



---

# **MACHINE LEARNING**

## **CHAPTER 1: PRELIMINARY**

---

# Learning Objectives

---

- 1、 What is pattern recognition and machine learning?
  - 2、 What are curve fitting and regularization?
  - 3、 What are ML and MAP Bayesian inferences?
  - 4、 How to deal with the curse of dimensionality?
  - 5、 What is the relationship between decision theory and machine learning?
  - 6、 What are generative and discriminative models?
  - 7、 How to use entropy、 KL divergence and mutual information for machine learning?
-

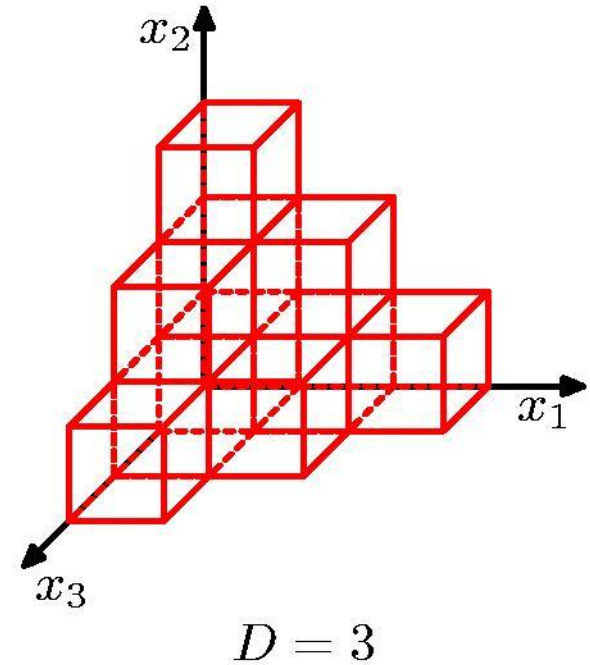
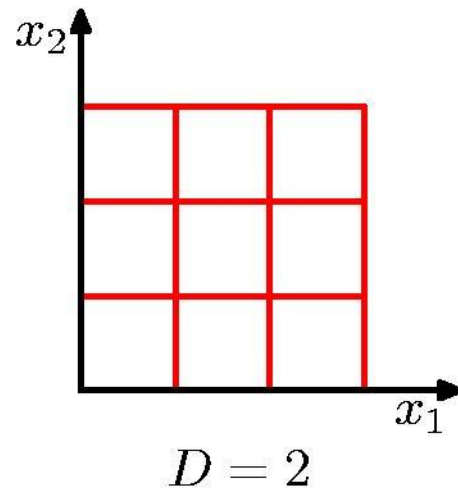
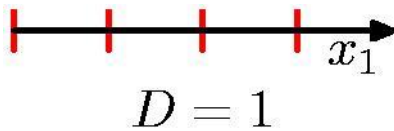
# Outlines

---

- Pattern Recognition and Machine Learning
  - Curve Fitting and Regularization
  - Probabilities and Gaussian Distributions
  - Bayesian Inferences (ML and MAP)
  - Curse of Dimensionality
  - Decision Theory
  - Entropy and Information
-

# Curse of Dimensionality

---



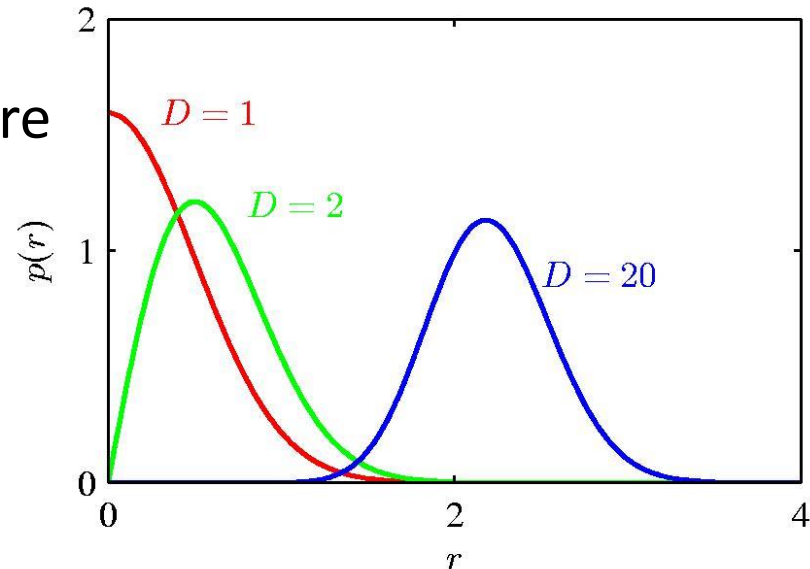
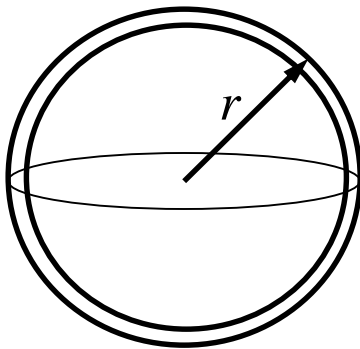
# Curse of Dimensionality

---

Polynomial curve fitting,  $M = 3$

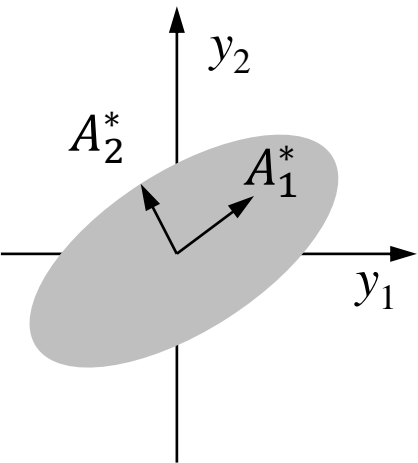
$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

Gaussian Densities in  
higher dimensions of a sphere



# Reduction of Dimensionality (PCA)

---

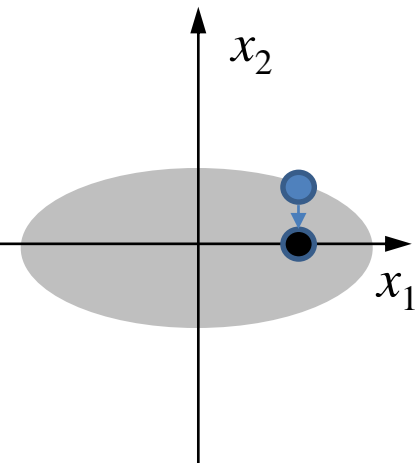


Basis  
Coefficients

Data  $\rightarrow$

$$Y = AX$$

principal component analysis



$$\max_{A_i} A_i^T \text{COV}(Y_n) A_i$$

$A$  : rotation ( $I \times I$ )

$Y$ : data ( $I \times N$ )

$A_i^*$  : optimal solution  
eigenvector

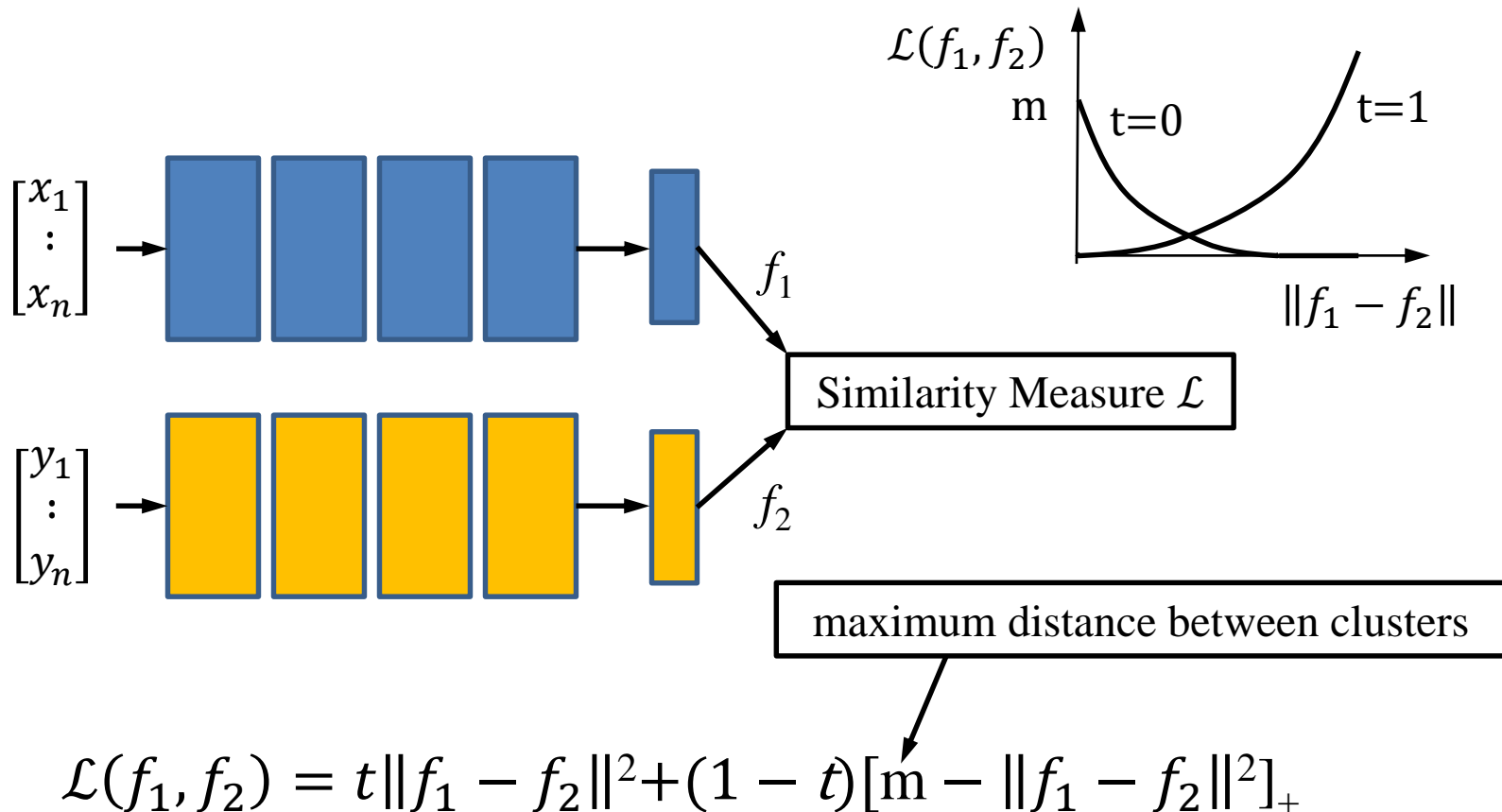
$$A_i^{*T} \text{COV}(Y_n) A_i^* = \lambda_i$$

$$s.t. \quad A_i^T A_i = 1 \quad A_i^T A_j = 0 \quad E[Y_n] = \mathbf{0}$$


---

# Feature Extraction (Contrastive Loss)

---

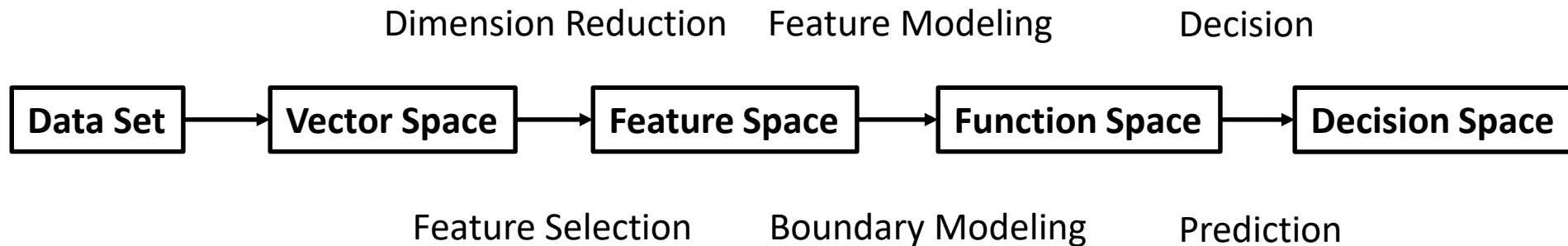


$t=1$ : two vectors belong to the same category;  $[ ]_+$ : non-negative

---

# Machine Learning Pipeline

---





# Outlines

---

- Pattern Recognition and Machine Learning
  - Curve Fitting and Regularization
  - Probabilities and Gaussian Distributions
  - Bayesian Inferences (ML and MAP)
  - Curse of Dimensionality
  - Decision Theory
  - Entropy and Information
-

# Decision Theory

---

## Inference step

Determine either  $p(t|\mathbf{x})$  or  $p(\mathbf{x}, t)$ .

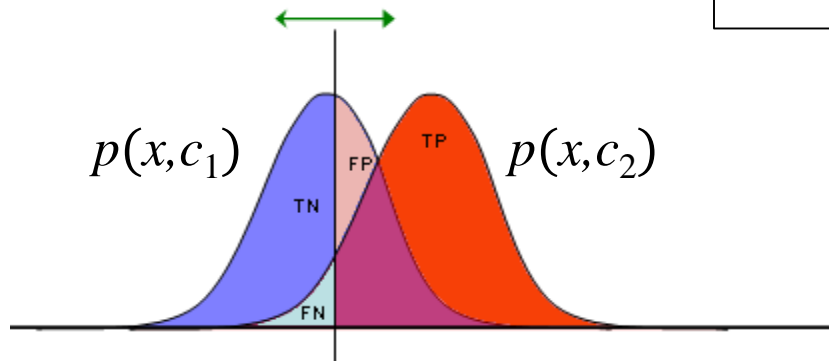
## Decision step

For given  $\mathbf{x}$ , determine optimal  $t$ .

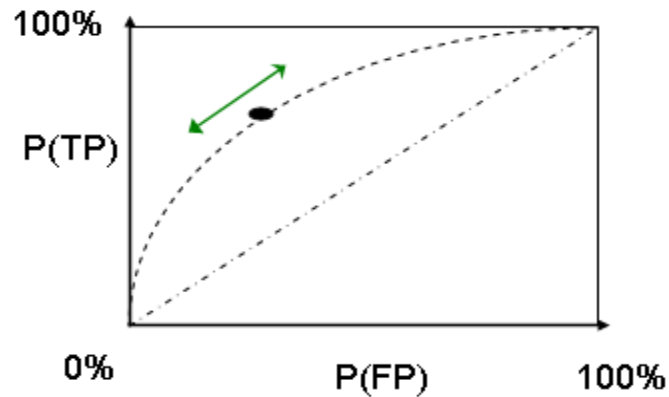
---

# Receiver Operating Characteristic Curve

$$\begin{aligned} p(x) &= p(x, c_1) + p(x, c_2) \\ &= p(x/c_1) p(c_1) + p(x/c_2) p(c_2) \end{aligned}$$

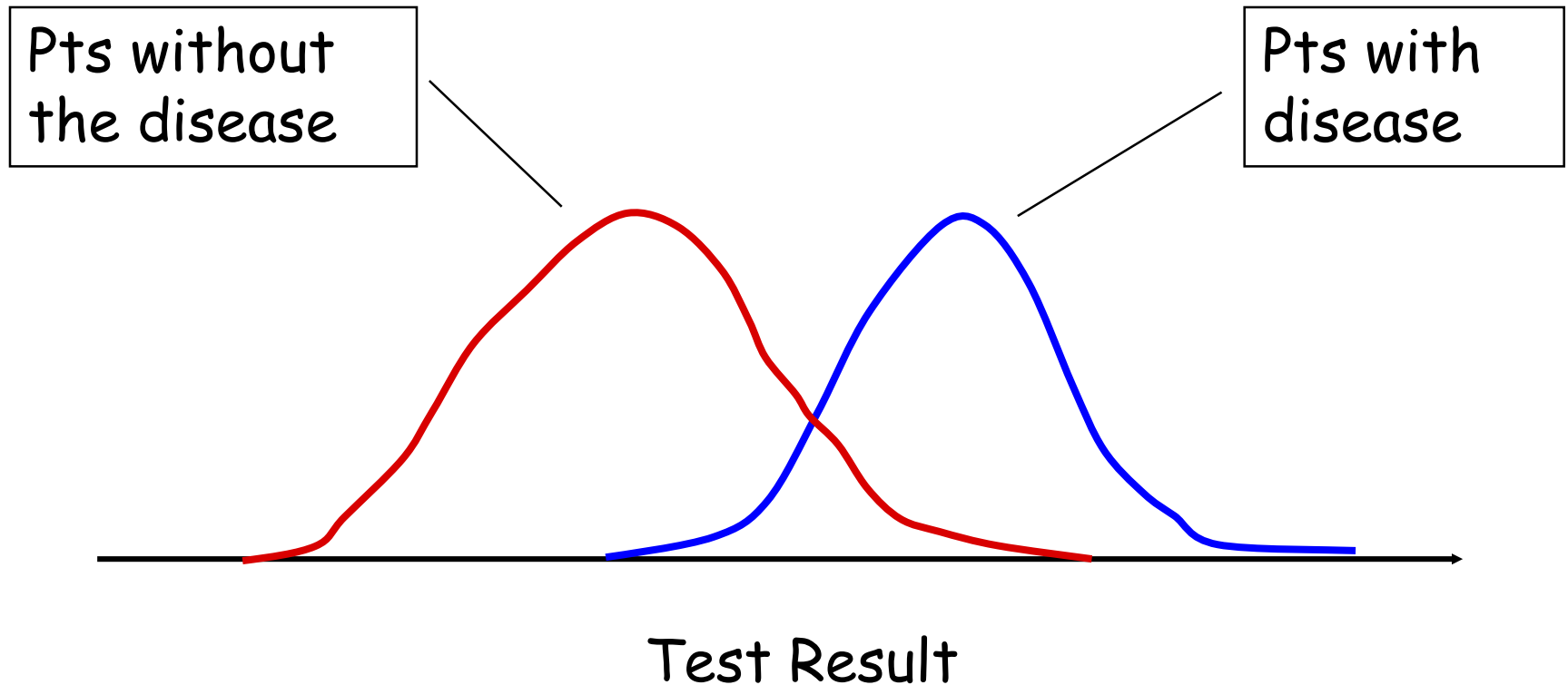


TP	FP
FN	TN
1	1



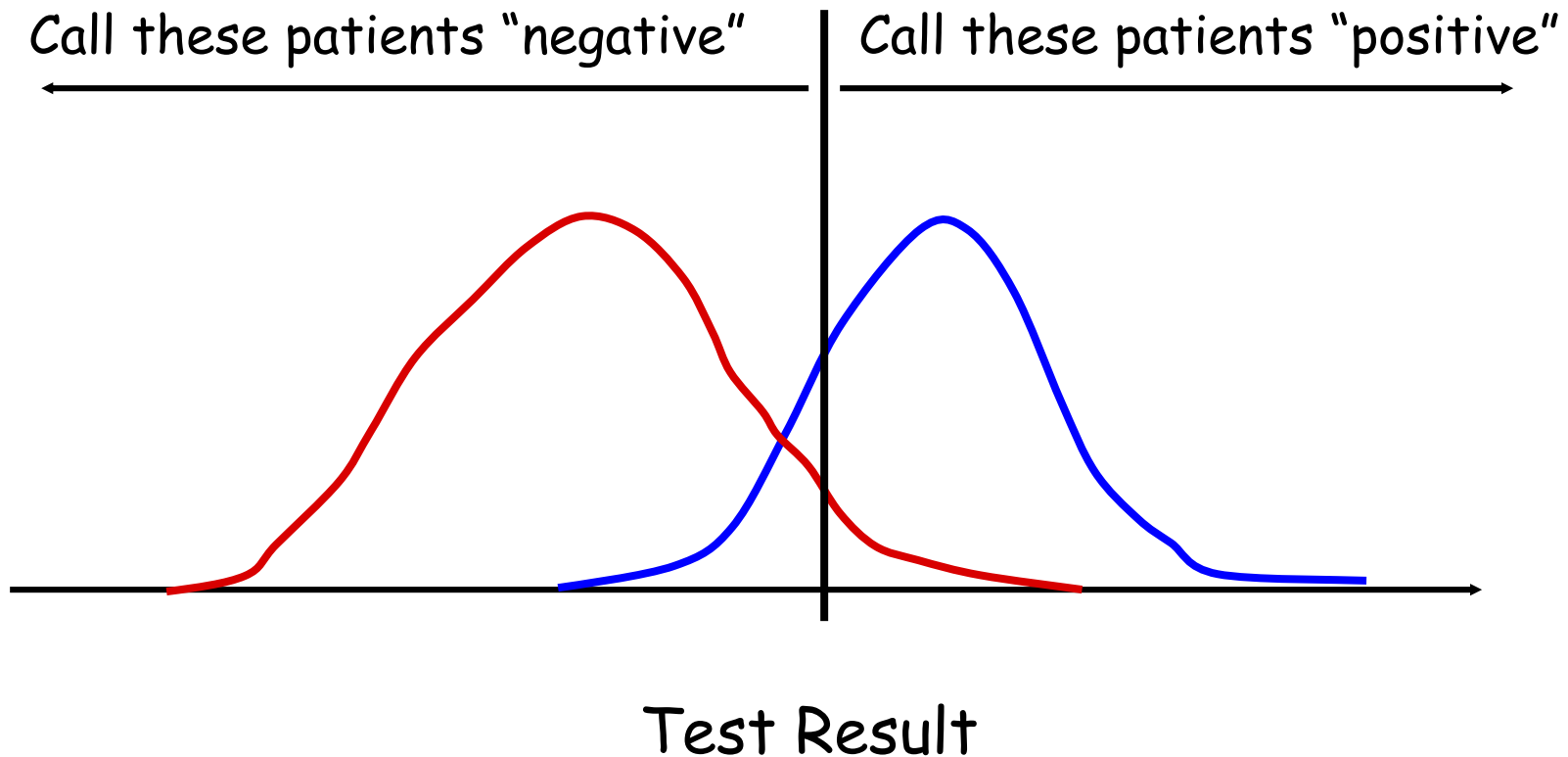
# Bimodal Distribution (Data Model)

---



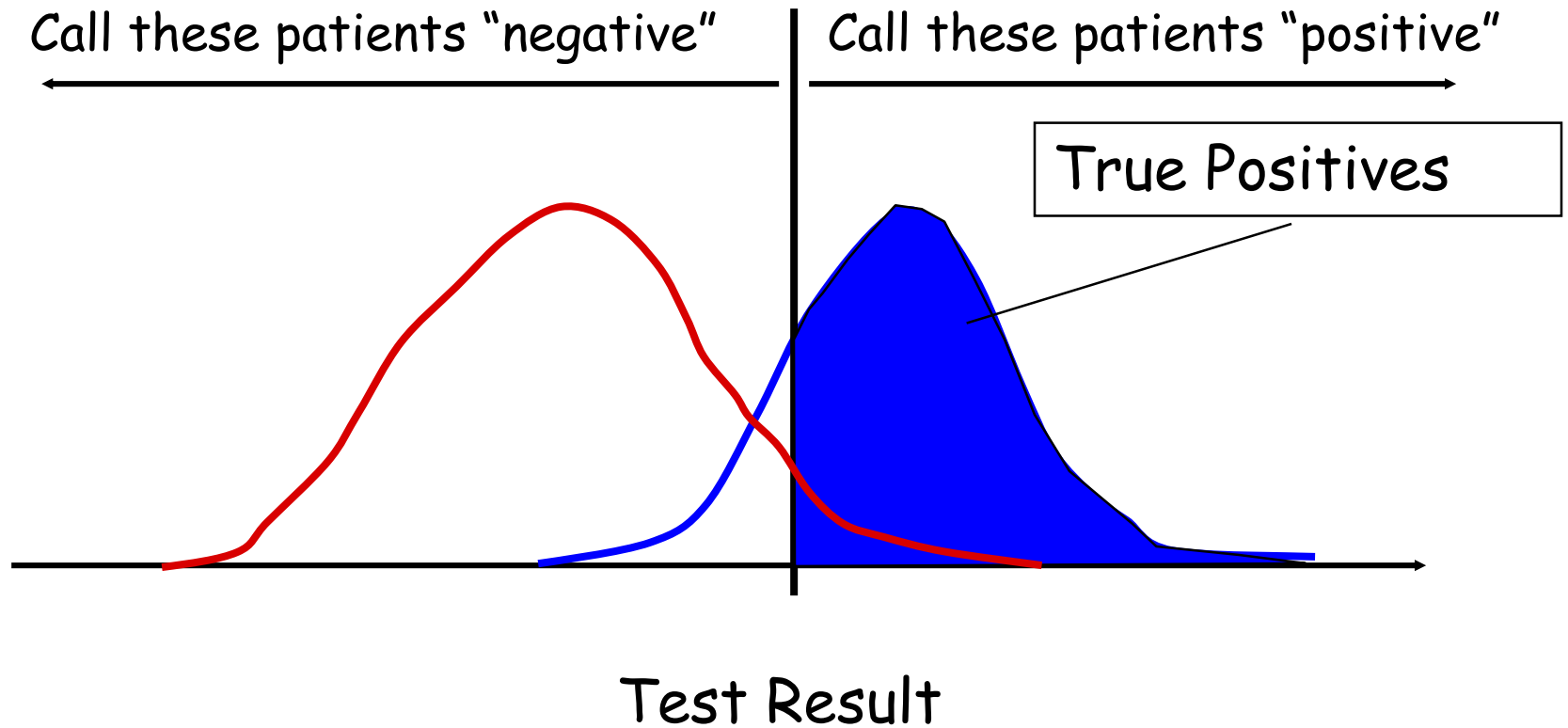
# Decision Threshold (Boundary Model)

---



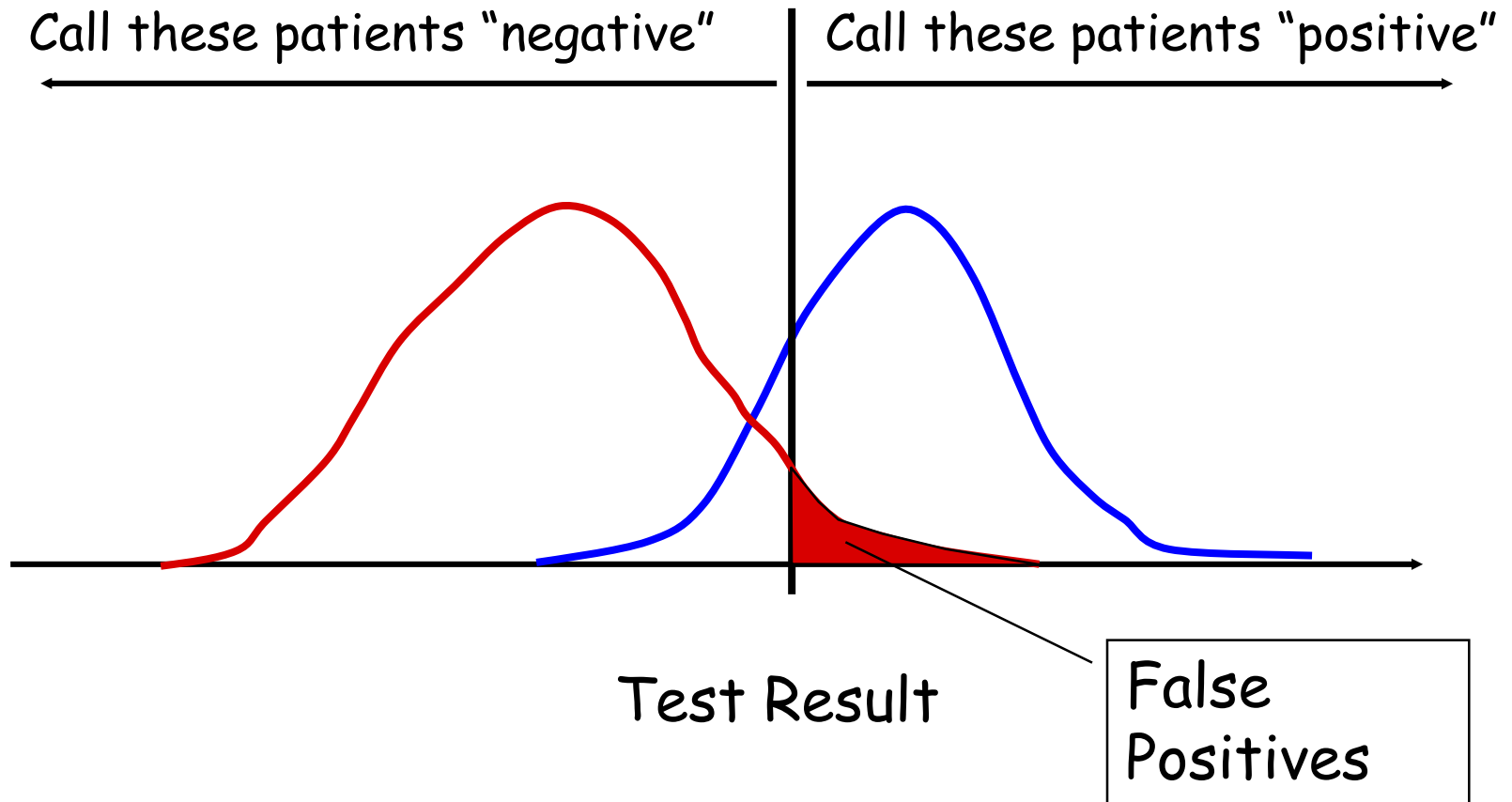
# True Positive

---



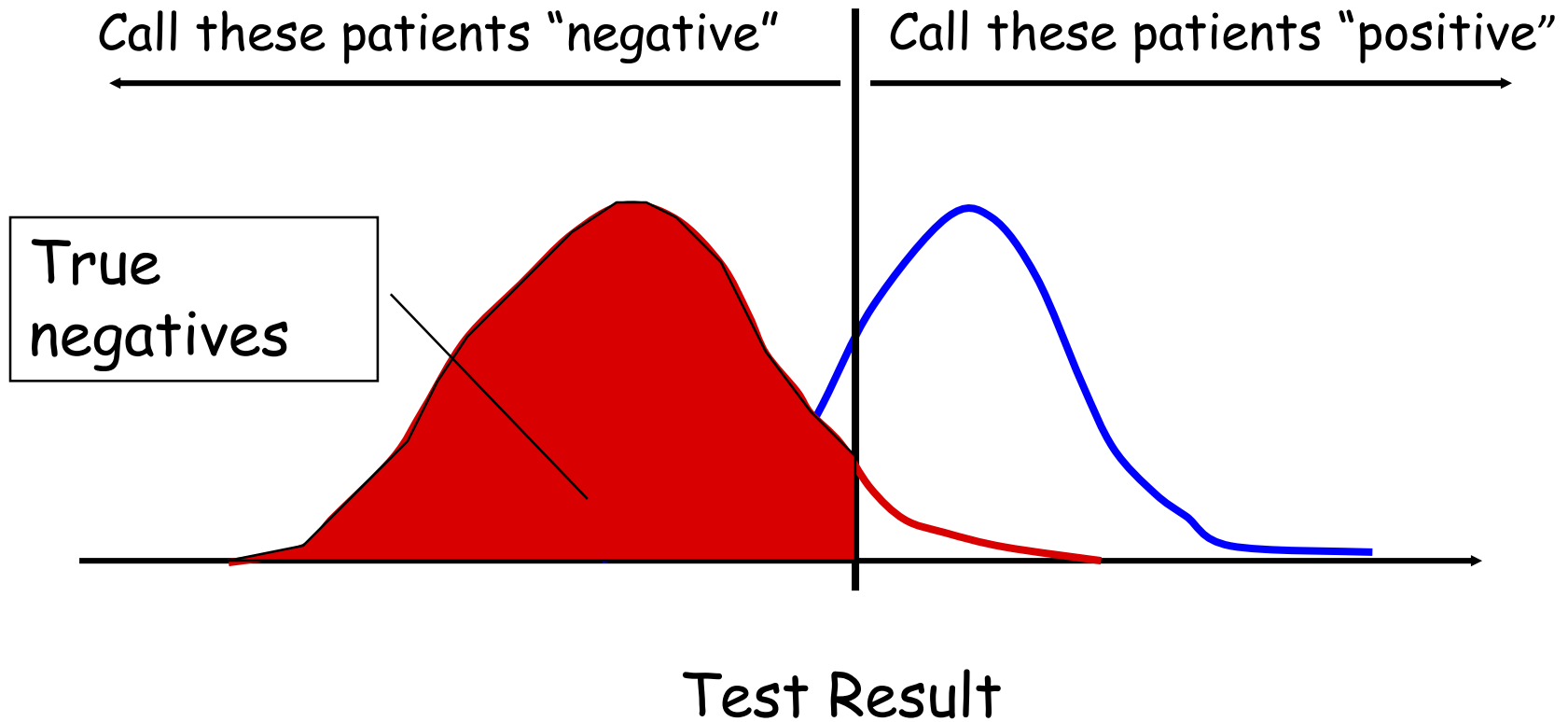
# False Positive

---



# True Negative

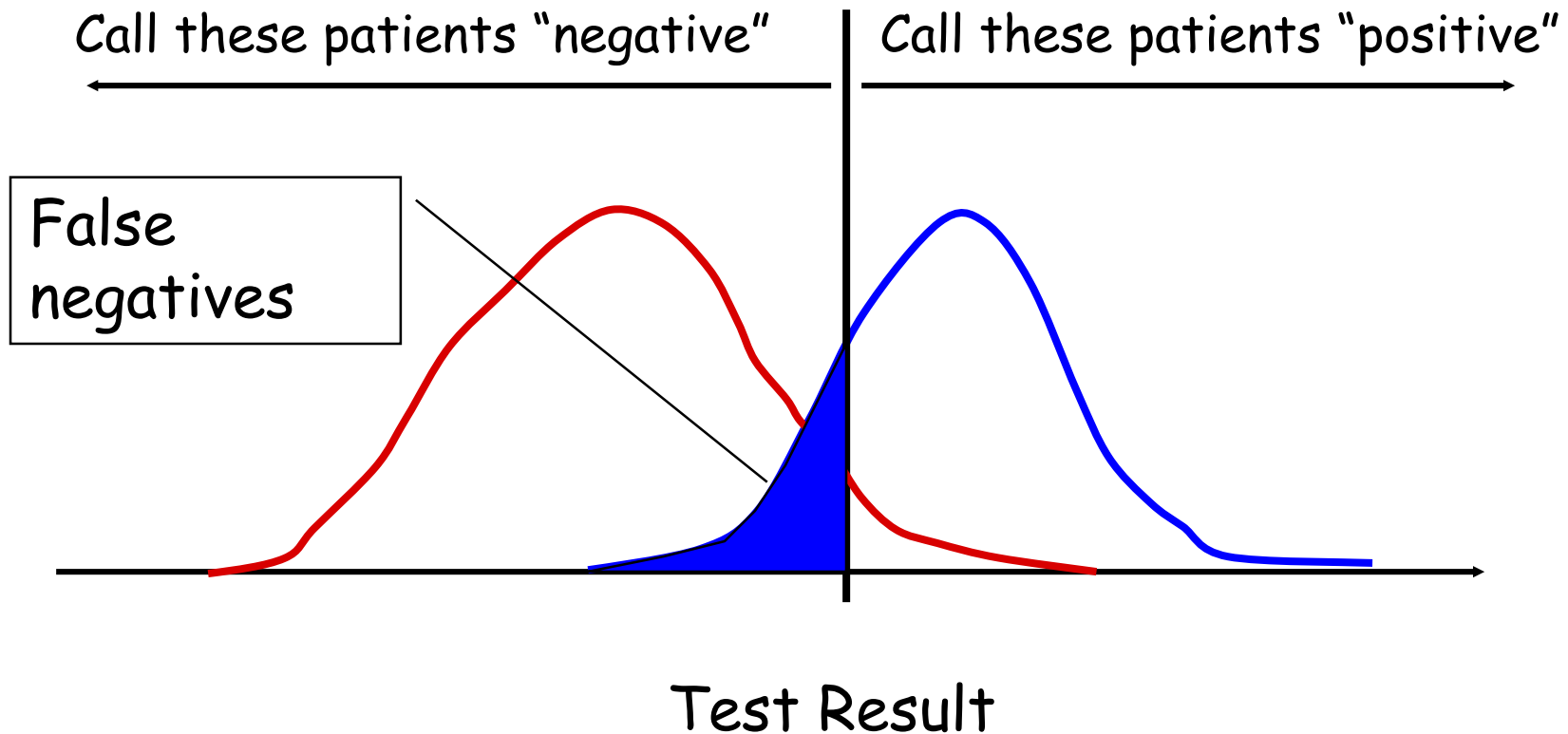
---





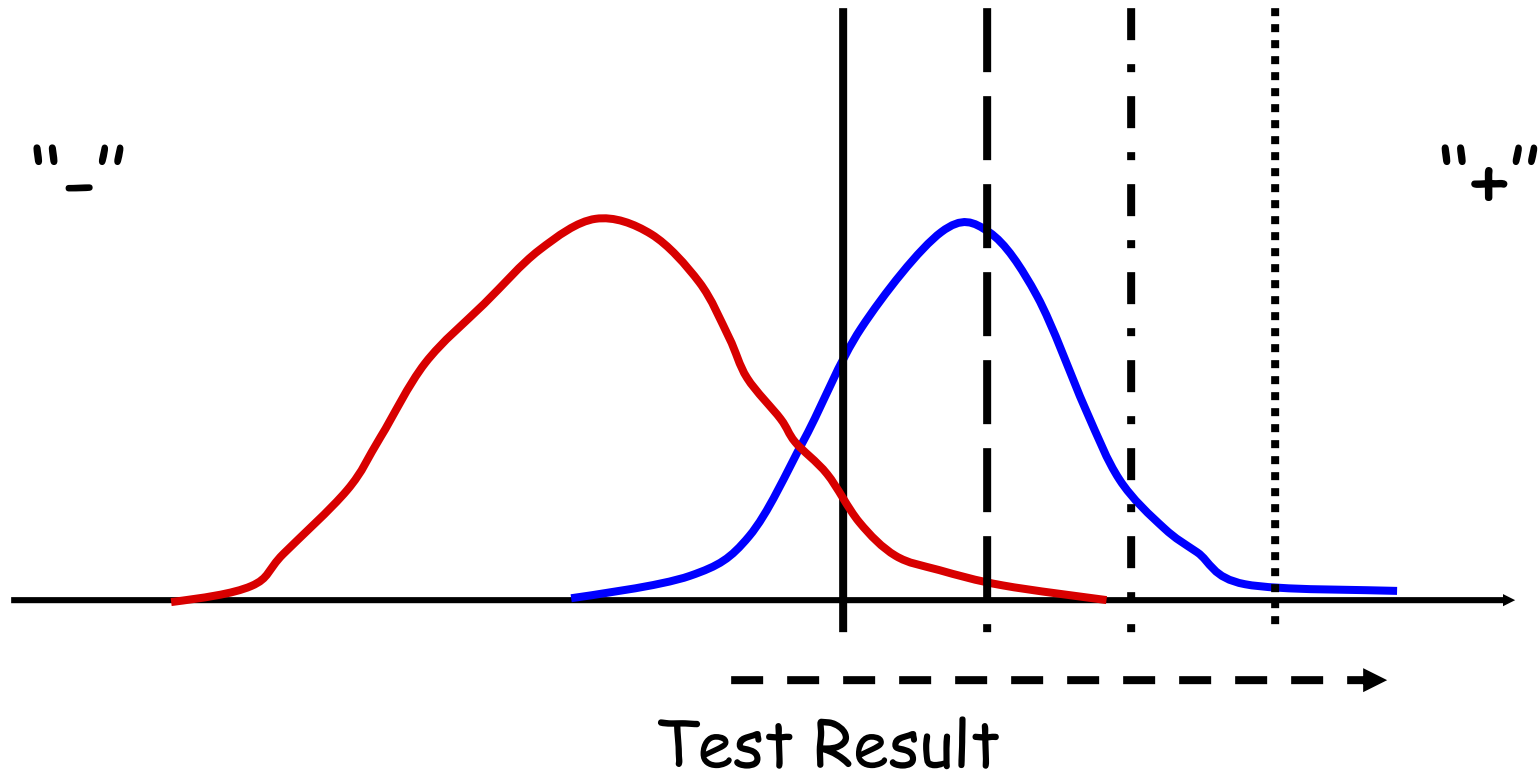
# False Negative

---



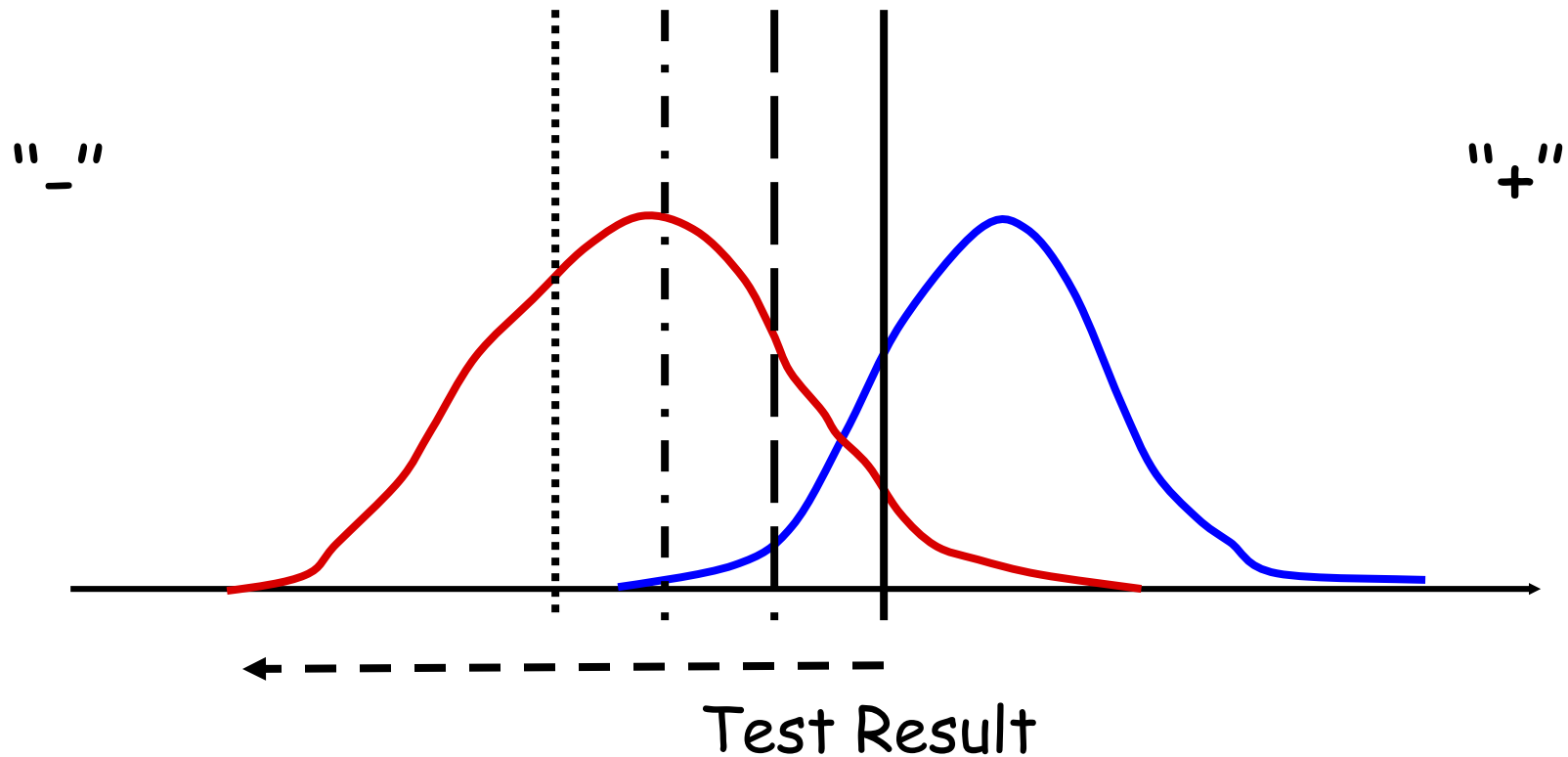
# Moving the Threshold: Right

---

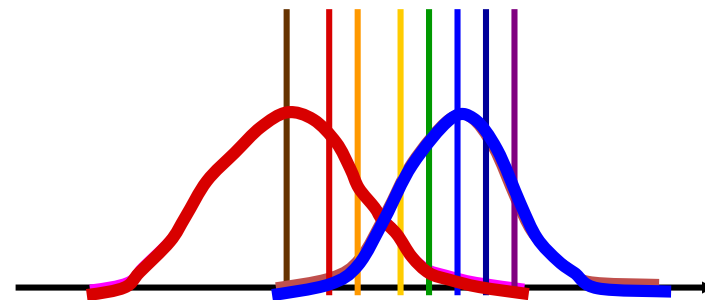
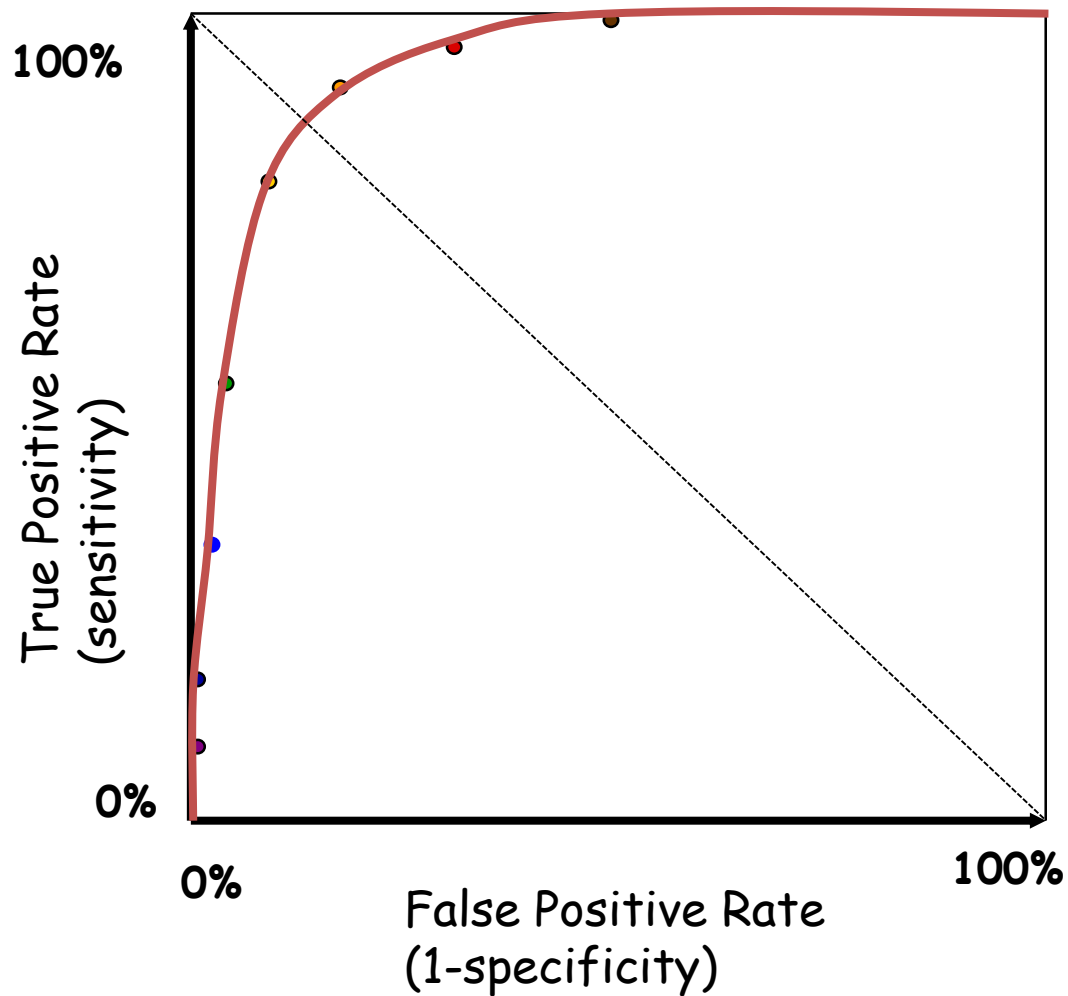


# Moving the Threshold: Left

---

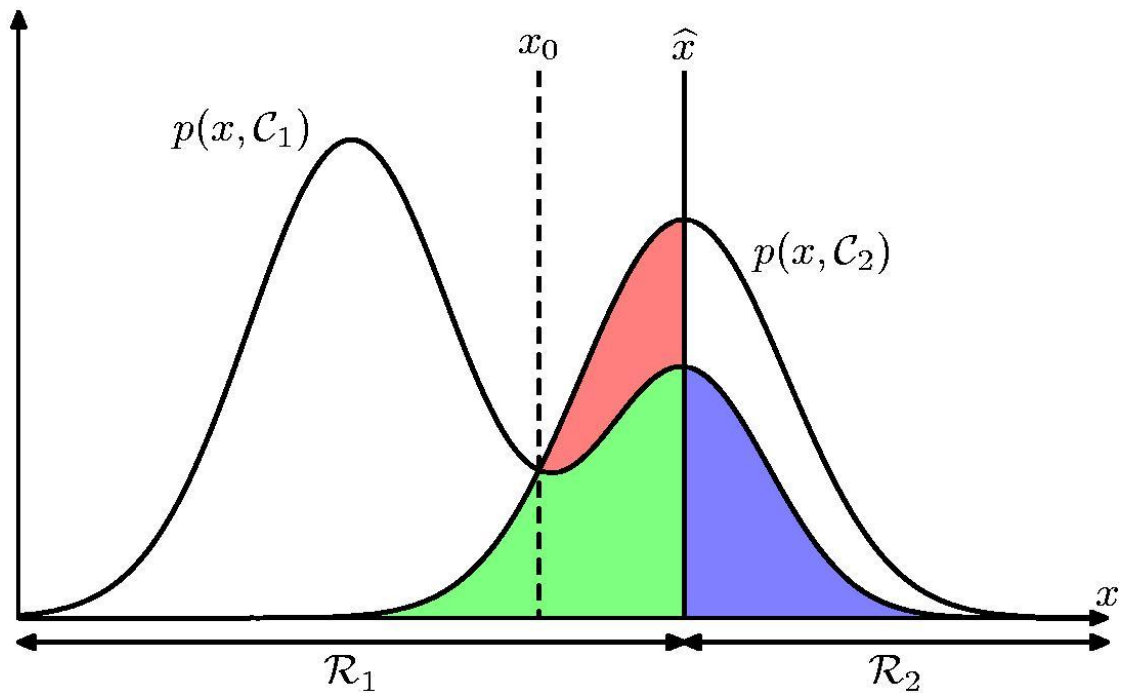


# ROC Curve



# Minimum Misclassification Rate

---



$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \end{aligned}$$

---

# Minimum Expected Loss

---

Example: classify medical images as 'cancer' or 'normal'

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

False  
Positive

False  
Negatives

The diagram illustrates the loss function for a binary classification task. The matrix shows that misclassifying a cancer image as normal (False Negative) incurs a loss of 1000, while misclassifying a normal image as cancer (False Positive) incurs a loss of 1. The boxes 'False Positive' and 'False Negatives' are connected to the corresponding cells in the matrix.

# Minimum Expected Loss

---

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

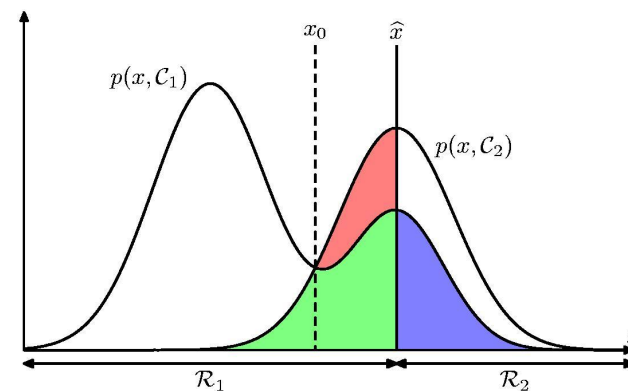
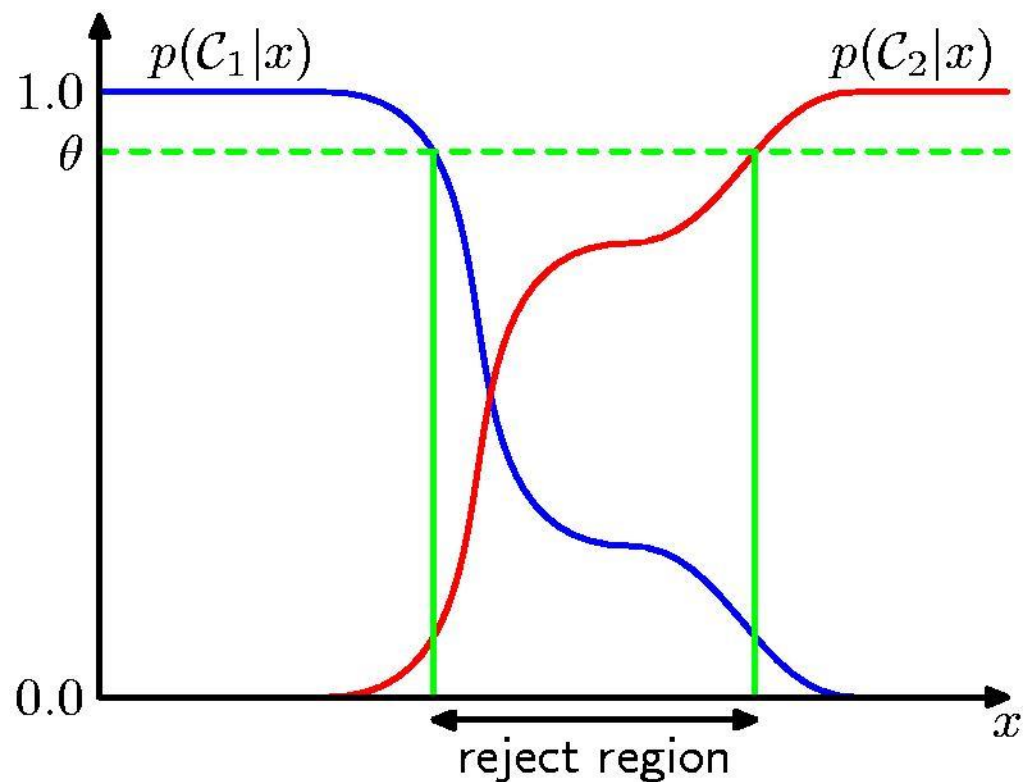
Regions  $\mathcal{R}_j$  are chosen to minimize

$$\mathbb{E}[L] = \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

---

# Reject Option

---





# Posterior Probability

---

The posterior probability of class 1 is given by:

$$p(C_1 | \mathbf{x}) = \frac{p(C_1)p(\mathbf{x} | C_1)}{p(C_1)p(\mathbf{x} | C_1) + p(C_0)p(\mathbf{x} | C_0)} = \frac{1}{1 + e^{-z}} = \sigma(z)$$

$$\text{where } z = \ln \frac{p(C_1)p(\mathbf{x} | C_1)}{p(C_0)p(\mathbf{x} | C_0)} = \ln \frac{p(C_1 | \mathbf{x})}{1 - p(C_1 | \mathbf{x})}$$



$z$  is called the logit and is given by the log odds

---

# Why Separate Inference and Decision?

---

- Minimizing risk (loss matrix may change over time)
- Reject option
- Unbalanced class priors
- Combining models

# Decision Theory for Regression

---

Inference step

Determine  $p(\mathbf{x}, t)$ .

Decision step

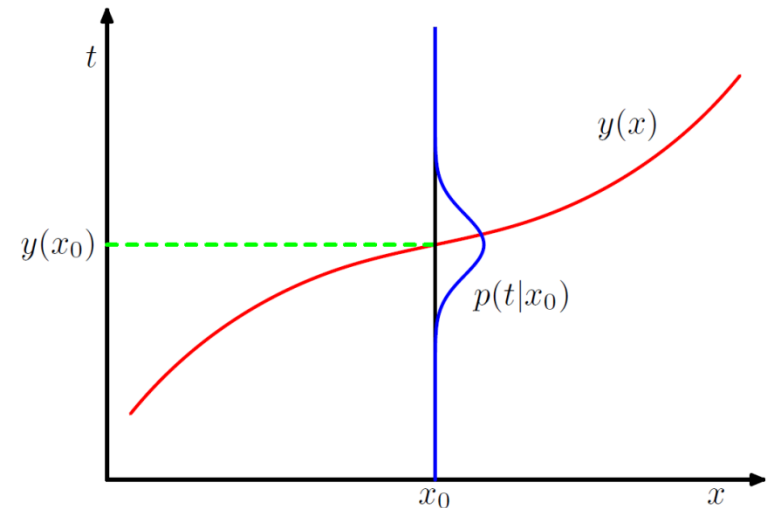
For given  $\mathbf{x}$ , make optimal prediction,  $y(\mathbf{x})$ , for  $t$ .

Loss function:  $\mathbb{E}[L] = \iint L(t, y(\mathbf{x}))p(\mathbf{x}, t) \, d\mathbf{x} \, dt$

---

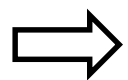
# The Expected Squared Loss Function

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$



$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 \end{aligned}$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) \, d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) \, d\mathbf{x}$$



$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$$

predictor

noise

$y(x)$  : an estimator of the mean of  $t$  for given  $\mathbf{x}$

# Generative vs Discriminative

---

## Generative approach:

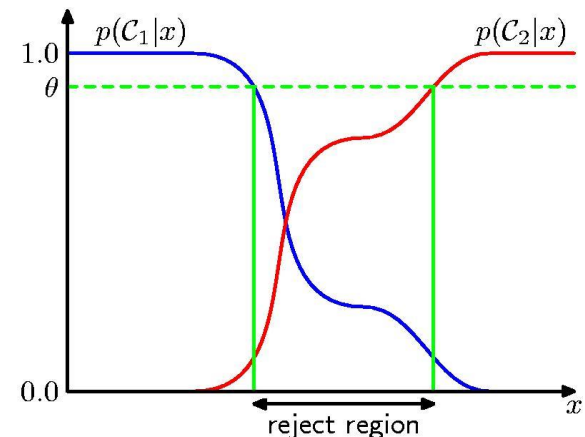
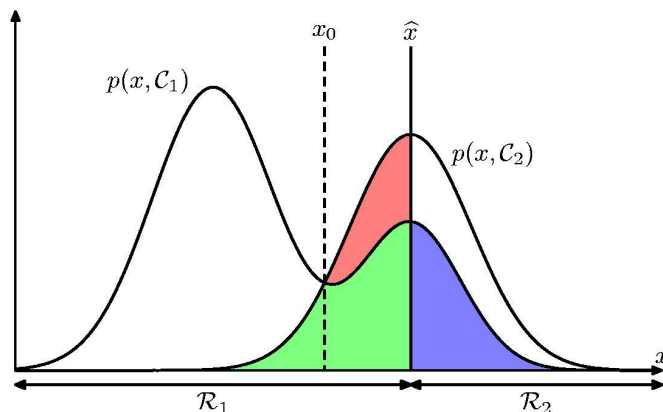
Model  $p(t, \mathbf{x}) = p(\mathbf{x}|t)p(t)$

Use Bayes' theorem  $p(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)p(t)}{p(\mathbf{x})}$

## Discriminative approach:

Model  $p(t|\mathbf{x})$  directly

$t$  : category



# Outlines

---

- Pattern Recognition and Machine Learning
  - Curve Fitting and Regularization
  - Probabilities and Gaussian Distributions
  - Bayesian Inferences (ML and MAP)
  - Curse of Dimensionality
  - Decision Theories
  - Entropy and Information
-

# Entropy

---

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning

# Entropy

---

Coding theory:  $x$  discrete with 8 possible states; how many bits to transmit the state of  $x$ ?

All states equally likely

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$



# Entropy

---

$x$	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$

---

# Entropy

---

In how many ways can  $N$  identical objects be allocated  $M$  bins?

$$W = \frac{N!}{\prod_i n_i!}$$

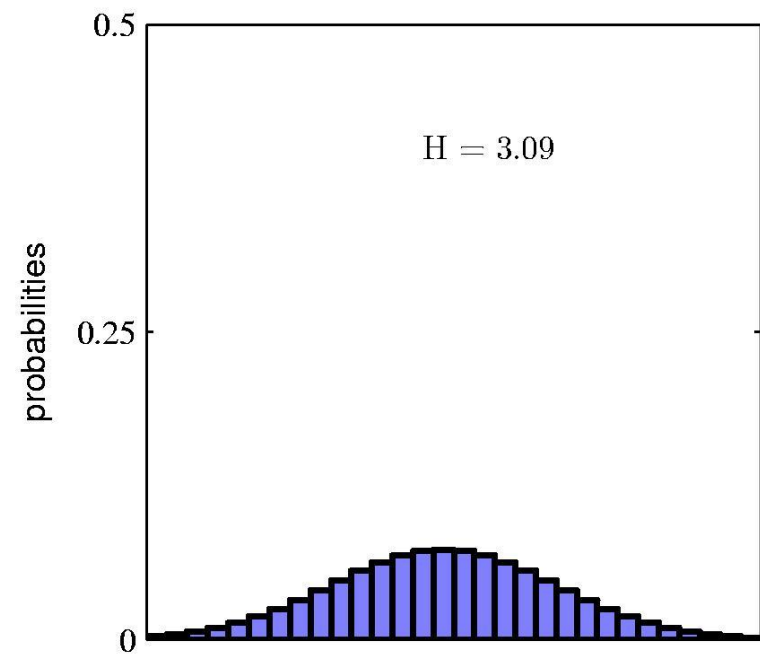
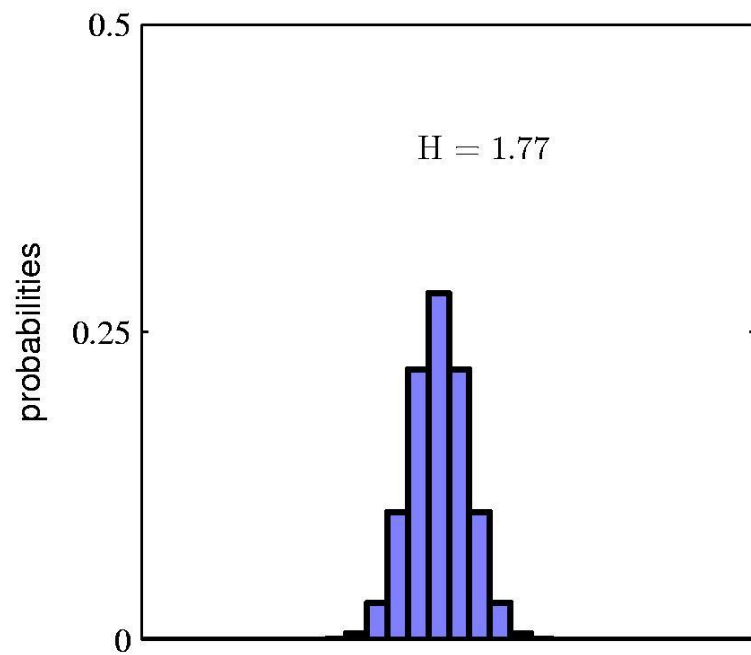
$$H = \frac{1}{N} \ln W \simeq - \lim_{N \rightarrow \infty} \sum_i \left( \frac{n_i}{N} \right) \ln \left( \frac{n_i}{N} \right) = - \sum_i p_i \ln p_i$$

Entropy maximized when  $\forall i : p_i = \frac{1}{M}$

---

# Entropy

---



# Differential Entropy

---

Put bins of width  $\Delta$  along the real line

$$\lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx$$

Differential entropy maximized (for fixed  $\sigma^2$ ) when

$$p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

in which case

$$H[x] = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\} .$$

---

# Conditional Entropy

---

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{y} \, d\mathbf{x}$$

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]$$

---

# The Kullback-Leibler Divergence

---

$$\begin{aligned} \text{KL}(p\|q) &= \overset{\text{Cross Entropy } C(p\|q)}{\underbrace{- \int p(\mathbf{x}) \ln q(\mathbf{x}) \, d\mathbf{x}}} - \overset{\text{Entropy } H(p)}{\underbrace{\left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} \right)}} \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} \, d\mathbf{x} \end{aligned}$$

$$\text{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^N \left\{ \overset{\text{Cross Entropy}}{\underbrace{- \ln q(\mathbf{x}_n|\boldsymbol{\theta})}} + \overset{\text{Negative Entropy}}{\underbrace{\ln p(\mathbf{x}_n)}} \right\}$$

$$\text{KL}(p\|q) \geq 0 \qquad \text{KL}(p\|q) \neq \text{KL}(q\|p)$$

KL divergence describes a distance between model  $p$  and model  $q$

---

# Cross Entropy for Machine Learning

---

Goal of Machine Learning:  $p(\text{real data}) \approx p(\text{model} | \theta)$

we assume:  $p(\text{training data}) \approx p(\text{real data})$

Operation of Machine Learning:  $p(\text{training data}) \approx p(\text{model} | \theta)$

$$\min_{\theta} \text{KL}(p(\text{training data}) || p(\text{model} | \theta))$$



$$\min_{\theta} C(p(\text{training data}) || p(\text{model} | \theta))$$

as  $H(p(\text{training data}))$  is fixed

# Cross Entropy for Machine Learning

---

$$C(p(\text{training data}) \parallel p(\text{model} | \theta))$$

Bernoulli model:  $p(\text{model} / \theta) = \rho^t (1 - \rho)^{1-t}$   $t_n$ : training data

Cross entropy :  $C = -\frac{1}{N} \sum_n t_n \ln \rho + (1 - t_n) \ln(1 - \rho)$   $\rho$ : model parameter

Gaussian model:  $p(\text{model} / \theta) \propto e^{-0.5(t-\mu)^2}$   $t_n$ : training data

Cross entropy :  $C \propto \frac{1}{N} \sum_n (t_n - \mu)^2$   $\mu$ : model parameter

---



# Mutual Information

---

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

Mutual information describes the degree of dependence between  $\mathbf{x}$  and  $\mathbf{y}$

---

# Information Gain

---



$H[\mathbf{x}]$ : uncertain of balls

$H[\mathbf{x}|\mathbf{y}]$ :  
uncertain of balls after  
weighing once

$\mathbf{x}$ : one ball lighter

$\mathbf{y}$ : weighing once

$\mathbf{x}|\mathbf{y}$ : one ball lighter  
after weighing once

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = \log_2 3$$

$$H[\mathbf{x}] = \log_2 N$$

After weighing  $\frac{N}{3}$  times, all the uncertainties can be removed

---

# Independent Signal Separation



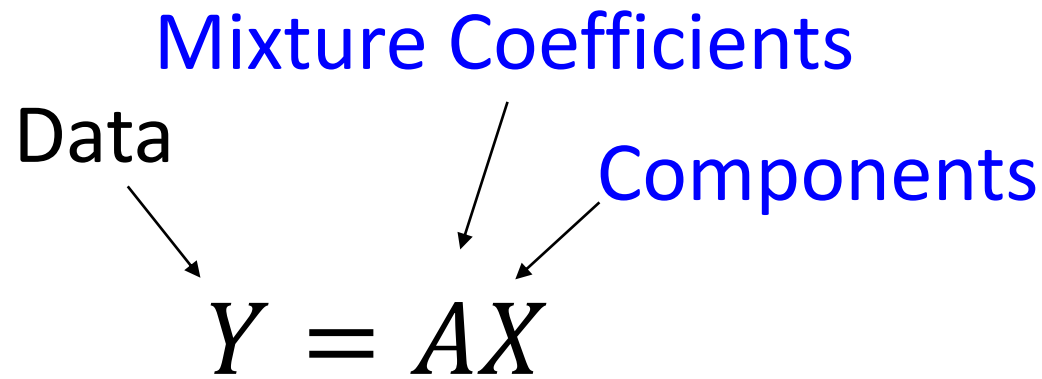
# Independent Component Analysis

---

Mixture Coefficients

Data

Components

$$Y = AX$$


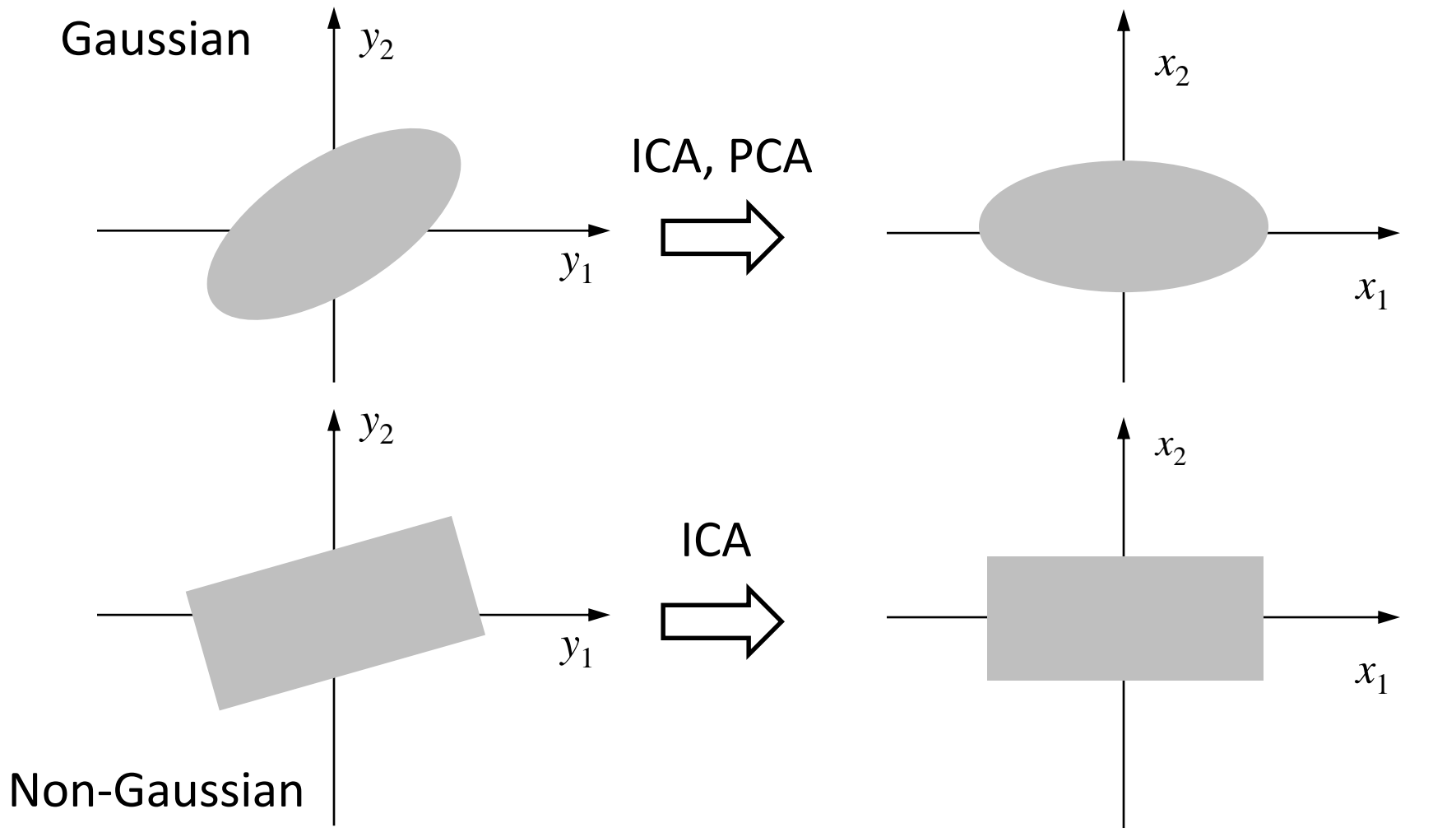
$$\min_A I([X_1, X_2, \dots, X_M] | A, Y)$$

After optimization, the components of  $X$  become as much independent as possible

---

# Illustration of ICA Operation

---



# Summary

---

- Pattern Recognition
  - Model Training and Regularization
  - Probabilities and Gaussian Distributions
  - Bayesian Inferences (ML and MAP)
  - Curse of Dimensionality
  - Decision Theory
  - Entropy and Information
-