# Learning Distributions
## Maximum Likelihood Estimate (MLE)
## Bayes Classifier
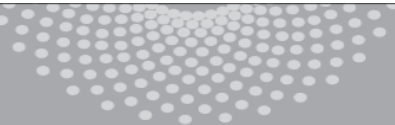
Aarti Singh

Machine Learning 10-315
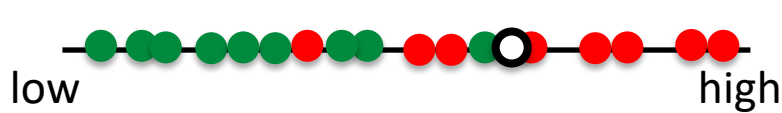Sept 14, 2020

# Modeling class distribution



low                                          high

X, average brain activity in "Amygdala"

f(X)

● Stress

● No Stress

○ Test
   subject

Modeling Class distribution P(Y=y)
                        = Bernoulli(θ)

$$P(Y = \textcolor{red}{\bullet}) = \theta \qquad\qquad P(Y = \textcolor{green}{\bullet}) = 1 - \theta$$

$$\underset{y}{\text{argmax}} \; P(x|\dot{y})\,P(\dot{y})$$
$$= \underset{y}{\text{argmax}} \; P(y|x)$$

Like a coin flip

# Bernoulli distribution

Data, D =



$X_1$  $X_2$  $X_3$  . . . .

- P(Heads) = $\theta$,  P(Tails) = 1-$\theta$

- Flips are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to Bernoulli distribution

$$P(D) = P(X_1, \ldots X_n) = \prod_{i=1}^{n} P(X_i)$$

$$\equiv P_\theta(D)$$

$$= P(D|\theta)$$

$$X_i \sim Ber(\theta)$$

Choose $\theta$ that maximizes the probability of observed data aka Likelihood

# Maximum Likelihood Estimation (MLE)

Choose θ that maximizes the probability of observed data (aka likelihood)

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

MLE of probability of head:

$$\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

#heads in D

= 3/5

# tails in D

"Frequency of heads"

# Derivation

$$D = \{X_1, \ldots X_n\}$$

$$\hat{\theta}_{MLE} = \arg\max_{\theta} \quad P(D \mid \theta)$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i \mid \theta) \qquad X_i \text{ ind}$$

$$\begin{array}{l} \arg \\ \max_{\theta} \, f(\theta) \\ = \arg\max_{\theta} \log f(\theta) \end{array}$$

$$= \arg\max_{\theta} \prod_{i=1}^{\alpha_H} \theta \prod_{j=1}^{\alpha_T} (1-\theta) \qquad \begin{array}{l} X_i \sim Ber(\theta) \\ P(X_i = H) = \theta \end{array}$$

$$= \arg\max_{\theta} \theta^{\alpha_H} (1-\theta)^{\alpha_T}$$

$$J(\theta) = \log\left(\theta^{\alpha_H}(1-\theta)^{\alpha_T}\right) = \alpha_H \log\theta + \alpha_T \log(1-\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\alpha_H}{\theta} + \frac{\alpha_T}{1-\theta} \cdot (-1) = \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} \Big|_{\hat{\theta}_{MLE}} = 0$$
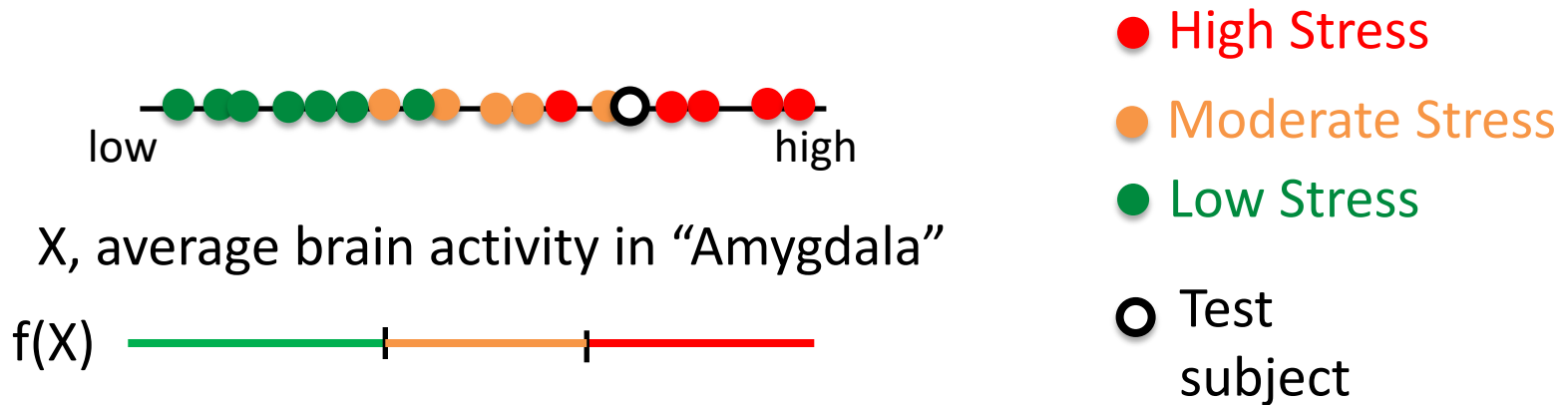
5

# Derivation

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

$$\frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} = 0 \implies \frac{\alpha_H}{\theta} = \frac{\alpha_T}{1-\theta}$$

$$\implies \alpha_H - \alpha_H\theta = \alpha_T\theta \implies \widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

$$X \sim Ber(\theta)$$

$$X_1 \ldots X_n \overset{iid}{\sim} Ber(\theta)$$

$$P(X_i) = P(X = X_i)$$

$$X_{ij} \leftarrow j^{th} \text{ feat}$$

$$X_i \leftarrow \text{data point}$$

# Modeling class distribution

● High Stress

● Moderate Stress

● Low Stress

low                high

X, average brain activity in "Amygdala"

f(X) ————————

○ Test subject

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

Modeling Class distribution P(Y) = Multinomial($p_H$, $p_M$, $p_L$)

$$P(Y = \bullet) = p_H \quad P(Y = \bullet) = p_M \quad P(Y = \bullet) = p_L$$

P(y=

Like a dice roll

$$p_H + p_M + p_L = 1$$

# Multinomial distribution

Data, D = rolls of a dice

$1, 6, 5, 2, 2, 1, 3, 4,$

- P(1) = $p_1$, P(2) = $p_2$, ..., P(6) = $p_6$    $p_1 + .... + p_6 = 1$

- Rolls are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to Multinomial($\theta$) distribution where

$$\theta = \{p_1, p_2, ..., p_6\}$$

Choose $\theta$ that maximizes the probability of observed data aka "Likelihood"

8

# Maximum Likelihood Estimation (MLE)

Choose $\theta$ that maximizes the probability of observed data

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} \ P(D \mid \theta)$$

*(handwritten)* $= p_1 \cdots p_6$
$p_1 + \cdots + p_6 = 1$
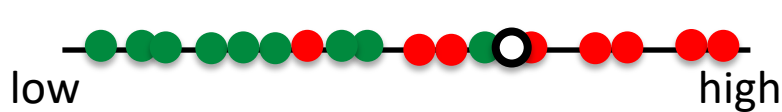
MLE of probability of rolls:

$$\widehat{\theta}_{MLE} = \hat{p}_{1,MLE}, \cdots, \hat{p}_{6,MLE}$$

$$\hat{p}_{y,MLE} = \frac{\alpha_y}{\sum_y \alpha_y}$$

*(handwritten)* ← # times dice rolls y
Rolls that turn up y

*(handwritten)* ← total # rolls
Total number of rolls

*(handwritten)* "Frequency of roll y"

# Bayes Classifier

low                           high

● Stress

● No Stress

○ Test subject

X, average brain activity in "Amygdala"

f(X)

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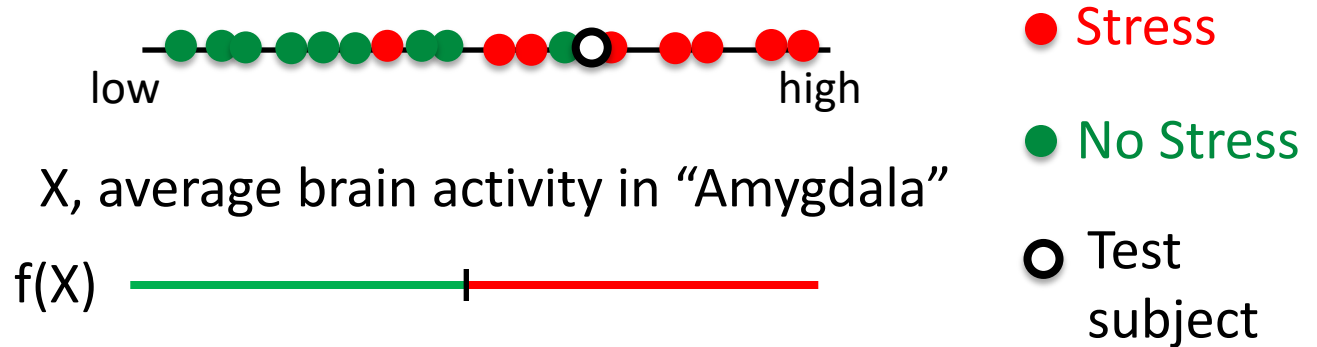
Class conditional Distribution of features        Class distribution

We can now consider appropriate distribution models for the two terms:

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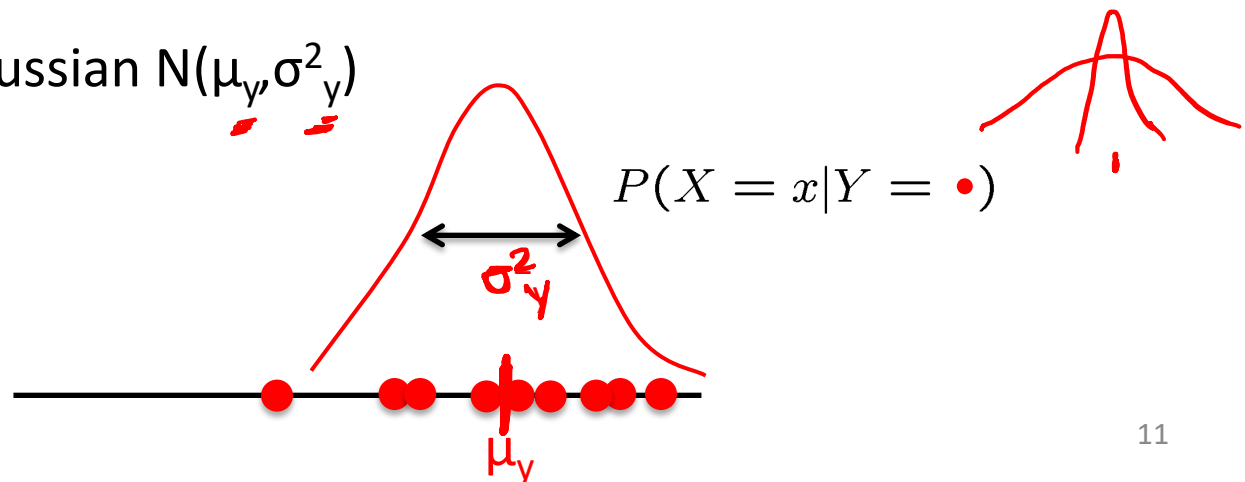
Class conditional distribution of features $P(X=x|Y=y)$

10

# Modeling class conditional distribution of features

Stress

No Stress

Test subject

low                                                      high

X, average brain activity in "Amygdala"

f(X)

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

➤ What distribution would you use?

E.g. $P(X=x|Y=y)$ = Gaussian $N(\mu_y, \sigma^2_y)$

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

$\sigma^2_y$

$\mu_y$

# 1-dim Gaussian distribution

X is Gaussian N(μ,σ²)

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

# Why Gaussian?

- Properties $\mu, \Sigma$ $\quad E[X_i^2] \quad E[X_i X_j]$

  – Fully Specified by first and second order statistics

    - Uncorrelated $\Leftrightarrow$ Independence $\quad P(X,Y) = P(X)P(Y)$

  – X, Y Gaussian => aX+bY Gaussian $\quad E[X,Y] = 0$

  – <u>Central limit theorem</u>: if $X_1, \ldots, X_n$ are any iid random variables with mean $\mu$ and variance $\sigma^2 < \infty$

  then

  $$\sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n} X_i - \mu\right) \to N(0, \sigma^2)$$

  $n \to \infty$

13

# d-dim Gaussian distribution

$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$ ← brain activity @ pixel i

X is Gaussian N(μ, Σ)          μ is d-dim vector, Σ is dxd dim matrix

$\mu_i = E[x_i]$

$\Sigma_{ji} = \Sigma_{ij} = E\left[(X_i - E[X_i])(X_j - E[X_j])\right]$

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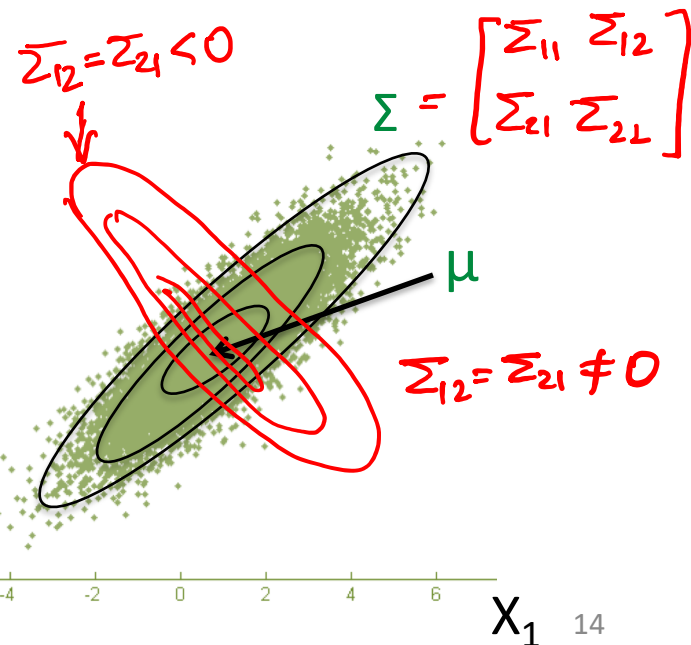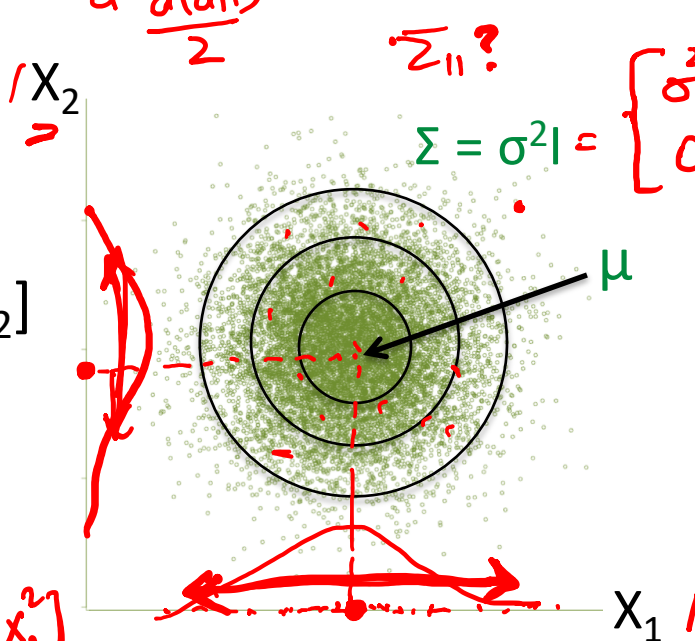
$= \sigma$

$d \quad \frac{d(d+1)}{2}$

$\Sigma_{11}?$

$\Sigma = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$

d=2
X = [X₁; X₂]

$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$

$E[X_2^2]$

μ

$\Sigma_{12} = \Sigma_{21} < 0$

$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

μ

$\Sigma_{12} = \Sigma_{21} \neq 0$



14

# How to learn parameters from data? MLE

# (Continuous case)

# Gaussian distribution

Data, D = _____

|  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|
| 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Sleep hrs

How many hours did you sleep last night?

➢ Poll

# Gaussian distribution



Data, D =

**Sleep hrs**

3    4    5    6    7    8    9

- Parameters: $\mu$ – mean, $\sigma^2$ - variance

- Sleep hrs are **i.i.d.**:
  – **Independent** events
  – **Identically distributed** according to Gaussian distribution

$X_i$ ~ sleep hrs of person $i$
~ avg brain activity of person $i$

# Maximum Likelihood Estimation (MLE)

Choose $\theta = (\mu, \sigma^2)$ that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i \mid \theta)$$

$D = \{X_1 ... X_n\}$

Independent draws

# Maximum Likelihood Estimation (MLE)

Choose $\theta = (\mu, \sigma^2)$ that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg\max_{\theta} \; P(D \mid \theta)$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i \mid \theta) \qquad \text{Independent draws}$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2 / 2\sigma^2}$$

$= \mu, \sigma^2$

Identically distributed

# Maximum Likelihood Estimation (MLE)

Choose $\theta = (\mu, \sigma^2)$ that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i \mid \theta) \quad \text{Independent draws}$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2 / 2\sigma^2} \quad \text{Identically distributed}$$

$$= \arg\max_{\theta=(\mu,\sigma^2)} \underbrace{\frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^{n}(X_i - \mu)^2 / 2\sigma^2}}_{J(\theta)}$$

$J(\mu)$

$$= \arg\min_{\mu} \sum_{i=1}^{n} \frac{(X_i - \mu)^2}{2\sigma^2}$$

$$\frac{\partial}{\partial \mu} J(\mu) = 0$$

# Derivation

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$
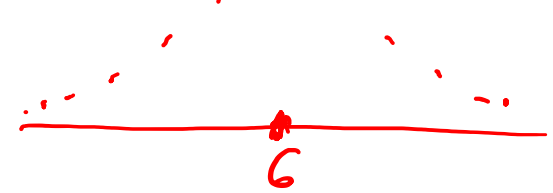
➢ Breakout

Groups 1-10: Jamboard_1_10
Groups 11-20: Jamboard_11_20

# MLE for Gaussian mean and variance

$$X_1 \ldots X_n \sim N(\mu, \sigma^2)$$

$$E[\hat{\mu}_{MLE}] = \mu \checkmark$$

$$\hat{\mu}_{MLE} = \frac{1}{n}\sum_{i=1}^{n} x_i$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n}\sum_{i=1}^{n} (x_i - \hat{\mu})^2$$

$$E[\hat{\sigma}^2_{MLE}] \neq \sigma^2$$

$$\hat{\mu} = \hat{\mu}_{MLE}$$

$$\frac{n-1}{n}\sigma^2$$

(unbiased)  $n-1$

# Gaussian Bayes classifier

*cont.*

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

Class conditional
Distribution of features

Class distribution

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

Bernoulli($\theta$)

How to learn parameters
$\theta$, $\mu_y$, $\Sigma_y$ from data?

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

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

# 1-dim Gaussian Bayes classifier

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

Class conditional
Distribution of features

Class distribution

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

Gaussian($\mu_y$, $\sigma^2_y$)

Bernoulli($\theta$)

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

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

$\mu_1 = \mu_2$