

# Bayes and Naïve Bayes Classifier

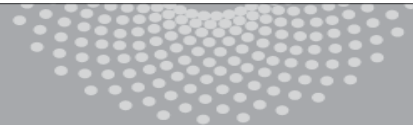
Aarti Singh

Machine Learning 10-701

Jan 23, 2023



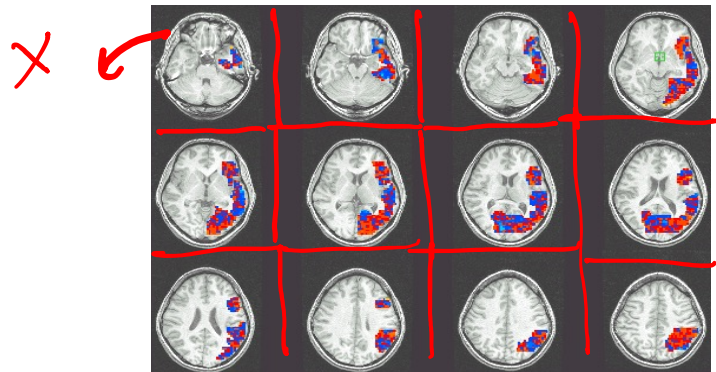
**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# Classification

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



High Stress  
Moderate Stress  
Low Stress

} Stressed  
] No stress

Input feature vector, X

Label, Y

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

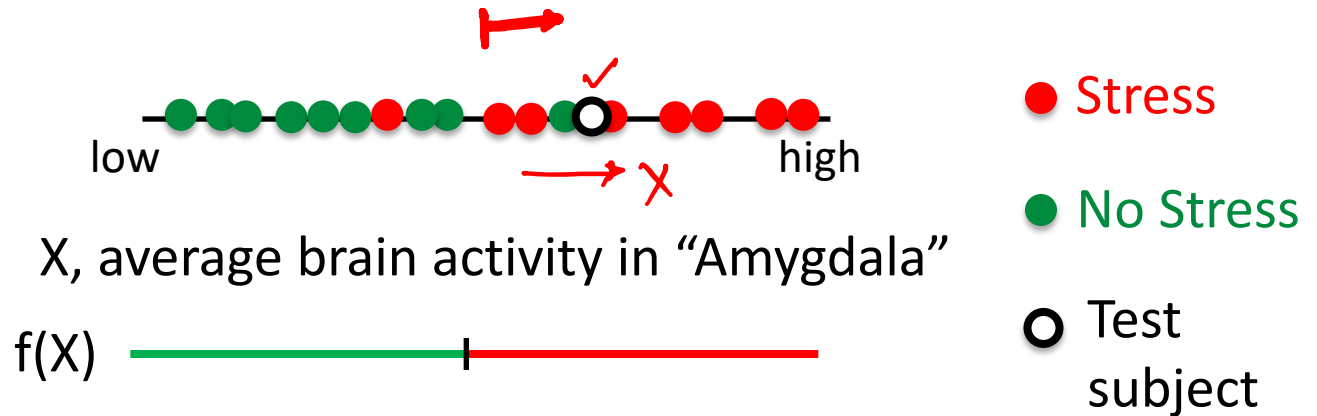
X is multi-dimensional (many features represent an input)

But lets start with a simple case:

label Y is binary (either “Stress” or “No Stress”)

X is average brain activity in the “Amygdala”

# 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

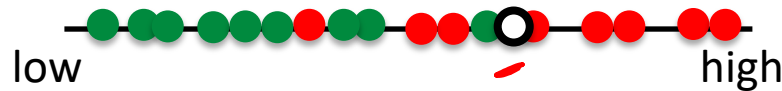
# Optimal classifier

Minimize loss in expectation (over random test data)

$$\min_f E_{XY}[\text{loss}(f(X), Y)] \quad (X, Y)$$

- Which classifier  $f$  is optimal for 0/1 loss, assuming we know data-generating distribution  $P(X, Y)$ ?

# Optimal Classifier

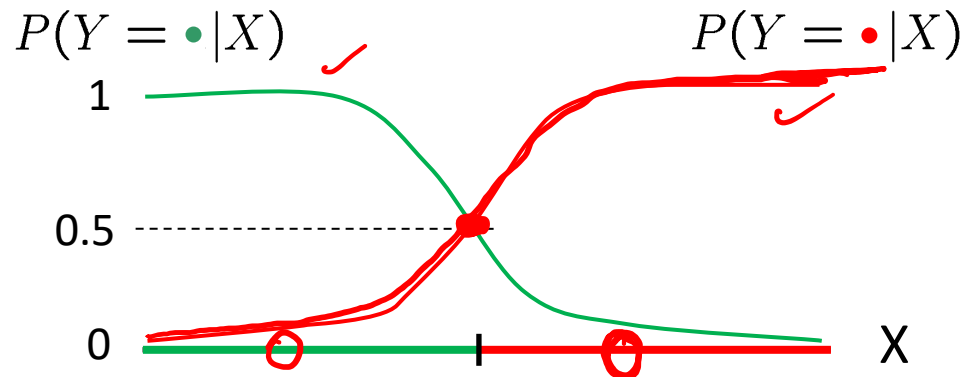


X, average brain activity in "Amygdala"



- Stress
- No Stress
- Test subject

Model X and Y as random variables



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

$$f(X) = \arg \max_y P(Y = y | X = x)$$

# Optimal classifier

$$E[\mathbb{1}_A] = P(A)$$

Minimize loss in expectation (over random test data)

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

$$\text{loss}(f(x), y) = \mathbb{1}_{f(x) \neq y}$$

- Which classifier  $f$  is optimal for 0/1 loss, assuming we know data-generating distribution  $P(X, Y)$ ?

$$\begin{aligned} P(f(x) \neq Y) &= \int P(f(x) \neq Y | x) P(x) dx \\ &= \int_{x: f(x)=1} P(Y=0|x) p(x) dx + \int_{x: f(x)=0} P(Y=1|x) p(x) dx \end{aligned}$$

$$\Rightarrow f(x) = \underset{y}{\text{argmax}} P(Y=y|x)$$

where  $f$  is  $\underset{y}{\text{argmin}} E[\text{loss}] = P(f(x) \neq Y)$

# Bayes Rule

**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:

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

||

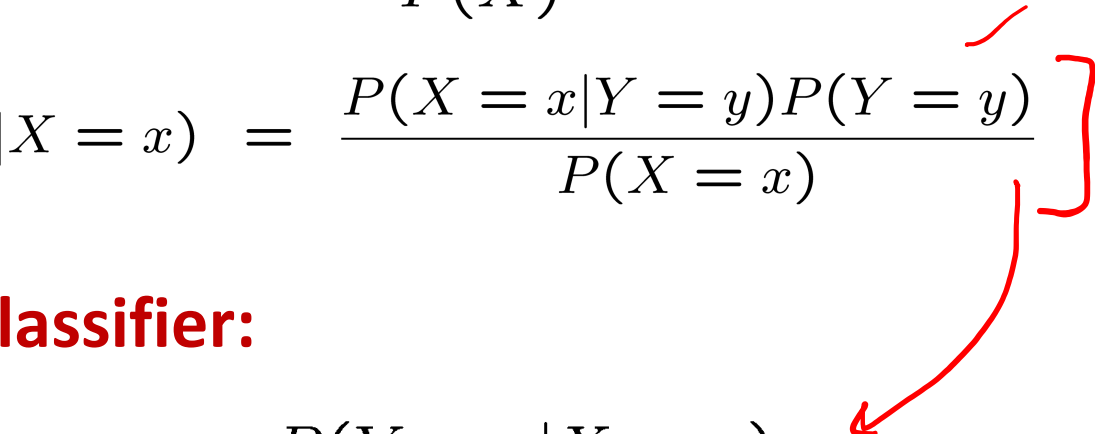
$$P(Y,X) = P(Y|X) P(X)$$



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:**

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

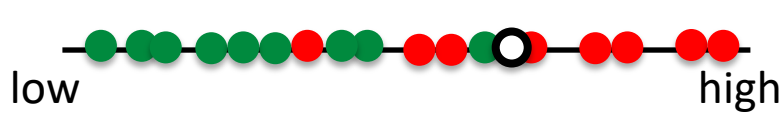
$$= \arg \max_{Y=y} \underbrace{P(X = x|Y = y)}_{\text{Class conditional Distribution of features}} \underbrace{P(Y = y)}_{\text{Distribution of class}}$$

Class conditional  
Distribution of features

Distribution of class



# Bayes Classifier



- Stress
- No Stress
- Test subject

X, average brain activity in “Amygdala”



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

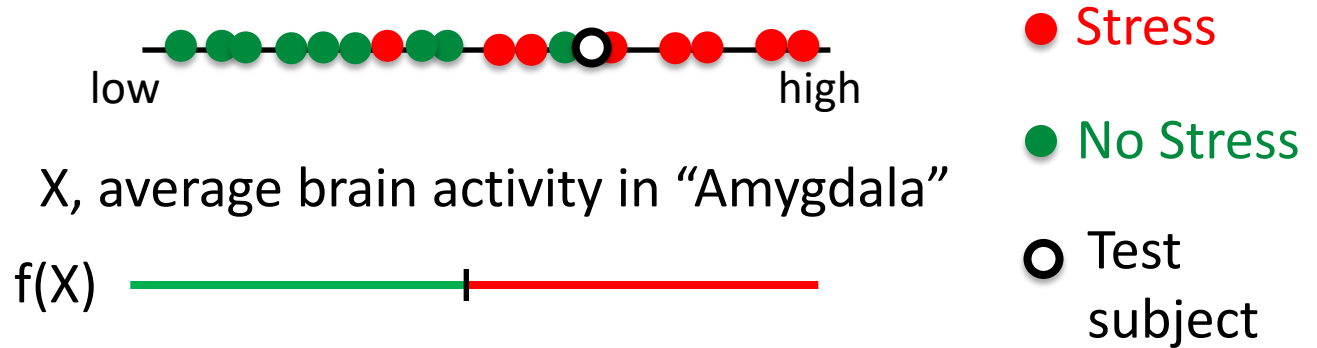
Class conditional  
Distribution of features

Class distribution

We can now consider distribution models to approximate ground truth:

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

# Modeling class distribution



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

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

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

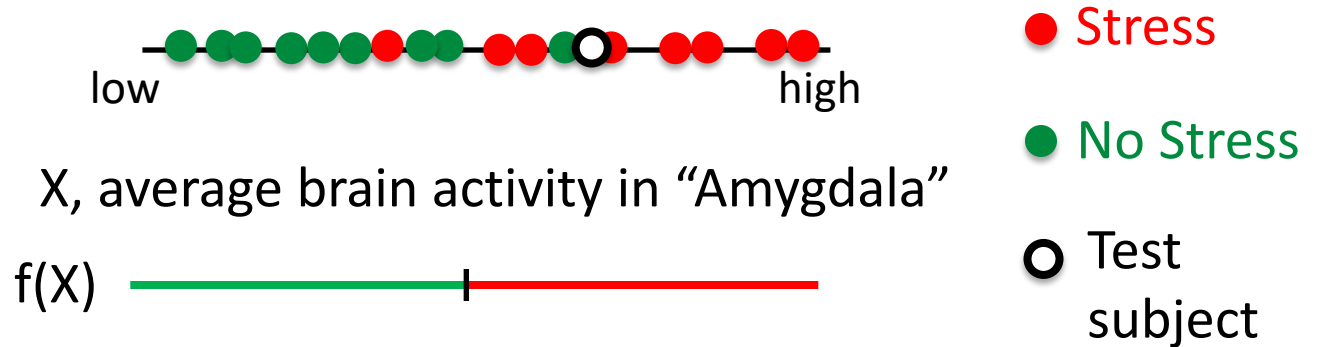
Like a coin flip



$P(Y=1)$   
 $P(Y=2)$   
 $\vdots$   
 $P(Y=6)$   
 $= 1 - P(Y=1) - P(Y=2) - \dots - P(Y=5)$

➤ How do we model multiple (>2) classes?

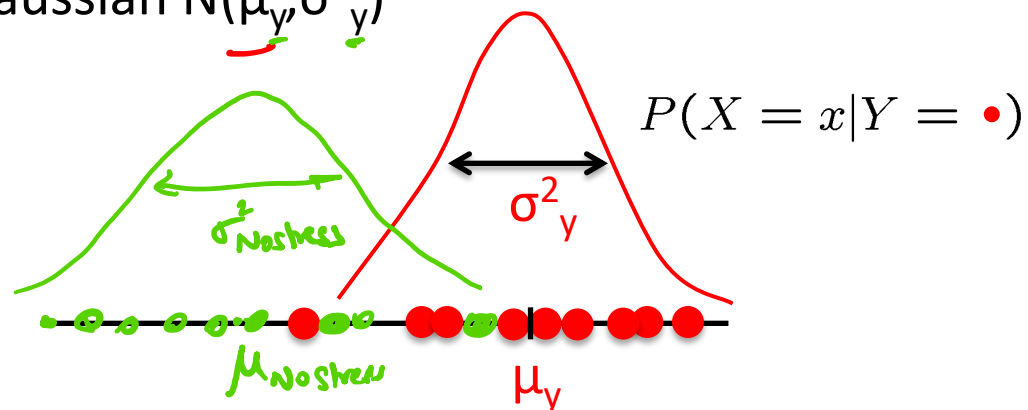
# Modeling class conditional distribution of features



Modeling class conditional distribution of feature  $P(X=x|Y=y)$

➤ What distribution would you use?

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



# Gaussian Bayes classifier

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

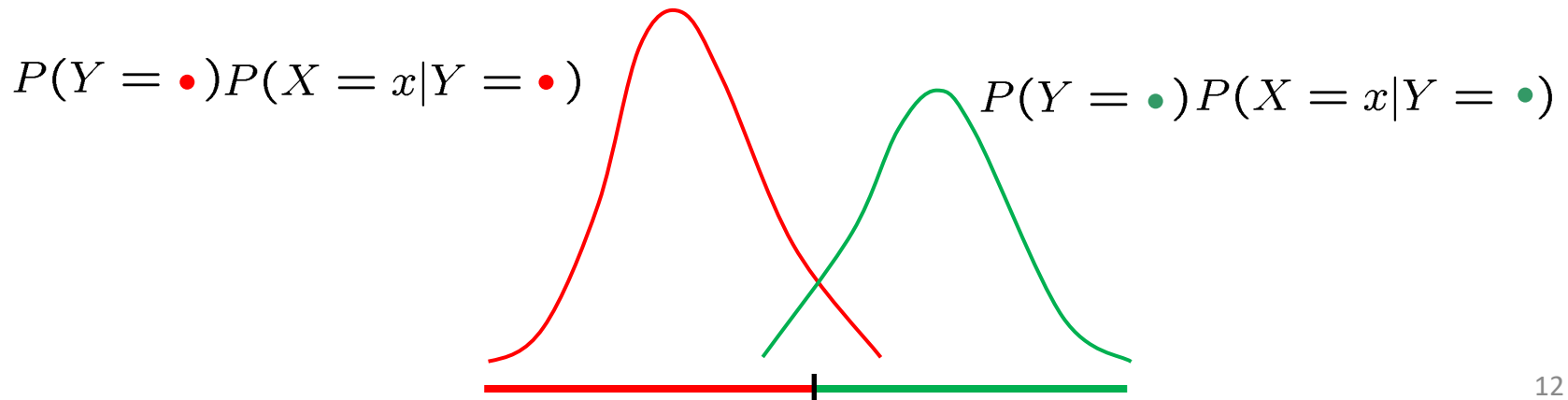
Use MLE/MAP to learn parameters  $\theta$ ,  $\mu_y$ ,  $\Sigma_y$  from data

Class conditional  
Distribution of features

Class distribution

Gaussian( $\mu_y$ ,  $\Sigma_y$ )

Bernoulli( $\theta$ )

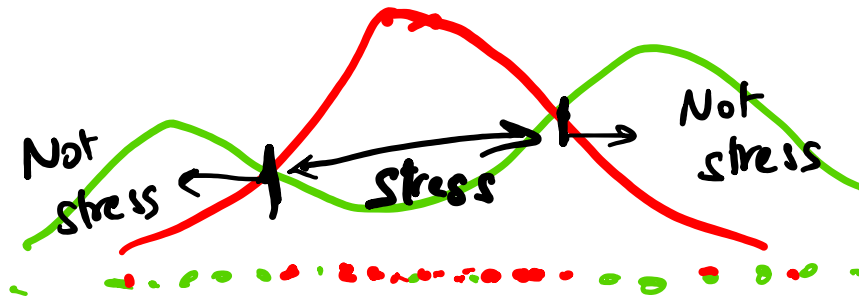


# Poll

- Is the Gaussian Bayes Classifier always optimal under 0/1 loss?

A. True

B. False



# 1-dim Gaussian Bayes classifier

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

Class conditional  
Distribution of features

Class distribution

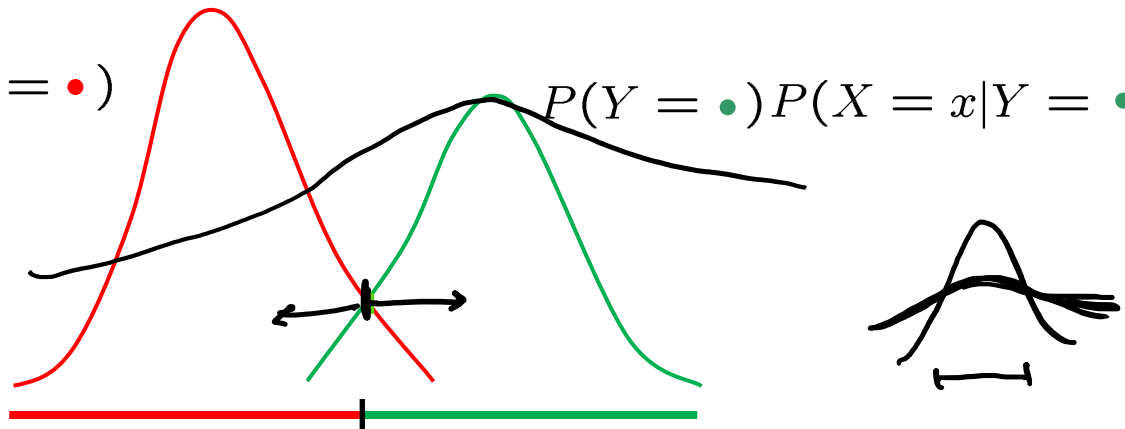
➤ What decision boundaries can we get in 1-dim?

Gaussian( $\mu_y, \sigma_y^2$ )

Bernoulli( $\theta$ )

$$P(Y = \bullet)P(X = x|Y = \bullet)$$

$$P(Y = \bullet)P(X = x|Y = \bullet)$$



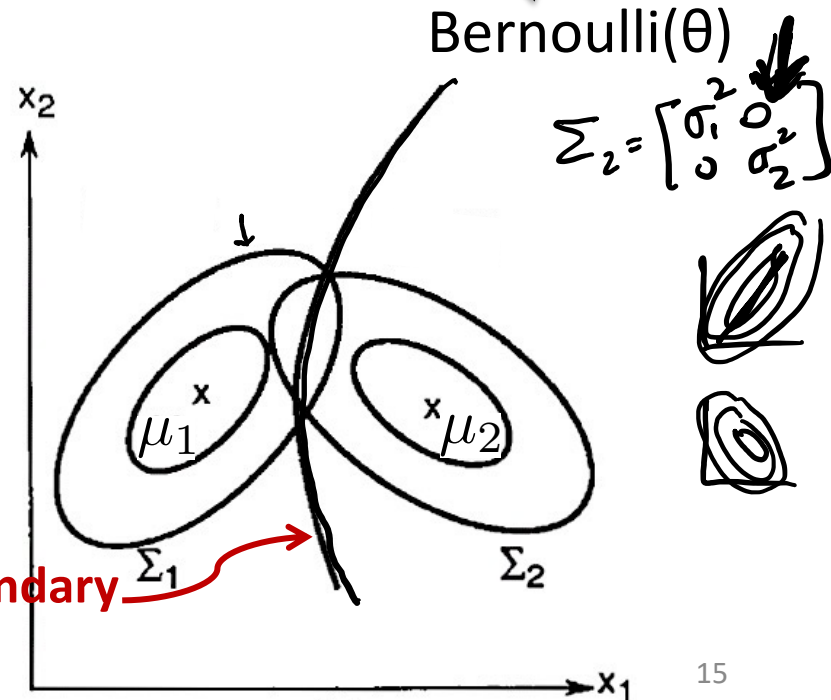
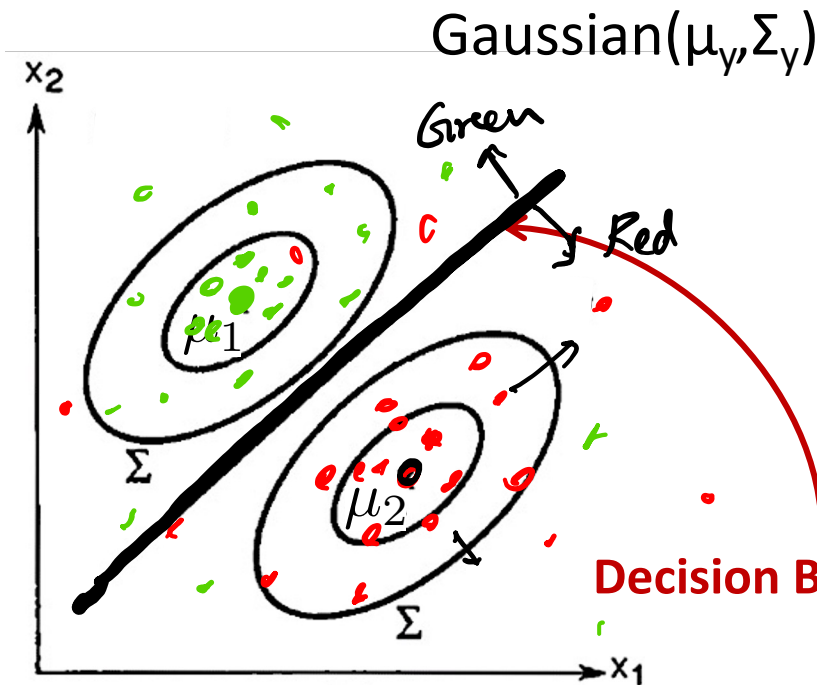
# d-dim Gaussian Bayes classifier

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

➤ What decision boundaries can we get in d-dim?

Class conditional  
Distribution of features

Class distribution



# Decision Boundary of Gaussian Bayes

- Decision boundary is set of points  $x$ :  $P(Y=1|X=x) = P(Y=0|X=x)$

Compute the ratio

$$1 = \frac{P(Y=1|X=x)}{P(Y=0|X=x)} = \frac{P(X=x|Y=1)P(Y=1)}{P(X=x|Y=0)P(Y=0)}$$

$\mathcal{N}(\mu_1, \Sigma_1)$   
 $\downarrow$   
 $\theta$   
 $\mathcal{N}(\mu_0, \Sigma_0)$   
 $\downarrow$   
 $1-\theta$

= quadratic eq<sup>n</sup> in  $x$  if  $\Sigma_1 \neq \Sigma_0$   
 $\exp(-x^T(\Sigma_1^{-1} - \Sigma_0^{-1})x + \dots)$   
 $0$  if  $\Sigma_1 = \Sigma_0$

In general, this implies a quadratic equation in  $x$ . But if  $\Sigma_1 = \Sigma_0$ , then quadratic part cancels out and decision boundary is linear.



# Recap

- **Bayes classifier** – assumes  $P_{XY}$  known, optimal for 0/1 loss

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

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

Class conditional

Class distribution

Distribution of features

- **Gaussian Bayes classifier** – assumes  
Class distribution is Bernoulli/Multinomial  
Class conditional distribution of features is Gaussian
- **Decision boundary** – (binary classification)

# How many parameters do we need to learn (continuous features)?

Class distribution:

$P(Y = y) = p_y$  for all  $y$  in  $H, M, L$

$p_H, p_M, p_L$  (sum to 1)

**K-1 if K labels**

$$p_L, p_M, p_H = 1 - p_L - p_M$$

Class conditional distribution of features:

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

*d features / attributes*

$\mu_y$  - d-dim vector

$\Sigma_y$  -  $d \times d$  matrix

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

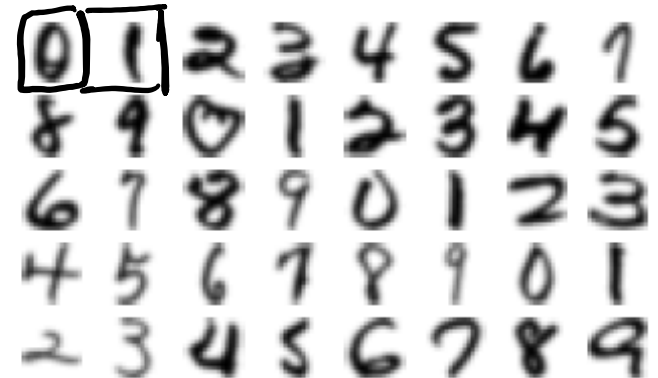
**Quadratic in dimension  $d$ ! If  $d = 256 \times 256$  pixels,  $\sim 13$  billion parameters!**

# How many parameters do we need to learn (discrete features)?

Class distribution:

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

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



**K-1 if K labels**

Class conditional distribution of (binary) features:

$P(X=x | Y = y) \sim$  For each label  $y$ , maintain probability table with  $2^d - 1$  entries

**$K(2^d - 1)$  if  $d$  binary features**

**Exponential in dimension  $d$ !**

# What's wrong with too many parameters?

- How many training data needed to learn one parameter (bias of a coin)?



- Need lots of training data to learn the parameters!
  - Training data  $>$  number of (independent) parameters

# Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:
  - Features are independent given class:

$$\begin{aligned} P(\underbrace{X^{(1)}}_{\text{feature}}, \underbrace{X^{(2)}}_{\text{feature}} | Y) &= P(X^{(1)} | X^{(2)}, Y) P(X^{(2)} | Y) \checkmark \\ &= \underbrace{P(X^{(1)} | Y)} \underbrace{P(X^{(2)} | Y)} \end{aligned} \quad X = \begin{bmatrix} X^{(1)} \\ X^{(2)} \end{bmatrix}$$

- More generally:

$$P(X^{(1)}, \dots, X^{(d)} | Y) = \prod_{i=1}^d P(X^{(i)} | Y) \quad X = \begin{bmatrix} X^{(1)} \\ \vdots \\ X^{(d)} \end{bmatrix}$$

- If conditional independence assumption holds, NB is optimal classifier! But worse otherwise.

# Conditional Independence

- X is **conditionally independent** of Y given Z:

probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall x, y, z) P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

- Equivalent to:

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

- e.g.,  $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

**Note:** does NOT mean Thunder is independent of Rain

# Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:
  - Features are independent given class:

$$\rightarrow \underbrace{P(X^{(1)}, \dots, X^{(d)} | Y)} = \prod_{i=1}^d P(X^{(i)} | Y)$$

$$X = \begin{bmatrix} X^{(1)} \\ \vdots \\ X^{(d)} \end{bmatrix}$$

$$\begin{aligned} f_{NB}(\mathbf{x}) &= \arg \max_y \underbrace{P(x^{(1)}, \dots, x^{(d)} | y)} \underbrace{P(y)} \\ &= \arg \max_y \underbrace{\prod_{i=1}^d P(x^{(i)} | y)} \underbrace{P(y)} \end{aligned}$$

Naïve assumption

- How many parameters now?

# How many parameters do we need to learn (continuous features)?

## ➤ Poll

Number of parameters for class distribution  $P(Y=y)$  for  $K$  classes?  $K-1$

Number of parameters for Class conditional distribution of features  $P(X = x|Y = y)$  for  $d$  features (using Gaussian Naïve Bayes assumption)?

A.  $K-1, Kd$  ✓

B.  $K-1, K(d + d(d+1)/2)$

C.  $K-1, Kd^2$

D.  $K-1, 2Kd$

$$p(X=x|Y=y) \sim \mathcal{N}(\mu_y, \Sigma_y)$$

$$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{bmatrix}$$

$$\Sigma_y = \begin{bmatrix} \circ & \circ \\ \circ & \circ \end{bmatrix}$$

$Kd$  vs.  $Kd^2$

↑ w/o Naive Bayes



# How many parameters do we need to learn (continuous features)?

Class probability:

$$P(Y = y) = p_y \text{ for all } y \text{ in } H, M, L \quad p_H, p_M, p_L \text{ (sum to 1)}$$

**K-1 if K labels**

Class conditional distribution of features (using Naïve Bayes assumption):

$$P(X^{(i)} = x^{(i)} | Y = y) \sim N(\mu_y^{(i)}, \sigma_y^{2(i)}) \text{ for each } y \text{ and each pixel } i$$

**2Kd if d features**

**Linear instead of Quadratic in dimension d!**

# How many parameters do we need to learn (discrete features)?

## ➤ Poll

Number of parameters for class distribution  $P(Y=y)$  for  $K$  classes?  $K-1$

Number of parameters for Class conditional distribution of features  $P(X = x | Y = y)$  for  $d$  binary features (using Naïve Bayes assumption)?

A.  $K-1, K2^d$

B.  $K-1, K(d-1)$

C.  $K-1, Kd$  ✓

D.  $K-1, 2Kd$

$$P(X_1=x_1 \dots X_d=x_d | Y=y)$$
$$= \prod_{i=1}^d P(X_i=x_i | Y=y)$$

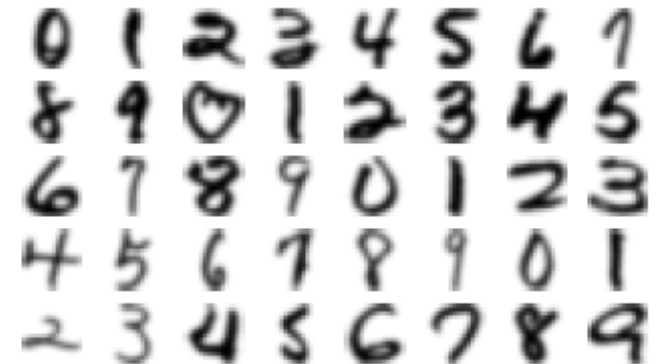
$Kd$   
vs.  $K2^d$  w/o Naive Bayes

# How many parameters do we need to learn (discrete features)?

Class probability:

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

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



**K-1 if K labels**

Class conditional distribution of (binary) features:

$P(X^{(i)} = x^{(i)} | Y = y)$  – one probability value for each  $y$ , pixel  $i$

**Kd if d binary features**

**Linear instead of Exponential in dimension d!**

# Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:
  - Features are independent given class:

$$\underbrace{P(X^{(1)}, \dots, X^{(d)} | Y)} = \prod_{i=1}^d \underbrace{P(X^{(i)} | Y)}$$

$$\begin{aligned} f_{NB}(\mathbf{x}) &= \arg \max_y P(x^{(1)}, \dots, x^{(d)} | y) P(y) \\ &= \arg \max_y \prod_{i=1}^d P(x^{(i)} | y) P(y) \end{aligned}$$

- Has fewer parameters, and hence requires fewer training data, even though assumption may be violated in practice

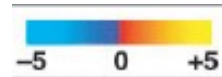
# Learned Gaussian Naïve Bayes Model

## Means for P(BrainActivity | WordCategory)

$$\mathcal{N}(\underline{\mu}_p, \underline{\Sigma}_p)$$

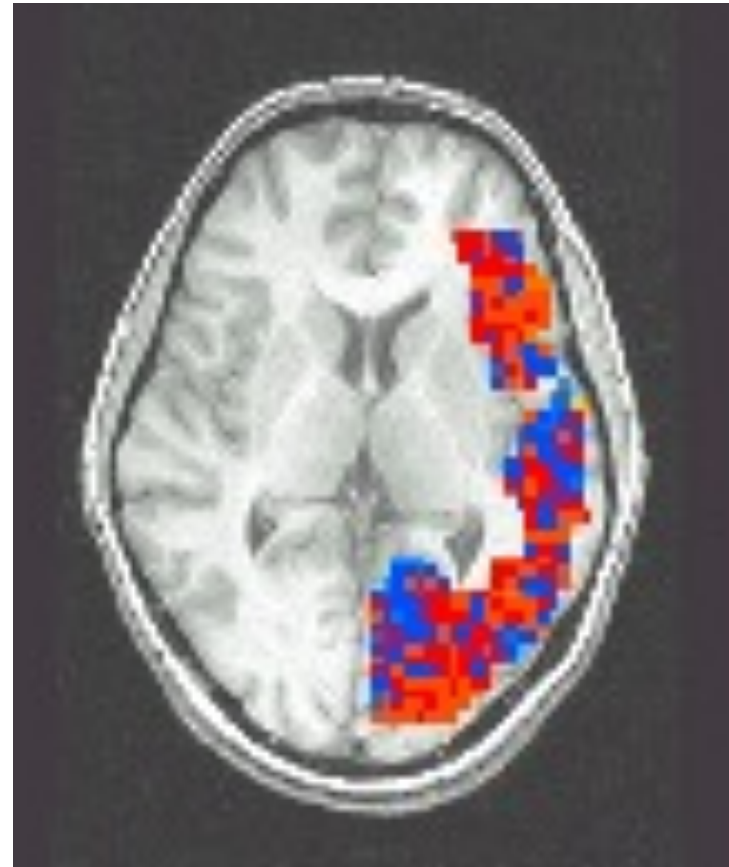
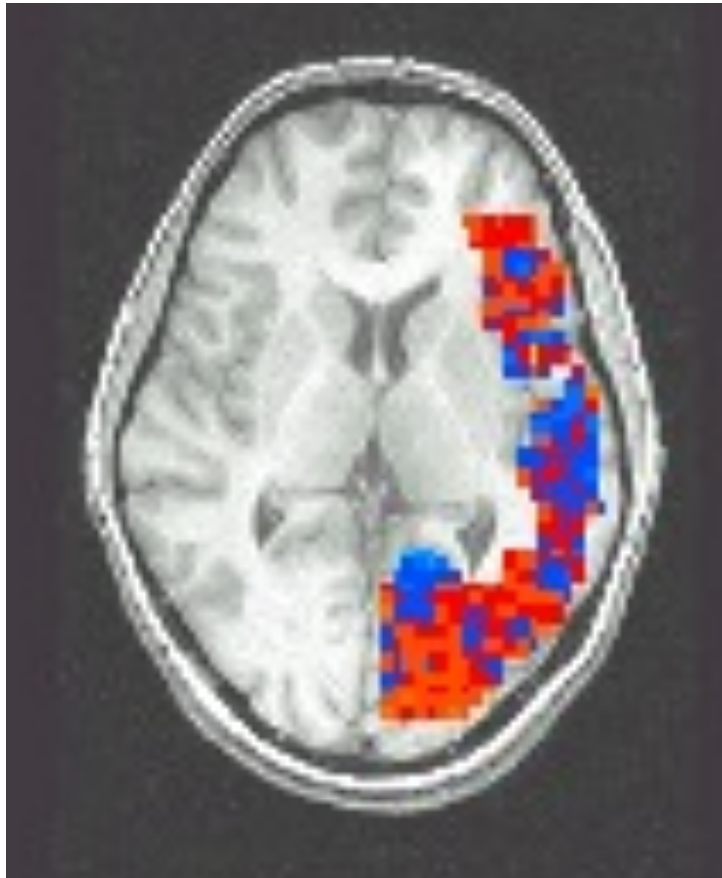
Pairwise classification accuracy: **85%** [Mitchell et al.03]

People words



Animal words

$$\mathcal{N}(\underline{\mu}_a, \underline{\Sigma}_a)$$



# Text classification

Input  $X \in \mathcal{X}$

Document/Article

remember to wake up when class ends  
=  
wake ends to class remember up when

## How to represent inputs mathematically?

- Document vector  $X$  ➤ Ideas?
  - list of words (different length for each document)
  - frequency of words (length of each document = size of vocabulary), also known as **Bag-of-words** approach ➤ Why might this be limited?
    - Misses out context!!
    - list of n-grams (n-tuples of words)

# Text classification

Raw input



Features



Model for input features



<i>basketball</i>	word1	5
<i>play</i>	word2	2
<i>refree</i>	word3	10
	word4	20
	word5	12
	word6	5
	word7	8
	word8	4
	.	.
	.	.
	.	.

$$P(X=x | Y=y) = P(\text{word1} = \underline{5}, \text{word2} = \underline{2}, \text{word3} = \underline{10}, \dots | Y=y)$$

Bayes classifier:

$$\arg \max_y P(x^{(1)}, \dots, x^{(d)} | y) P(y)$$

*Be Categorical*

Naïve Bayes classifier:

$$\arg \max_y \prod_{i=1}^d P(x^{(i)} | y) P(y)$$

# Glossary of Machine Learning

- iid random variables
- Class prior  $p(Y)$
- Class conditional distribution of inputs  $p(\underline{x}|Y)$
- Optimal classifier under 0/1 loss  $\arg \max_y P(Y|X)$
- Bayes rule  $P(X|Y)P(Y)$
- Gaussian Bayes classifier
- Naïve Bayes classifier
- Decision boundary