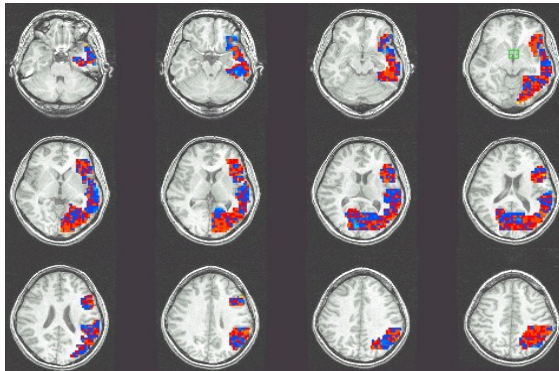


# Announcements

- Recitation on Friday Jan 28 – Convexity review
- QnA1 due TODAY
- HW1 to be released TODAY

# Recap – Bayes classifier



High Stress  
Moderate Stress  
Low Stress

$(X, Y)$  - random variables with joint distribution  $P_{XY}$

**Input feature vector,  $X$**

**Label,  $Y$**

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

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

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

Class conditional

Class distribution

Distribution of features

# Recap – Gaussian Bayes classifier

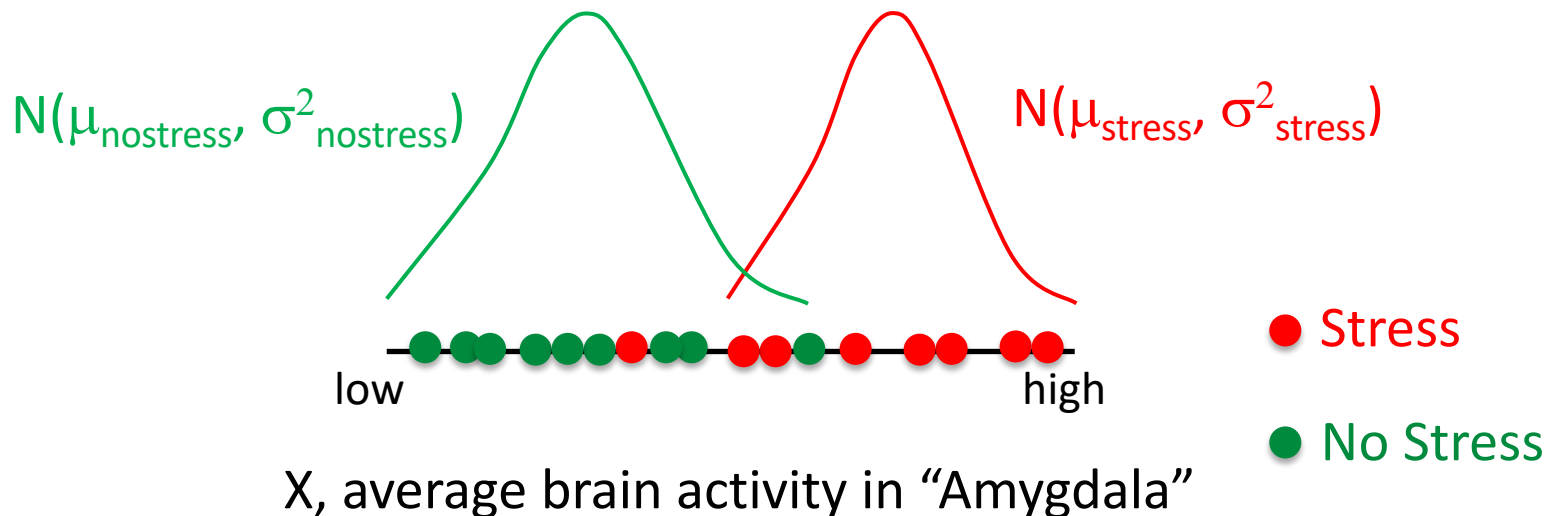
In practice  $P_{XY}$  unknown, use a distribution model to approximate

**Gaussian Bayes classifier** – assumes

Class distribution  $P(Y)$  is Bernoulli( $\theta$ )

[Categorical if multiple classes]

Class conditional distribution of features  $P(X|Y)$  is Gaussian



# d-dim Gaussian Bayes classifier

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

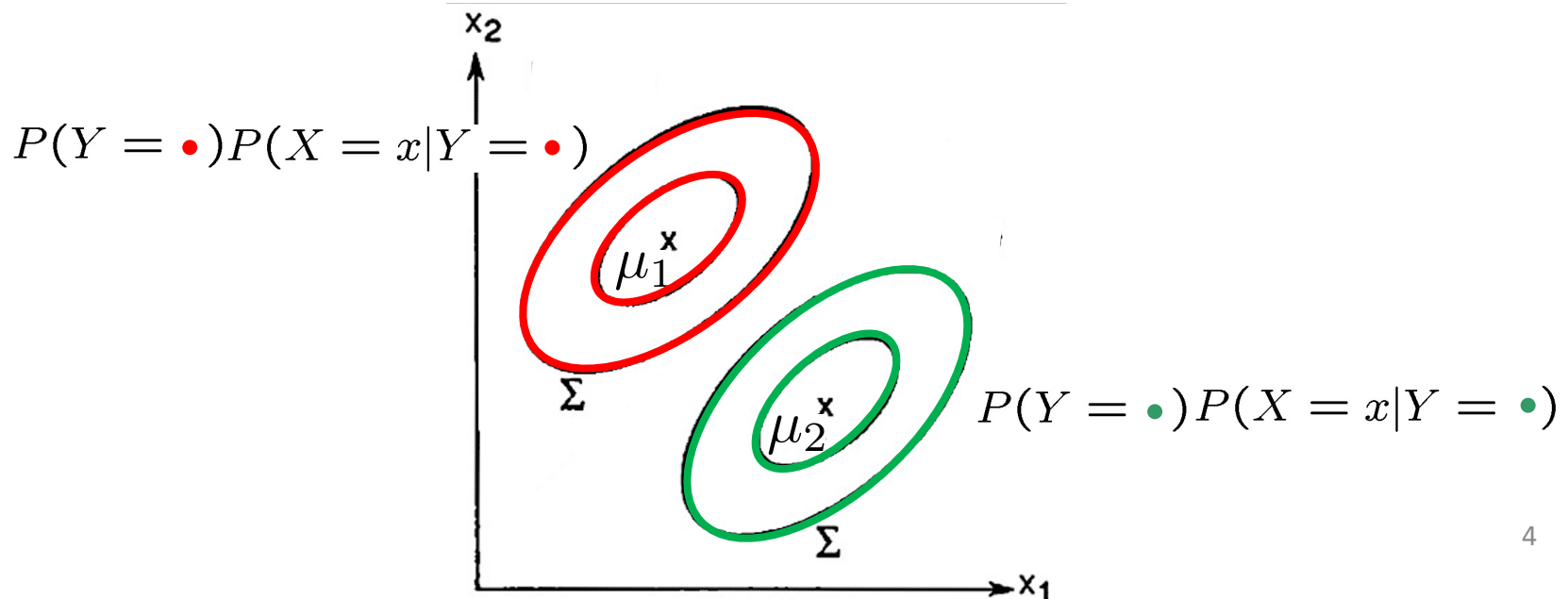
Learn parameters  $\theta, \mu_y, \Sigma_y$  from data

Class conditional  
Distribution of inputs

Class distribution

Gaussian( $\mu_y, \Sigma_y$ )

Bernoulli( $\theta$ )



# 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?

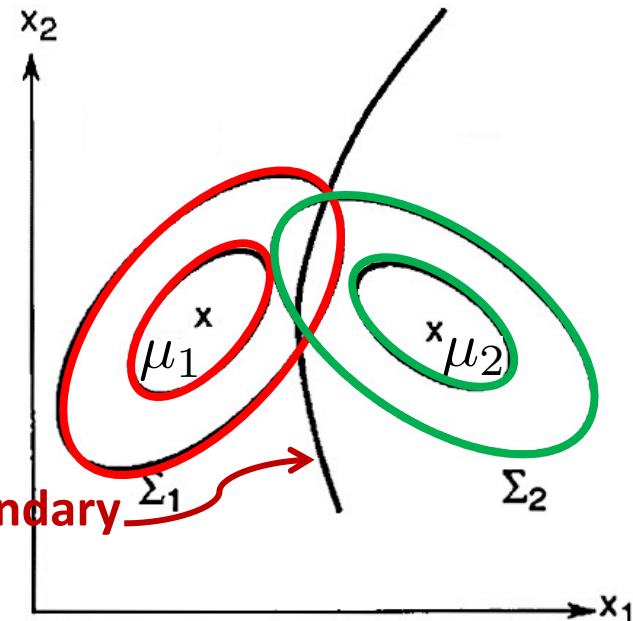
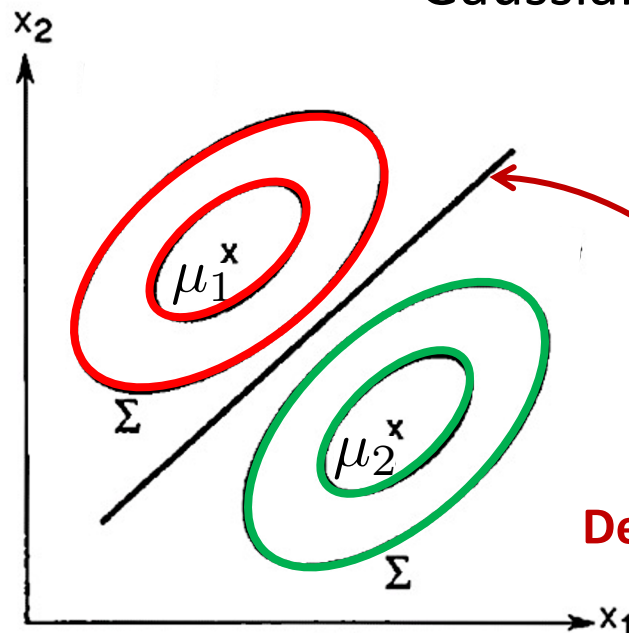
Class conditional

Class distribution

Distribution of inputs

Gaussian( $\mu_y, \Sigma_y$ )

Bernoulli( $\theta$ )



# 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

$$\begin{aligned} 1 &= \frac{P(Y = 1 | X = x)}{P(Y = 0 | X = x)} = \frac{P(X = x | Y = 1)P(Y = 1)}{P(X = x | Y = 0)P(Y = 0)} \\ &= \sqrt{\frac{|\Sigma_0|}{|\Sigma_1|}} \exp \left( -\frac{(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}{2} + \frac{(x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0)}{2} \right) \frac{\theta}{1 - \theta} \end{aligned}$$

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.

# d-dim Gaussian Bayes classifier

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

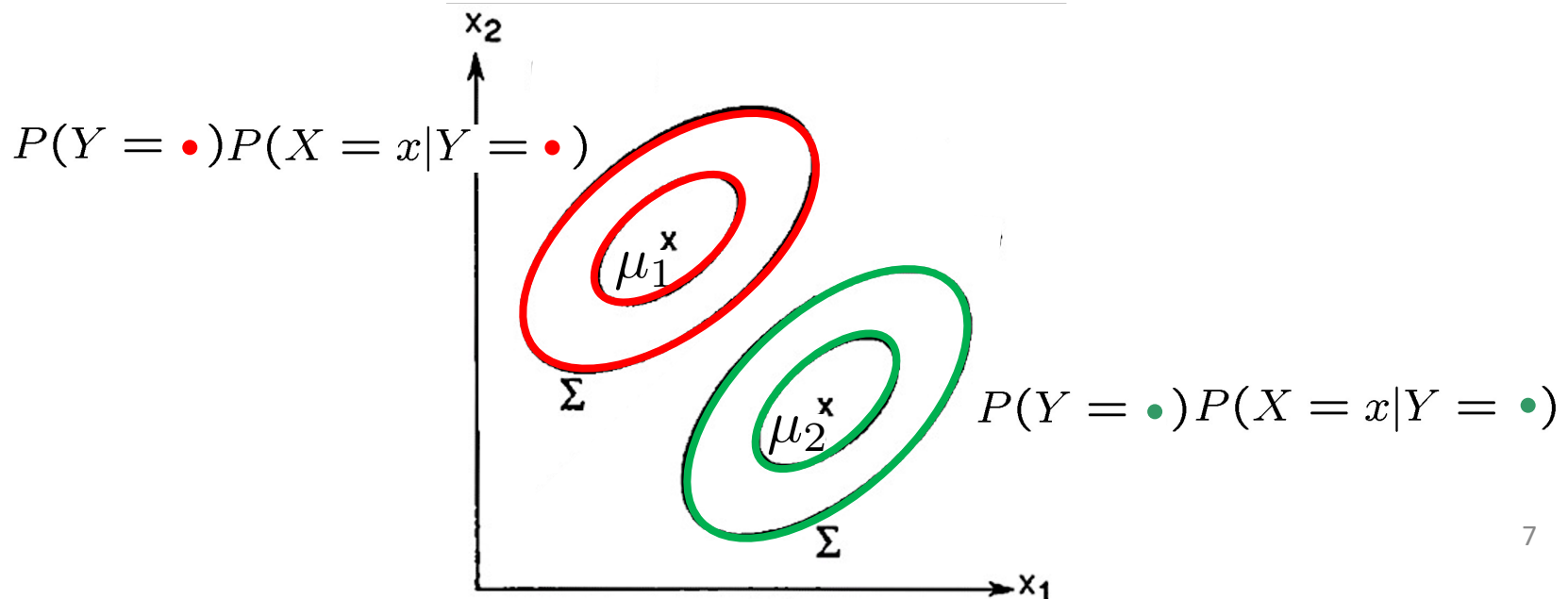
Learn parameters  $\theta, \mu_y, \Sigma_y$  from data

Class conditional  
Distribution of inputs

Class distribution

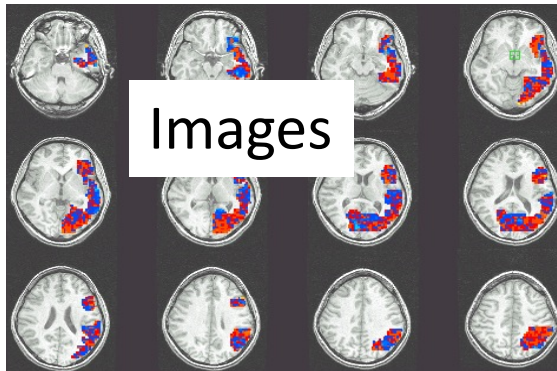
Gaussian( $\mu_y, \Sigma_y$ )

Bernoulli( $\theta$ )

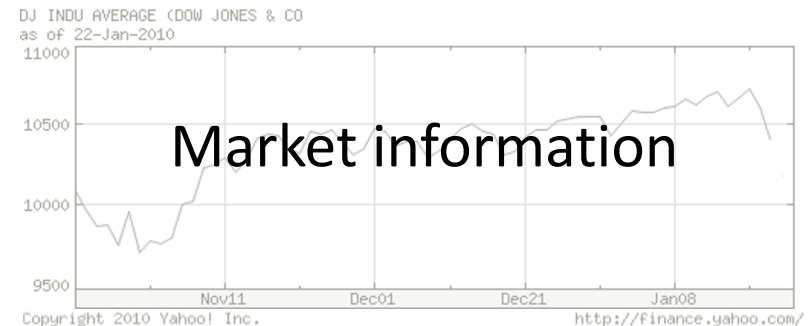


# Notion of “Features aka Attributes”

Input  $X \in \mathcal{X}$



Input  $X \in \mathcal{X}$



## How to represent inputs mathematically?

- Image  $X$  = intensity/value at each pixel, fourier transform values, SIFT etc.
- Market information  $X$  = daily/monthly? price of share for past 10 years



# Notion of “Features aka Attributes”

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?
    - list of n-grams (n-tuples of words)
- Misses out context!!

# Text classification

Raw input



Features



Model for input features



word1	5
word2	2
word3	10
word4	20
word5	12
word6	5
word7	8
word8	4

.	.
.	.
.	.

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

HW1!

# Glossary of Machine Learning

- Task
- Supervised learning
  - Classification
  - Regression
- Unsupervised learning
  - Learning distribution
  - Clustering
  - Dimensionality reduction/Embedding
- Input,  $X$
- Label,  $Y$
- Prediction,  $f(X)$
- Experience = Training data
- Test data
- Overfitting
- Generalization
- Performance measure/loss – 0/1, squared
- iid
- Class conditional distribution of inputs
- Bayes rule
- Bayes Optimal classifier
- Decision boundary
- Feature/Attribute

# Maximum Likelihood Estimation (MLE)

Aarti Singh

Machine Learning 10-315

Jan 26, 2022



**MACHINE LEARNING** DEPARTMENT



**How to learn parameters from data?**

**MLE**

**(Discrete case)**

# Learning parameters in distributions

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

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

Learning  $\theta$  is equivalent to learning probability of head in coin flip.

➤ How do you learn that?

Data =



Answer: 3/5

➤ Why??

# Bernoulli distribution

Data,  $D =$



- Parameter  $\theta$  :  $P(\text{Heads}) = \theta$ ,  $P(\text{Tails}) = 1-\theta$
- Flips are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to Bernoulli distribution

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

# Maximum Likelihood Estimation (MLE)

Choose  $\theta$  that maximizes the probability of observed data (aka likelihood)

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

MLE of probability of head:

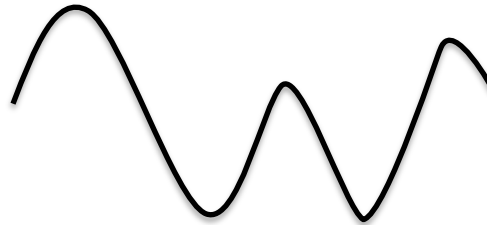
$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} = 3/5$$

"Frequency of heads"

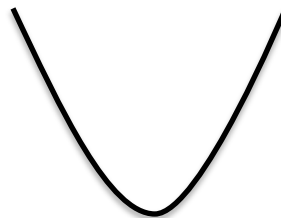


# Short detour - Optimization

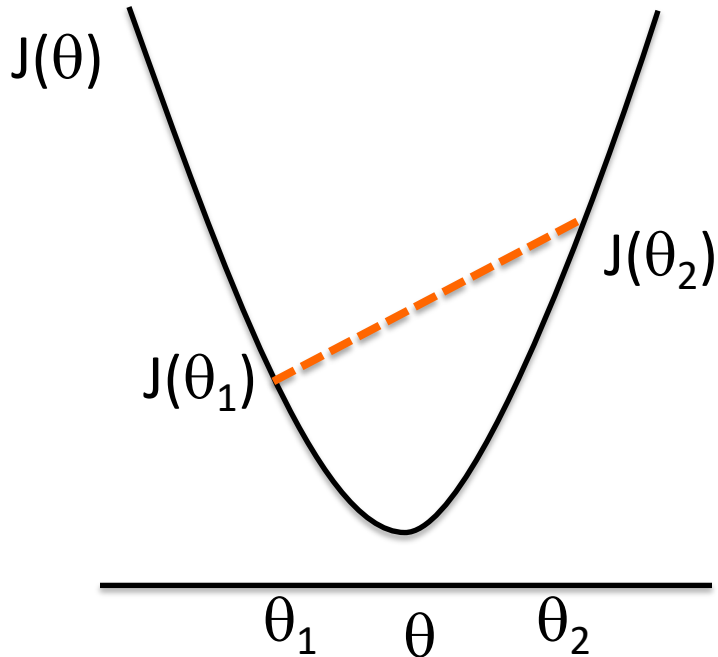
- Optimization objective  $J(\theta)$
- Minimum value  $J^* = \min_{\theta} J(\theta)$
- Minima (points at which minimum value is achieved) may not be unique



- If function is strictly convex, then minimum is unique

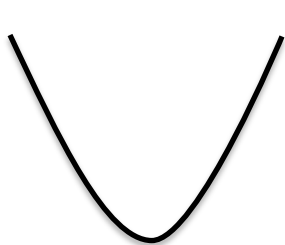


# Convex functions

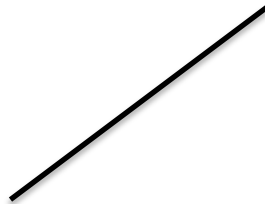


A function  $J(\theta)$  is called **convex** if the line joining two points  $J(\theta_1), J(\theta_2)$  on the function does not go below the function on the interval  $[\theta_1, \theta_2]$

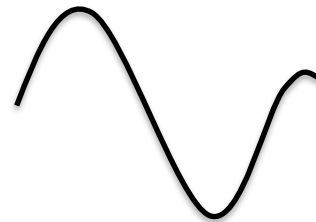
(Strictly) Convex functions  
have a unique minimum!



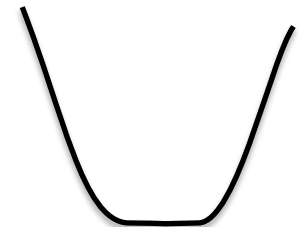
Convex



Both Concave  
& Convex



Neither



Convex but not  
strictly convex<sup>18</sup>

# Optimizing convex (concave) functions

- Derivative of a function
- Derivative is zero at minimum of a convex function
- Second derivative is positive at minimum of a convex function

# Bernoulli MLE Derivation

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

# Categorical distribution

Data,  $D$  = rolls of a dice



- $P(1) = p_1, P(2) = p_2, \dots, P(6) = p_6 \quad p_1 + \dots + p_6 = 1$
- Rolls are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to  $\text{Categorical}(\theta)$  distribution where

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

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

# Maximum Likelihood Estimation (MLE)

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

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

MLE of probability of rolls:

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

$$\hat{p}_{y,MLE} = \frac{\alpha_y \longleftarrow \text{Rolls that turn up } y}{\sum_y \alpha_y \longleftarrow \text{Total number of rolls}}$$

“Frequency of roll  $y$ ”

**How to learn parameters from data?**

**MLE**

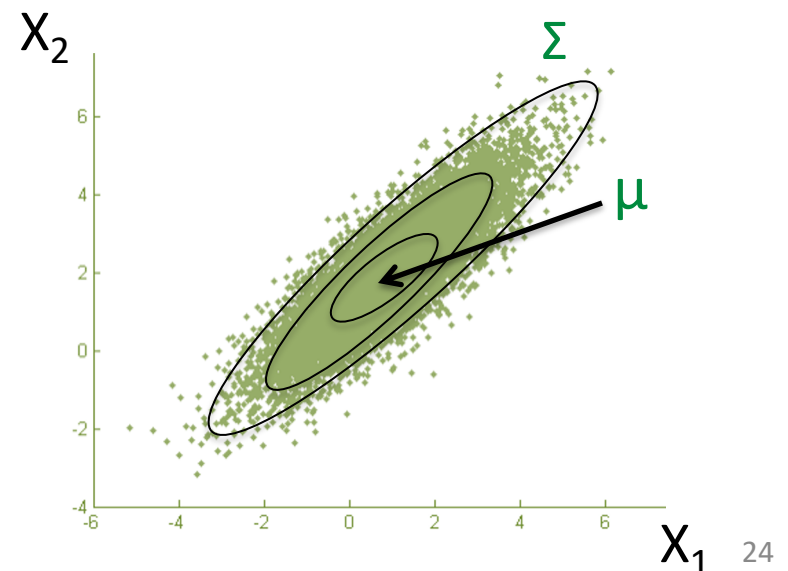
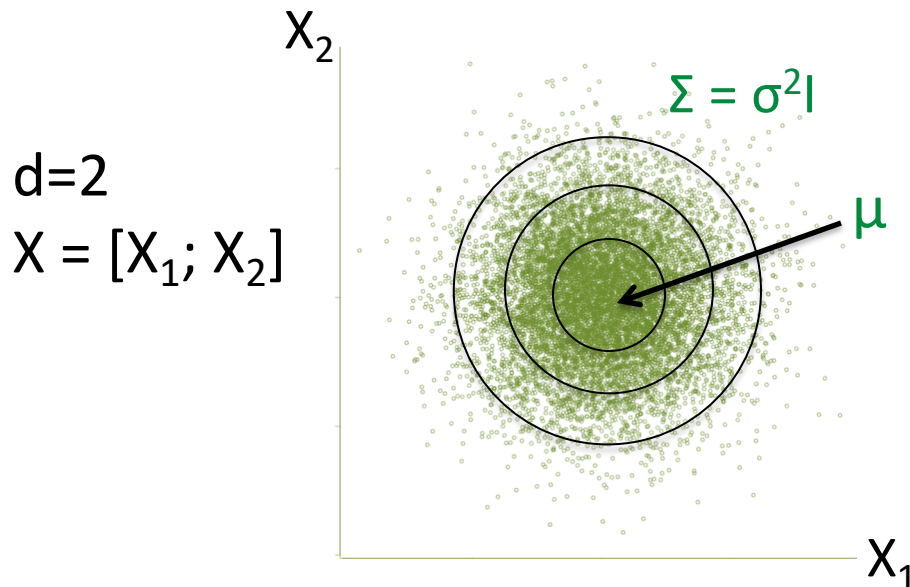
**(Continuous case)**

# d-dim Gaussian distribution

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

$\mu$  is d-dim vector,  $\Sigma$  is dxd dim matrix

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



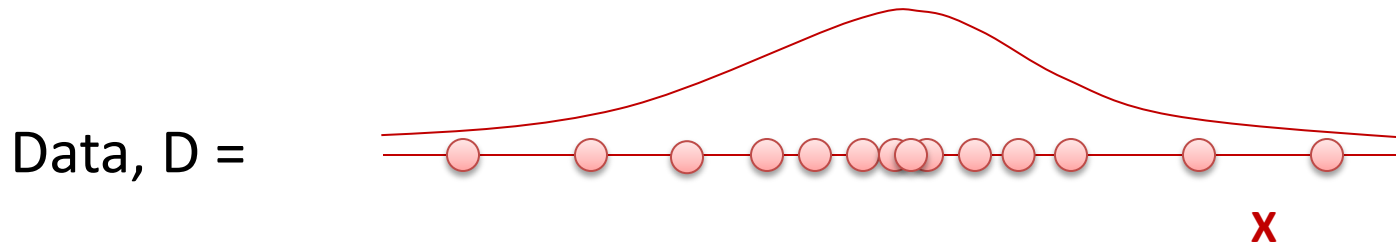


**How to learn parameters from data?**

**MLE**

**(Continuous case)**

# Gaussian distribution



- Parameters:  $\mu$  – mean,  $\sigma^2$  – variance
- Data are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to Gaussian distribution

# Maximum Likelihood Estimation (MLE)

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

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D \mid \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n P(X_i \mid \theta) \quad \text{Independent draws}\end{aligned}$$

# Maximum Likelihood Estimation (MLE)

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

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n P(X_i | \theta) \quad \text{Independent draws} \\ &= \arg \max_{\theta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2 / 2\sigma^2} \quad \text{Identically distributed}\end{aligned}$$

# Maximum Likelihood Estimation (MLE)

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

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D | \theta) \\&= \arg \max_{\theta} \prod_{i=1}^n P(X_i | \theta) \quad \text{Independent draws} \\&= \arg \max_{\theta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2 / 2\sigma^2} \quad \text{Identically distributed} \\&= \arg \max_{\theta = (\mu, \sigma^2)} \underbrace{\frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^n (X_i - \mu)^2 / 2\sigma^2}}_{J(\theta)}\end{aligned}$$

# MLE for Gaussian mean

➤ Poll

$$P(D|\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^n (X_i - \mu)^2 / 2\sigma^2}$$

A.  $\max_{\mu} \sum_{i=1}^n (X_i - \mu)^2$

C.  $\max_{\mu} \mu^2 - 2\mu \sum_{i=1}^n X_i$

B.  $\min_{\mu} \sum_{i=1}^n (X_i - \mu)^2$

D.  $\max_{\mu} n\mu^2 - 2\mu \sum_{i=1}^n X_i$

# MLE for Gaussian mean and variance

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

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

Self exercise:

Derive MLE of variance?

d-dimensional versions?

MLE for uniform or  
exponential  
distribution?

More coming up in HW1