# Moving the Threshold: left

"−"

"+"

Test Result

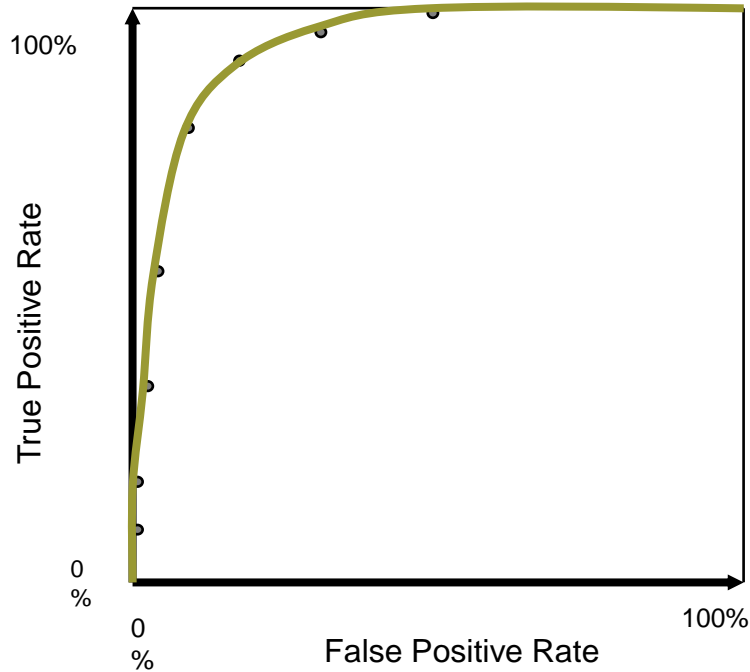without the disease
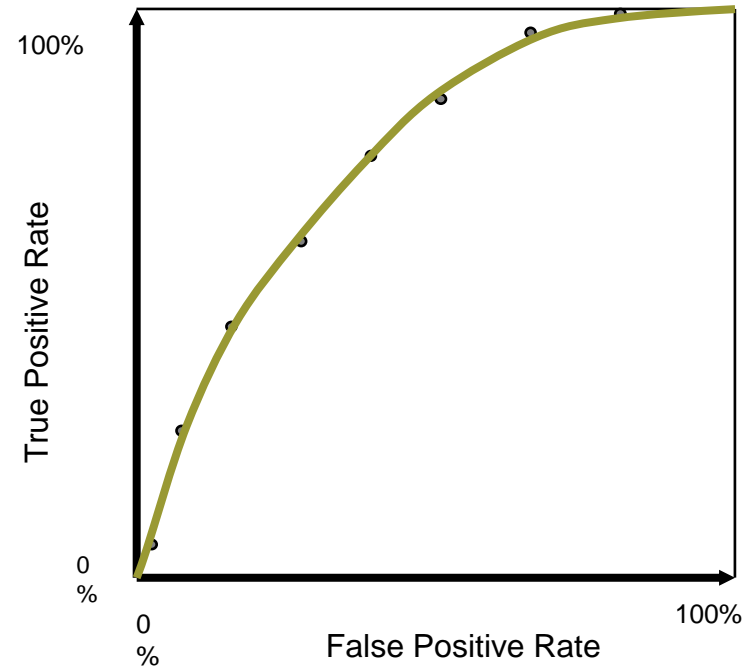with the disease

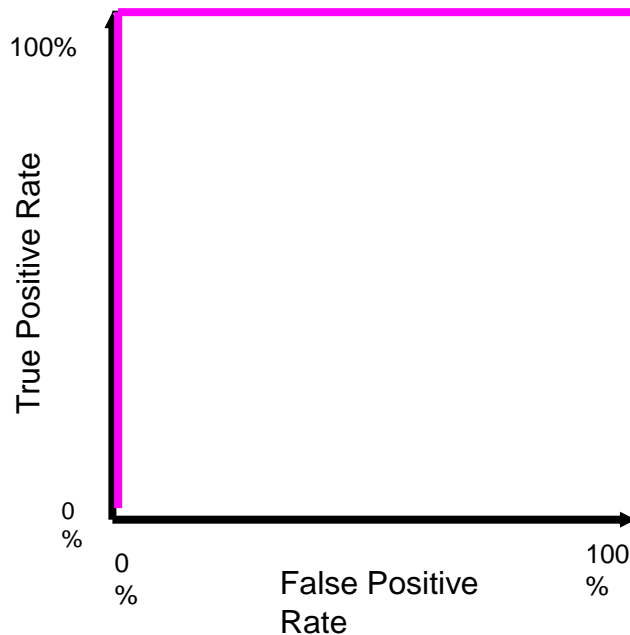# ROC curve

# ROC curve comparison
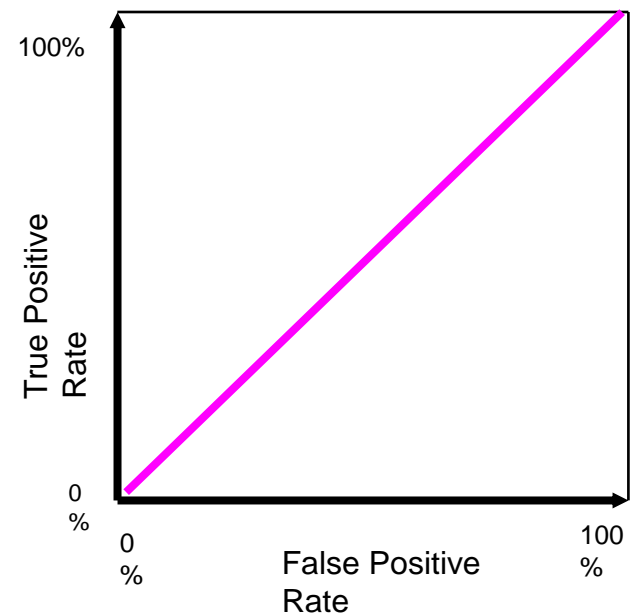
A good test:

A poor test:

# ROC curve extremes

Best Test:

Worst test:



The distributions don't overlap at all

The distributions overlap completely

# ROC Algorithm Conceptually

**Algorithm 1** Conceptual method for calculating an ROC curve. See algorithm 2 for a practical method.

**Inputs:** $L$, the set of test instances; $f(i)$, the probabilistic classifier's estimate that instance $i$ is positive; $min$ and $max$, the smallest and largest values returned by $f$; $increment$, the smallest difference between any two $f$ values.

```
1:  for t = min to max by increment do
2:      FP ⇐ 0
3:      TP ⇐ 0
4:      for i ∈ L do
5:          if f(i) ≥ t then                    /* This example is over threshold */
6:              if i is a positive example then
7:                  TP ⇐ TP + 1
8:              else                             /* i is a negative example, so this is a false positive */
9:                  FP ⇐ FP + 1
10:     Add point (FP/N, TP/P) to ROC curve
11: end
```

- What is the **computational complexity** of this algorithm?

# Practical ROC Algorithm

**Algorithm 2** Practical method for calculating an ROC curve from a test set

**Inputs:** $L$, the set of test instances; $f(i)$, the probabilistic classifier's estimate that instance $i$ is positive.

**Outputs:** $R$, a list of ROC points from (0,0) to (1,1)

1: $L_{sorted} \Leftarrow L$ sorted decreasing by $f$ scores
2: $FP \Leftarrow 0$
3: $TP \Leftarrow 0$
4: $R \Leftarrow \langle \rangle$
5: $f_{prev} \Leftarrow -\infty$
6: **for** $i \in L_{sorted}$ **do**
7:   **if** $f(i) \neq f_{prev}$ **then**
8:     ADD_POINT$\left(\left(\frac{FP}{N}, \frac{TP}{P}\right), R\right)$
9:     $f_{prev} \Leftarrow f(i)$
10:   **if** i is a positive example **then**
11:     $TP \Leftarrow TP + 1$
12:   **else**                                    /* i is a negative example, so this is a false positive */
13:     $FP \Leftarrow FP + 1$
14: ADD_POINT$\left(\left(\frac{FP}{N}, \frac{TP}{P}\right), R\right)$
15: **end**

1: **subroutine** ADD_POINT$(P, R)$
2: push $P$ onto $R$
3: **end subroutine**

- What is the **computational complexity** of this algorithm?
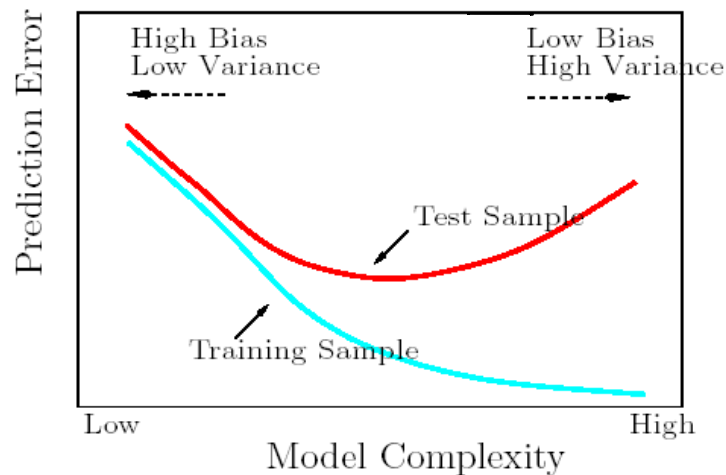
# Model Assessment and Selection

# Introduction

- **_Generalization_** performance of a learning method:
    - Measure of prediction capability on independent test data
    - Guide model selection
    - Give quantitative assessment of chosen model
- Chapter shows key methods and how to apply them to model selection

# Bias, Variance, and Model Complexity



Figure 7.1: *Behavior of test sample and training sample error as the model complexity is varied.*

**Generalization**: test sample vs. training sample performance

- Training data usually monotonically increasing performance with model complexity

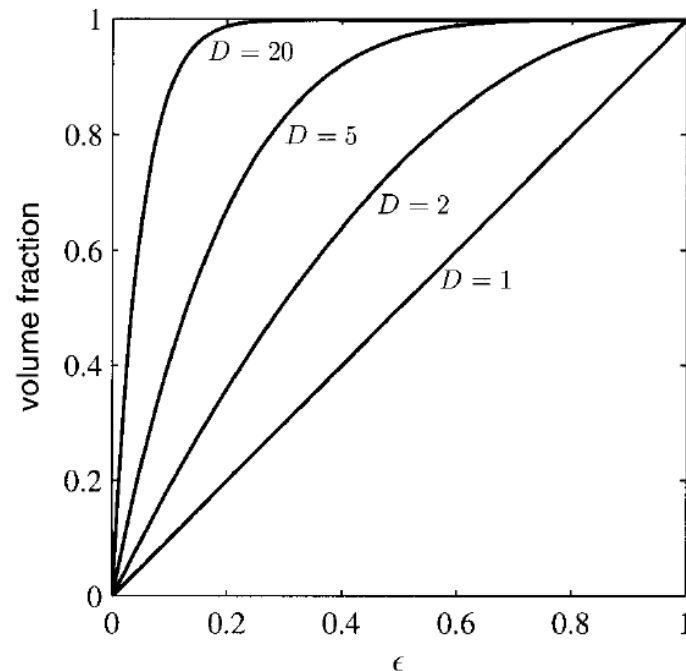- Bias and Variance are a **tradeoff**
- **Bias** in this context refers to the difference between the model prediction and the actual data
- **Variance** in this context refers to the sensitivity of the model's predictions to realizations of the model

# The Curse of Dimensionality

- As the dimensionality of the problem increases, most of the volume (and therefore data) lies within a thin shell at the surface of a hypersphere

Plot of the fraction of the volume of a sphere lying in the range $r = 1-\epsilon$ to $r = 1$ for various values of the dimensionality $D$.



- This implies that data is **very sparse** throughout most of the volume!
- Therefore it is hard to find a **representative data sample** for a high-dimensional space.

# Measuring Performance

■ Target variable $Y$

■ Vector of inputs $X$

■ Prediction model $\hat{f}(X)$

■ Typical Choices of Loss function

$$L\left(Y, \hat{f}(X)\right) = \begin{cases} \left(Y - \hat{f}(X)\right)^2 & \textit{squared error} \\ \left|Y - \hat{f}(X)\right| & \textit{absolute error} \end{cases}$$

# Generalization Error

- ## Test error aka. Generalization error

$$Err = E\left[ L\left( Y, \hat{f}\left( X \right) \right) \right]$$

- Note: This **expectation** averages anything that is random, including the randomness in the training sample that it produced

- ## Training error

$$\overline{err} = \frac{1}{N} \sum_{i=1}^{n} L\left( y_i, \hat{f}\left( x_i \right) \right)$$

- **average loss** over training sample
- not a good estimate of **test error** (next slide)

# Training Error



Figure 7.1: *Behavior of test sample and training sample error as the model complexity is varied.*

■ Choosing a model by examining **training error** leads to **overfitting**

❑ not a good estimate of test error

❑ consistently decreases with model complexity
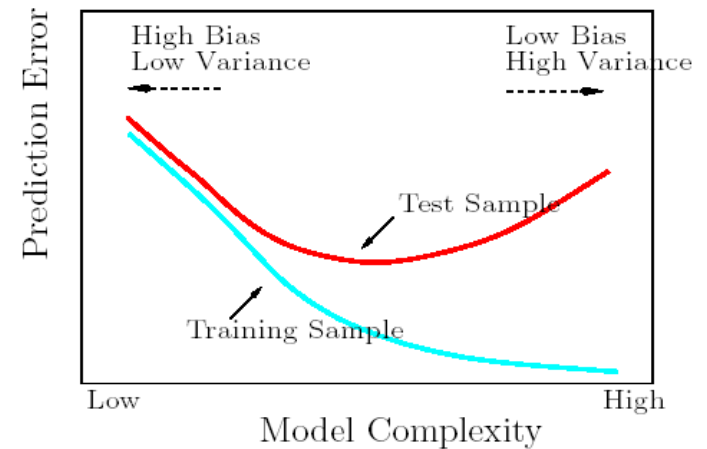
❑ drops to zero with high enough complexity

■ **Overfitted model**:
the model describes **noise** in the data rather than the underlying relationship

# Two separate goals

- **Model selection**:
  - Estimating the performance of different models in order to choose the (approximate) best one
- **Model assessment**:
  - Having chosen a final model, estimating its prediction error (generalization error) on new data

- Ideal situation: Data Rich:
  - split data into the 3 parts for *training*, *validation (est. prediction error+select model)*, and *testing (assess model)*
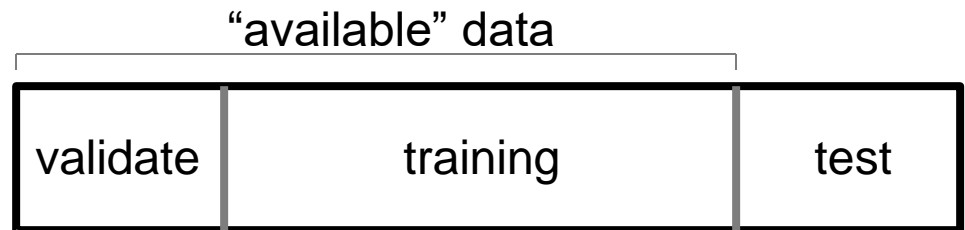  - Typical split*: 50% / 25% / 25%*

- Data-Poor situation
  - *Approximation of validation* step either analytically (AIC, BIC, MDL, SRM)
  - AIC/BIC - http://www.youtube.com/watch?v=YkD7ydzp9_E
  - or by efficient sample reuse (cross-validation, bootstrap)

# Cross Validation

- Used for determining the values of parameters of a model
- Also used for selecting a model
- The idea is simple: split all available data into 3 pieces – training, validation and testing

"available" data

| validate | training | test |
|---|---|---|

The technique of $S$-fold cross-validation, illustrated here for the case of $S = 4$, involves taking the available data and partitioning it into $S$ groups (in the simplest case these are of equal size). Then $S - 1$ of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all $S$ possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the $S$ runs are then averaged.

run 1
run 2
run 3
run 4