

Classification – Bayes optimal classifier

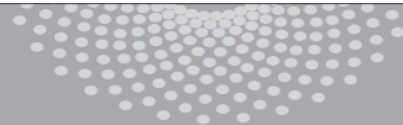
Aarti Singh

Machine Learning 10-315

Aug 28, 2019



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Logistics

- Add yourself to 10-315 on Piazza
- Recitation on Friday – Probability review
- Office hours
 - Mon Siddharth 1-2 pm
 - Tues Yue TBA
 - Wed Aarti 9:30-10:30 am outside classroom
 - Thurs Fabricio 11 am-12 noon
- QnA1 to be released TODAY on Canvas

Performance Measure

For a random test data X , measure of closeness between true label Y and prediction $f(X)$

Binary Classification $\text{loss}(Y, f(X)) = 1_{\{f(X) \neq Y\}}$ **0/1 loss**

Regression $\text{loss}(Y, f(X)) = (f(X) - Y)^2$ **square loss**

- What if overestimating stock price is 10 times more costly than underestimating it?
- What if missing a tumor is 10 times more costly than falsely detecting it?

Performance Measure

For a random test data X , measure of closeness between true label Y and prediction $f(X)$

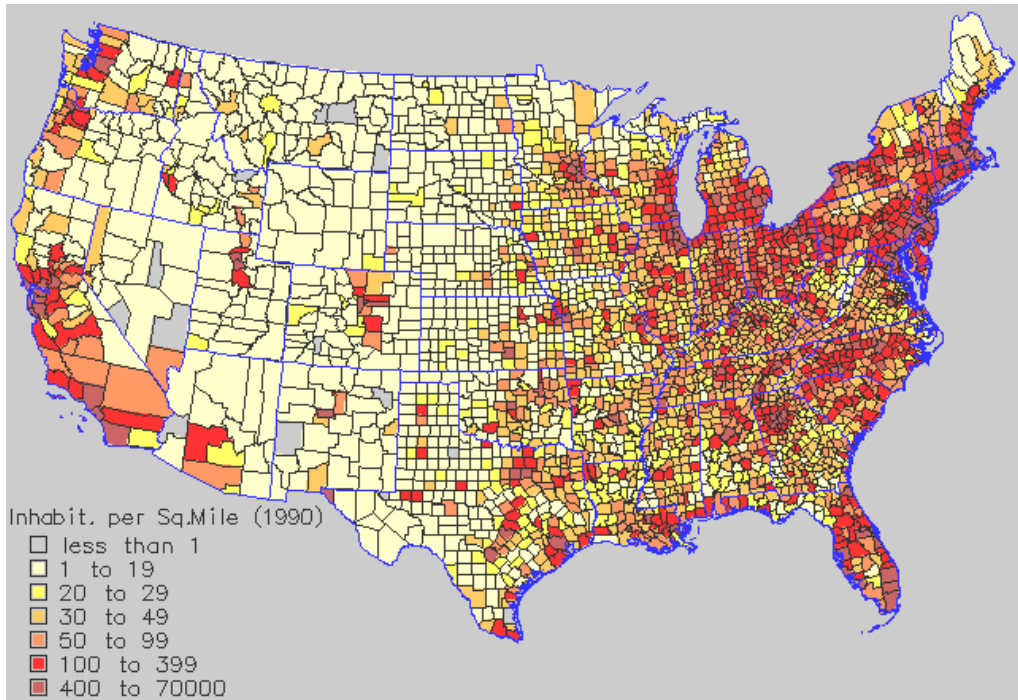
Binary Classification $\text{loss}(Y, f(X)) = 1_{\{f(X) \neq Y\}}$ **0/1 loss**

Regression $\text{loss}(Y, f(X)) = (f(X) - Y)^2$ **square loss**

Density Estimation?

Unsupervised Learning

Density/Distribution Estimation



Population density



Bias of a coin

Performance Measure

For a random test data X , measure of closeness between true label Y and prediction $f(X)$

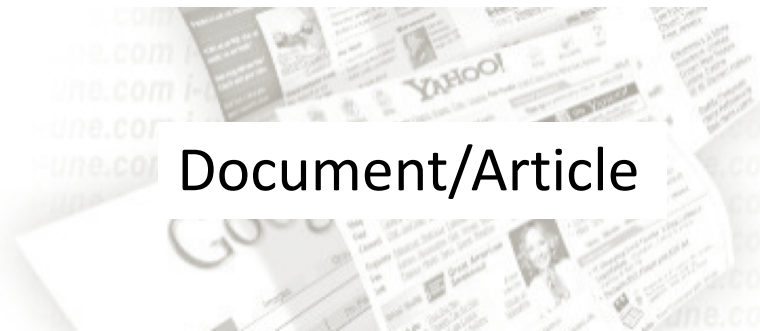
Binary Classification $\text{loss}(Y, f(X)) = \mathbf{1}_{\{f(X) \neq Y\}}$ **0/1 loss**

Regression $\text{loss}(Y, f(X)) = (f(X) - Y)^2$ **square loss**

Density Estimation $\text{loss}(f(X)) = -\log(\mathbb{P}_f(X))$ **Negative log likelihood loss**

Notion of “Features”

Input $X \in \mathcal{X}$



Input $X \in \mathcal{X}$



- How to represent inputs mathematically?
- Document vector X = list of words (different length for each document)
frequency of words (length of each document = size of vocabulary)
- Market information X = daily/monthly? price of share for past 10 years
- Image X = intensity at each pixel, fourier transform values, SIFT etc.

Classification

Goal: Construct prediction rule $f : \mathcal{X} \rightarrow \mathcal{Y}$



Sports
Science
News

Input feature vector, X

Label, Y

In general: label Y can belong to more than two classes

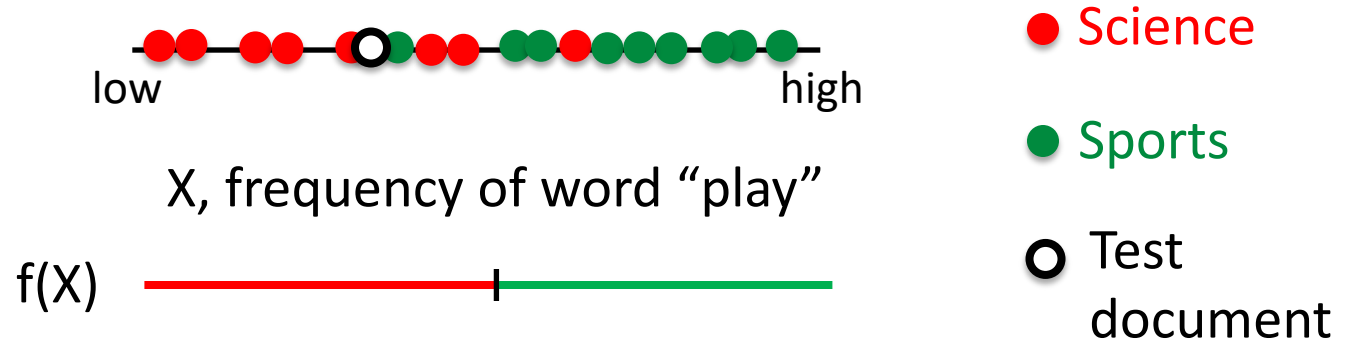
X is multi-dimensional (many features represent an input)

But lets start with a simple case:

label Y is binary (either “Sports” or “Science”)

X is frequency of word “play” = count/total length of document

Binary Classification



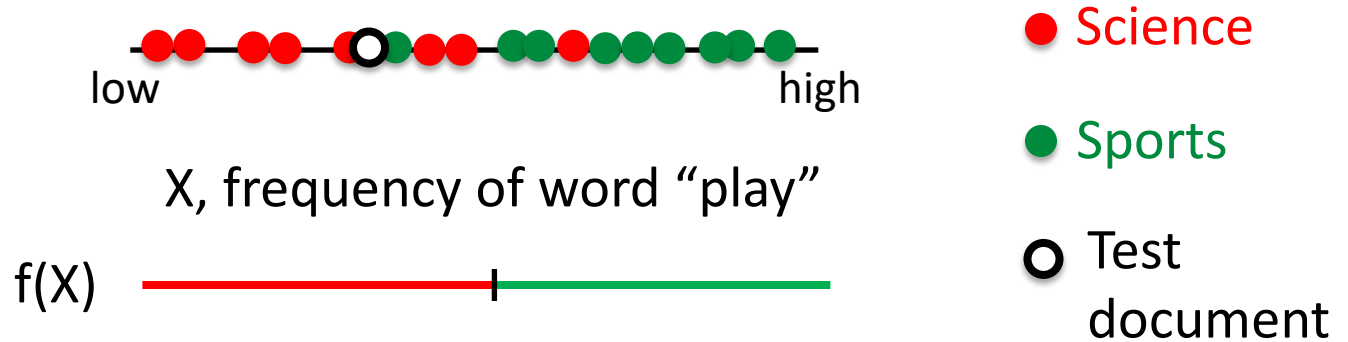
Model X and Y as random variables with joint distribution P_{XY}

Training data $\{X_i, Y_i\}_{i=1}^n \sim \text{iid}$ (independent and identically distributed) samples from P_{XY}

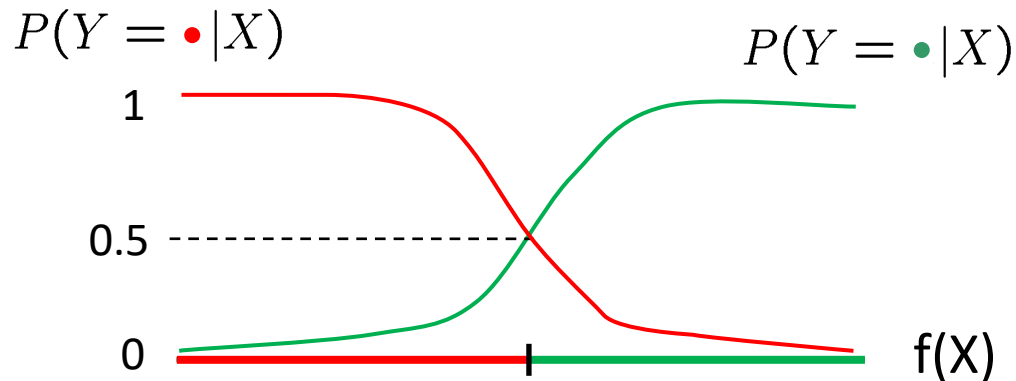
Test data $\{X, Y\} \sim \text{iid}$ sample from P_{XY}

Training and test data are independent draws from same distribution

Binary Classification



Model X and Y as random variables



For a given X, $f(X) = \text{label } Y \text{ which is more likely}$

$$f(X) = \arg \max_{Y=y} P(Y = y | X = x)$$

Optimal Classifier

Optimal classifier: $f^*(x) = \arg \max_{Y=y} P(Y = y|X = x)$

Why??

Goal:

Construct **prediction rule** $f^* : \mathcal{X} \rightarrow \mathcal{Y}$
that minimizes $\text{loss}(Y, f(X))$ for a randomly drawn
test data (X, Y)

$$\min_f \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

$$= \min_f \mathbb{E}_{XY} [\mathbf{1}_{\{f(X) \neq Y\}}]$$

0/1 loss

$$= \min_f \mathbb{P}_{XY}(f(X) \neq Y)$$

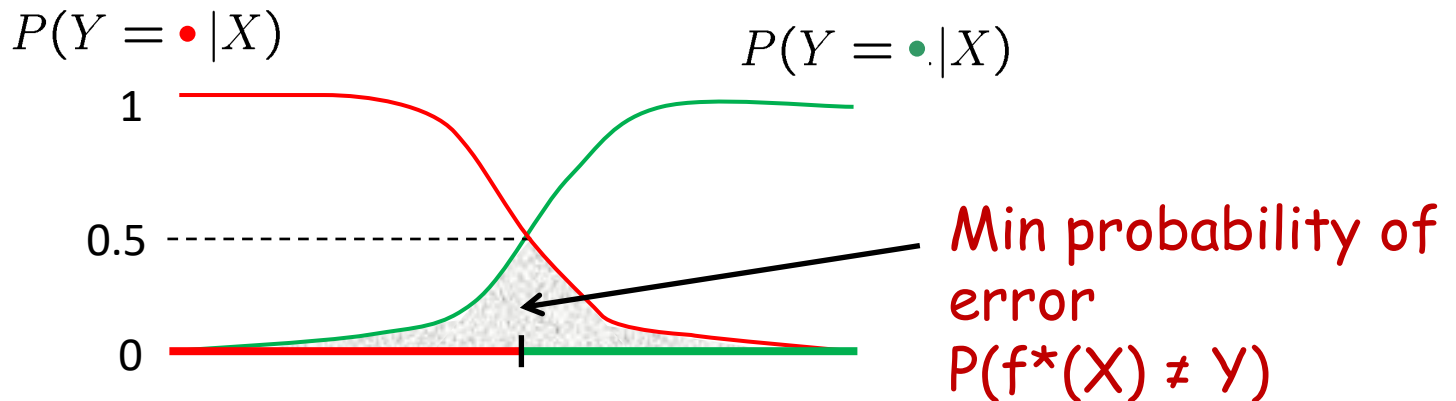
Probability of Error

Minimizer is indeed f^* !!

HW1!

Error of Optimal Classifier

Optimal classifier: $f^*(x) = \arg \max_{Y=y} P(Y = y|X = x)$



- Even the optimal classifier makes mistakes: min probability of error > 0

Bayes Optimal Classifier

Bayes Rule:
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

To see this, recall:

$$P(X,Y) = P(X|Y) P(Y)$$

$$P(Y,X) = P(Y|X) P(X)$$



Thomas Bayes

Bayes Optimal Classifier

Bayes Rule:
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

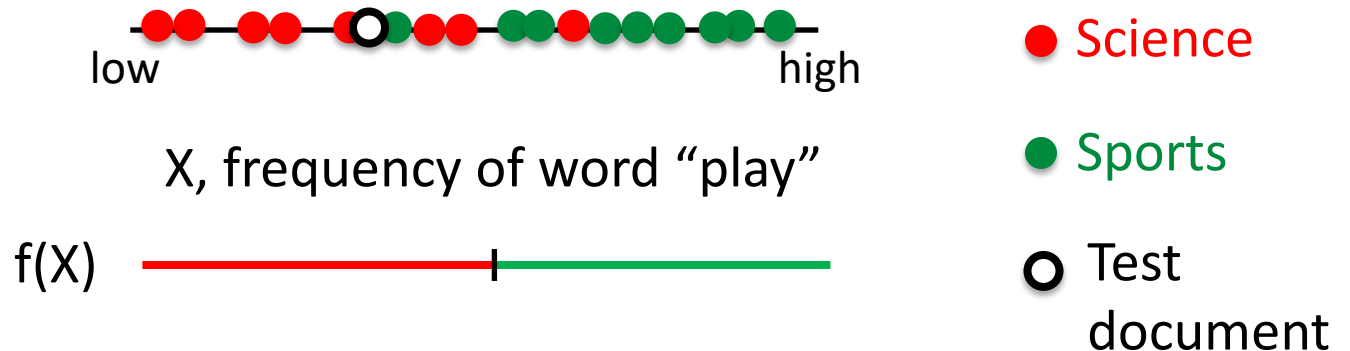
Bayes Optimal classifier:

$$f^*(x) = \arg \max_{Y=y} P(Y = y|X = x)$$

$$= \arg \max_{Y=y} \underbrace{P(X = x|Y = y)}_{\text{Class conditional distribution}} \underbrace{P(Y = y)}_{\text{Class probability}}$$

Class conditional distribution Class probability

Bayes Optimal Classifier



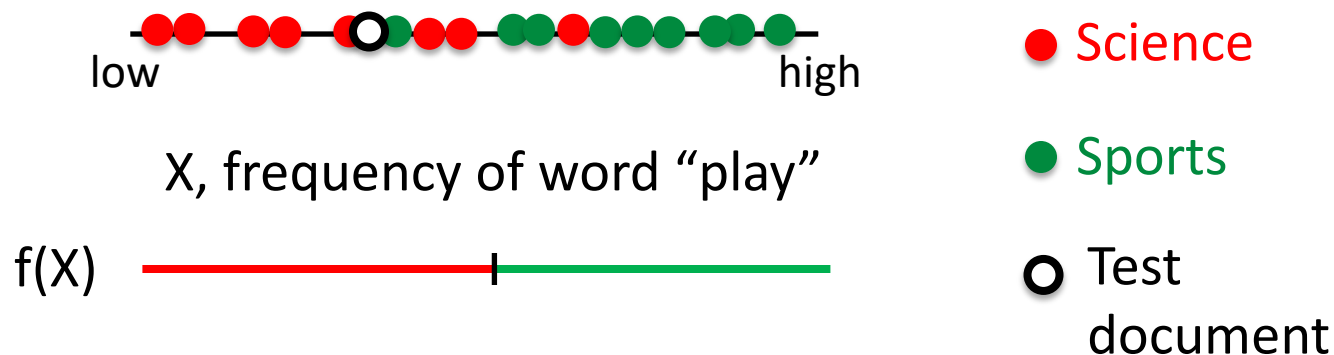
$$f^*(x) = \arg \max_{Y=y} \underbrace{P(X = x | Y = y)}_{\text{Class conditional distribution}} \underbrace{P(Y = y)}_{\text{Class probability}}$$

We can now consider appropriate models for the two terms:

Class probability $P(Y=y)$

Class conditional distribution of features $P(X=x | Y=y)$

Modeling class probability



Modeling Class probability $P(Y=y) = \text{Bernoulli}(\theta)$

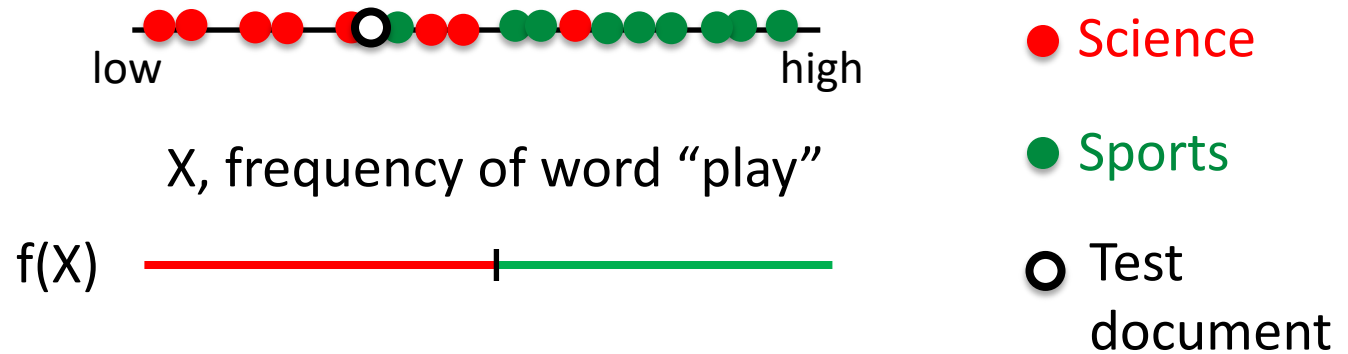
$$P(Y = \bullet) = \theta$$

$$P(Y = \bullet) = 1 - \theta$$

Like a coin flip

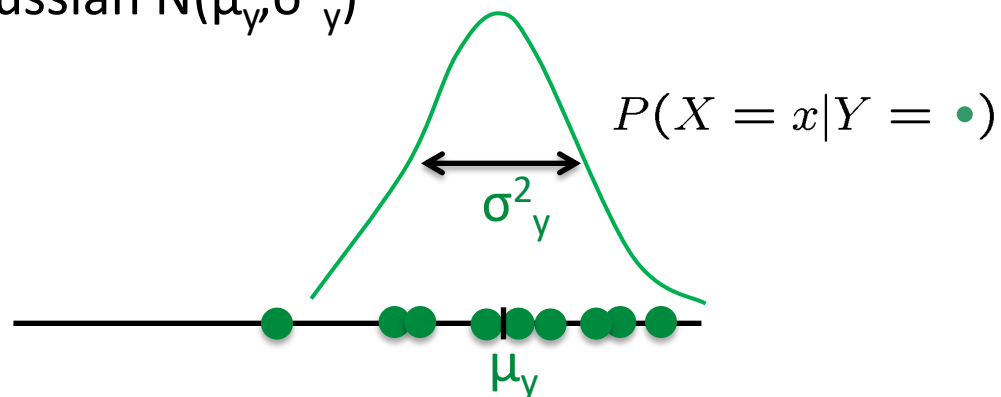


Modeling class conditional distribution of features



Modeling Class Conditional distribution of features $P(X=x|Y=y)$

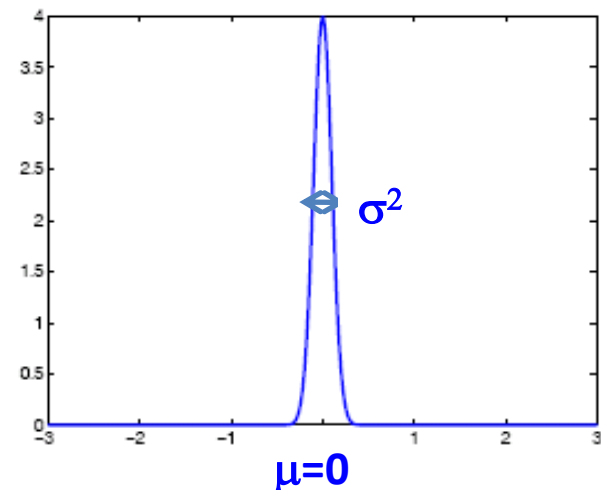
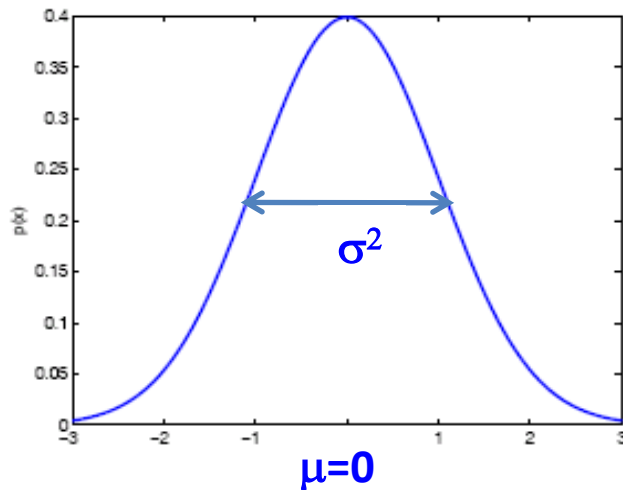
E.g. $P(X=x|Y=y) = \text{Gaussian } N(\mu_y, \sigma_y^2)$



1-dim Gaussian distribution

X is Gaussian $N(\mu, \sigma^2)$

$$P(X = x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

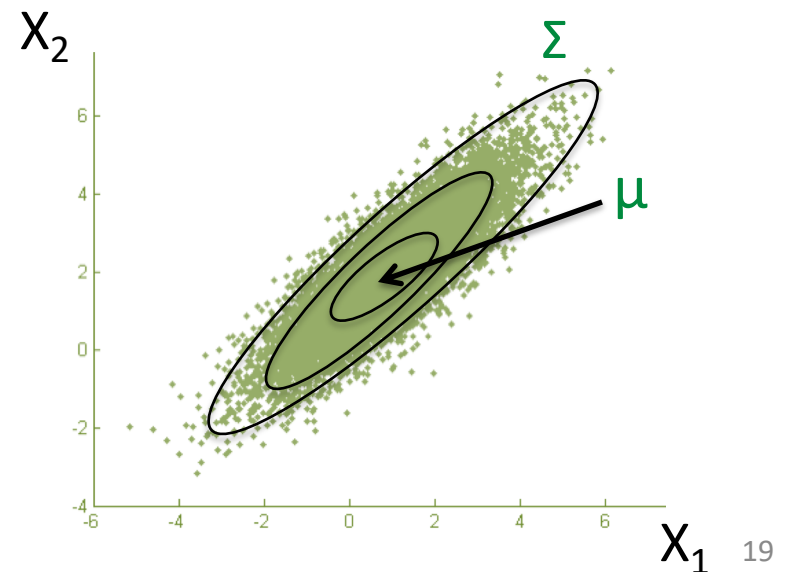
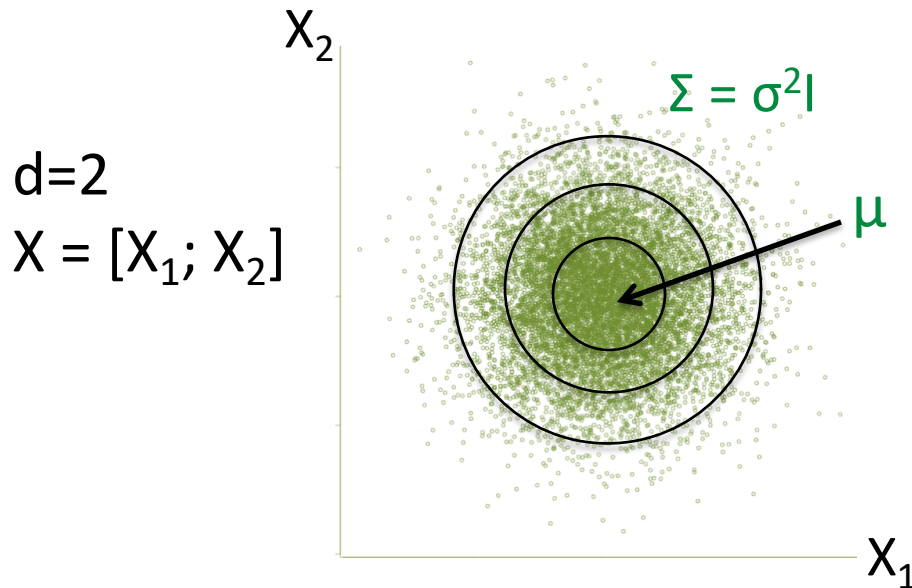


d-dim Gaussian distribution

X is Gaussian $N(\mu, \Sigma)$

μ is d-dim vector, Σ is dxd dim matrix

$$P(X = x | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$



Gaussian Bayes classifier

$$f^*(x) = \arg \max_{Y=y} \underbrace{P(X = x|Y = y)}_{\text{Class conditional density}} \underbrace{P(Y = y)}_{\text{Class probability}}$$

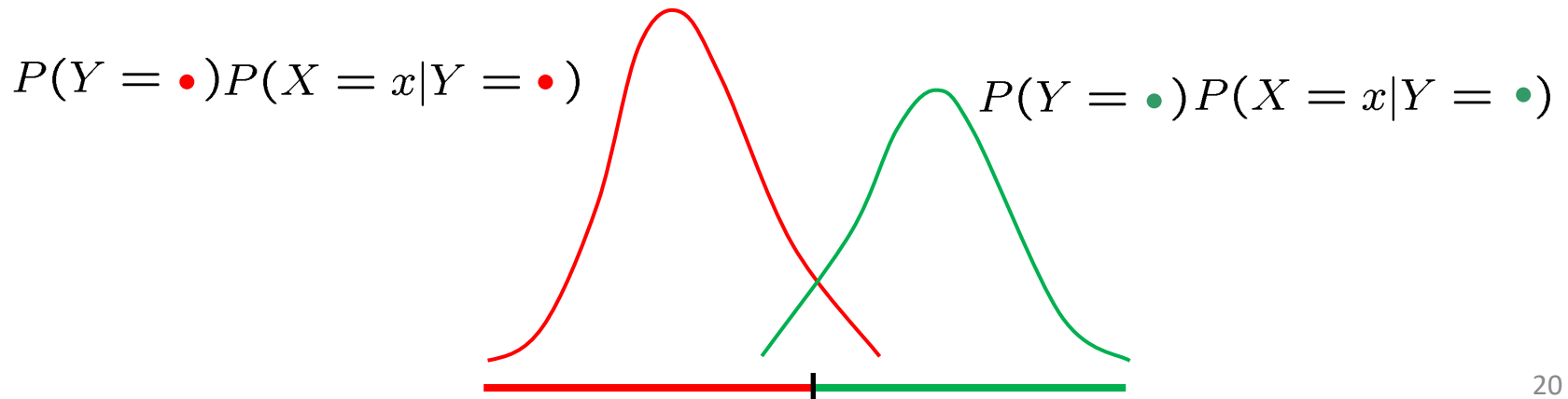
How to learn parameters
 θ, μ_y, Σ_y from data?

Class conditional
density

Class probability

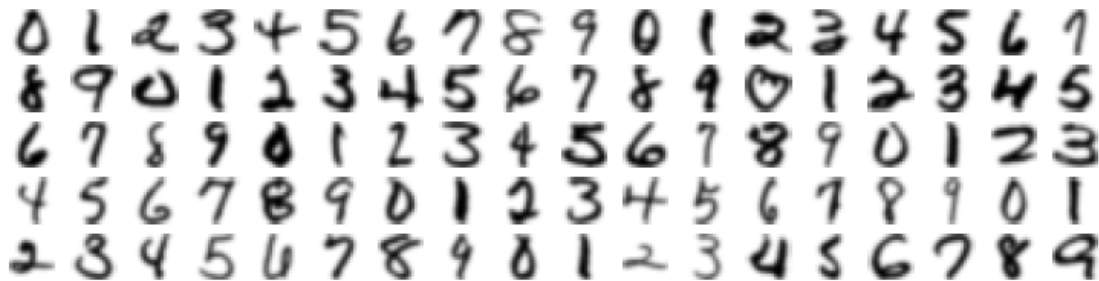
Gaussian(μ_y, Σ_y)

Bernoulli(θ)

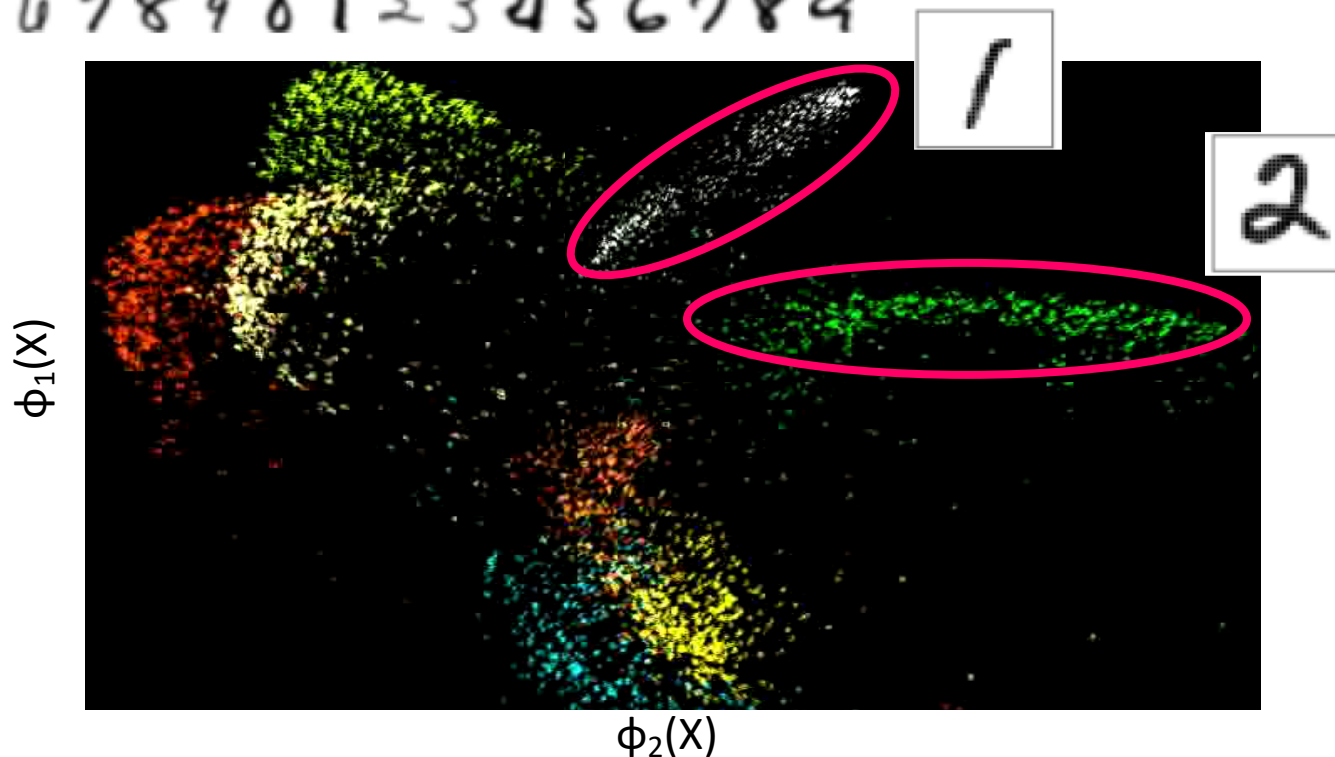


Multi-class problem
Multi-dimensional input X

Handwritten digit recognition



Multi-class
classification



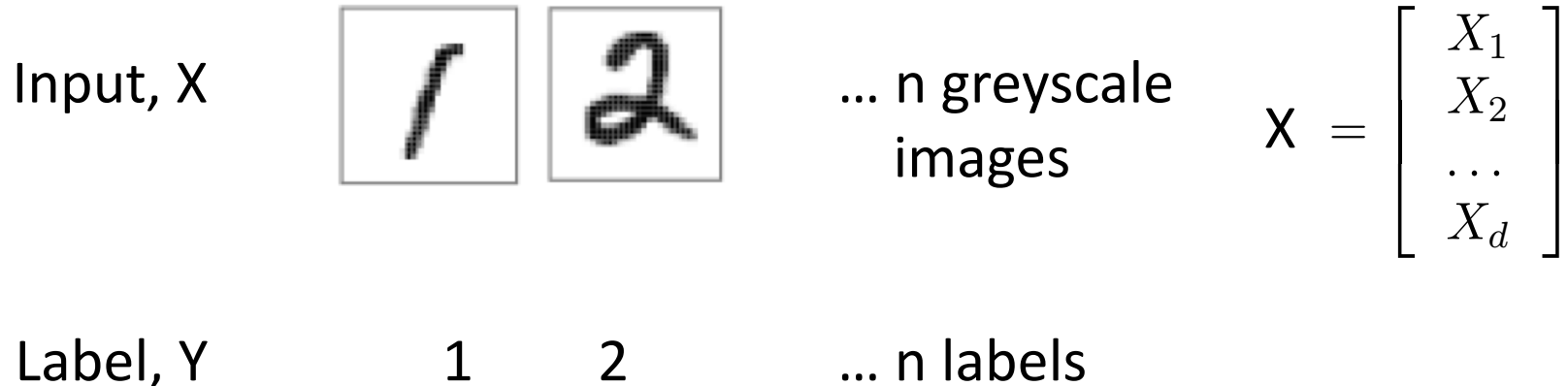
Note: 8 digits shown out of 10 (0, 1, ..., 9);

Axes are obtained by nonlinear dimensionality reduction (later in course)

Handwritten digit recognition

Training Data:

Each image represented as a vector of **intensity values** at the **d pixels (features)**



Gaussian Bayes model:

$P(Y = y) = p_y$ for all y in 0, 1, 2, ..., 9

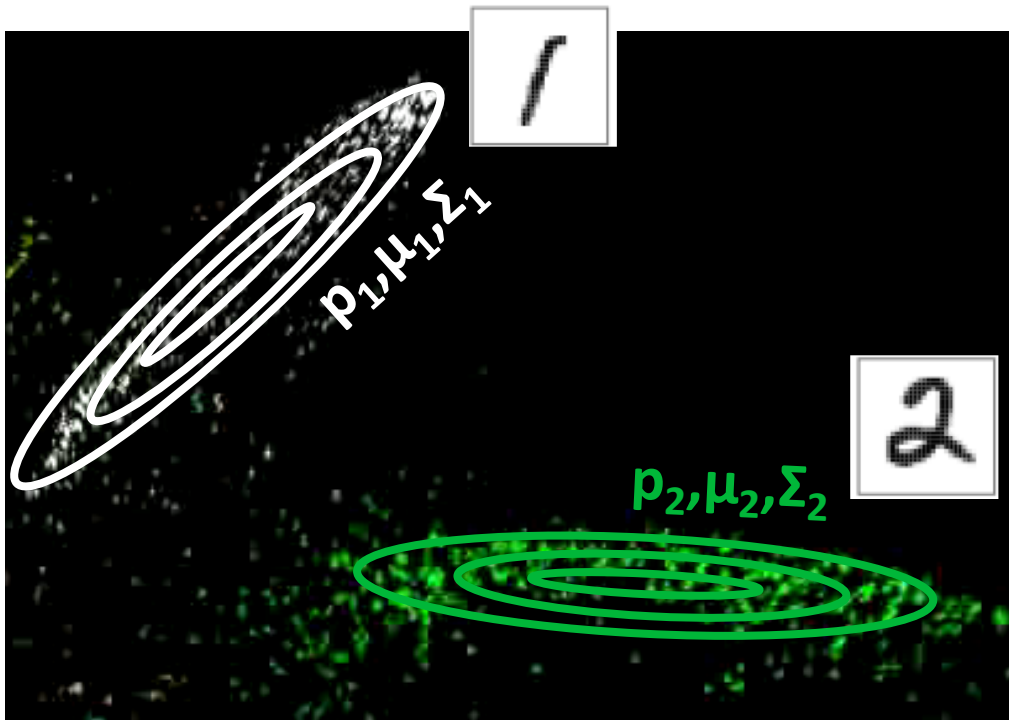
p_0, p_1, \dots, p_9 (sum to 1)

$P(X=x | Y = y) \sim N(\mu_y, \Sigma_y)$ for each y

μ_y - d-dim vector

Σ_y - dxd matrix

Gaussian Bayes classifier



How to learn parameters p_y, μ_y, Σ_y from data?

$P(Y = y) = p_y$ for all y in $0, 1, 2, \dots, 9$

$P(X=x | Y = y) \sim N(\mu_y, \Sigma_y)$ for each y

p_0, p_1, \dots, p_9 (sum to 1)

μ_y – d-dim vector

Σ_y - dxd matrix

How many parameters do we need to learn?

Class probability:

$$P(Y = y) = p_y \text{ for all } y \text{ in } 0, 1, 2, \dots, 9 \quad p_0, p_1, \dots, p_9 \text{ (sum to 1)}$$

K-1 if K labels

Class conditional distribution of features:

$$P(X=x | Y = y) \sim N(\mu_y, \Sigma_y) \text{ for each } y \quad \begin{array}{l} \mu_y - d\text{-dim vector} \\ \Sigma_y - d \times d \text{ matrix} \end{array}$$

$Kd + Kd(d+1)/2 = O(Kd^2)$ if d features

Quadratic in dimension d ! If $d = 256 \times 256$ pixels, ~ 21.5 billion parameters!