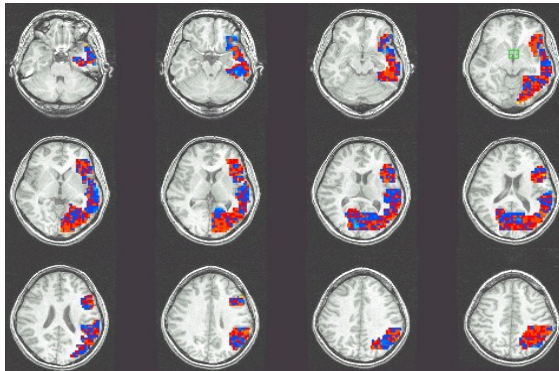


Announcements

- 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

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

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

Class conditional

Class distribution

Distribution of features

Recap – Gaussian Bayes classifier

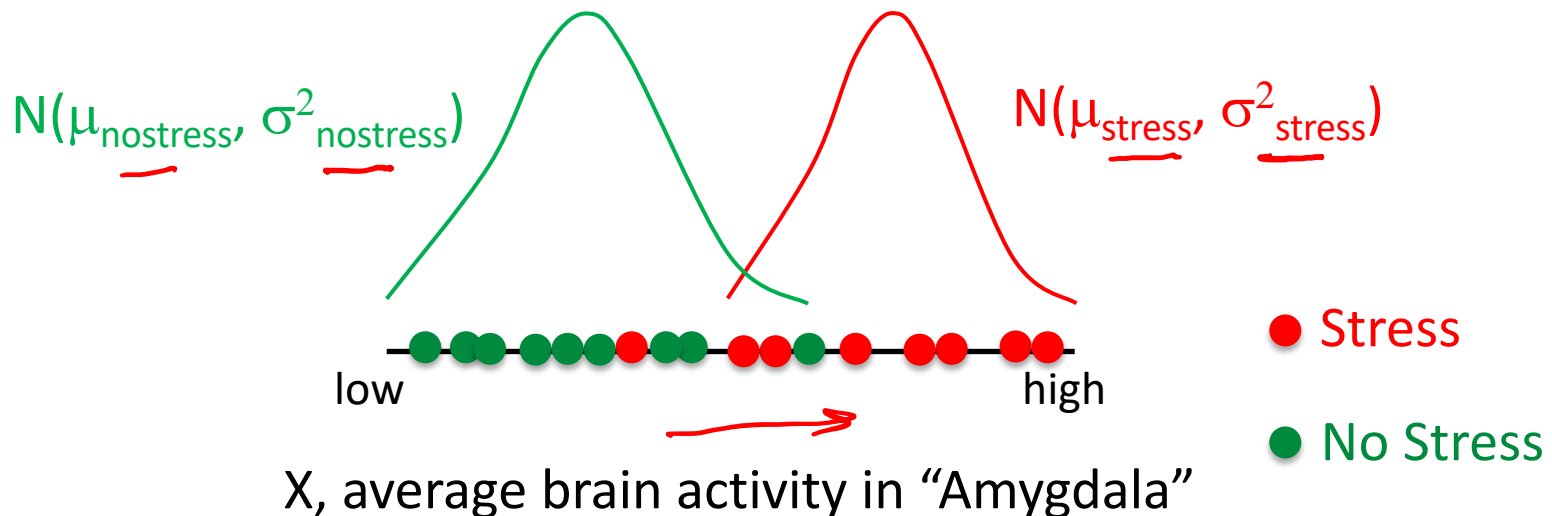
In practice P_{XY} unknown, use a distribution model to approximate

Gaussian Bayes classifier – assumes

Class distribution $P(Y)$ is Bernoulli(θ)

[Categorical if multiple classes]

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



$$X = \begin{bmatrix} x_{(1)} \\ x_{(2)} \\ \vdots \\ x_{(d)} \end{bmatrix}$$

d-dim Gaussian Bayes classifier

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

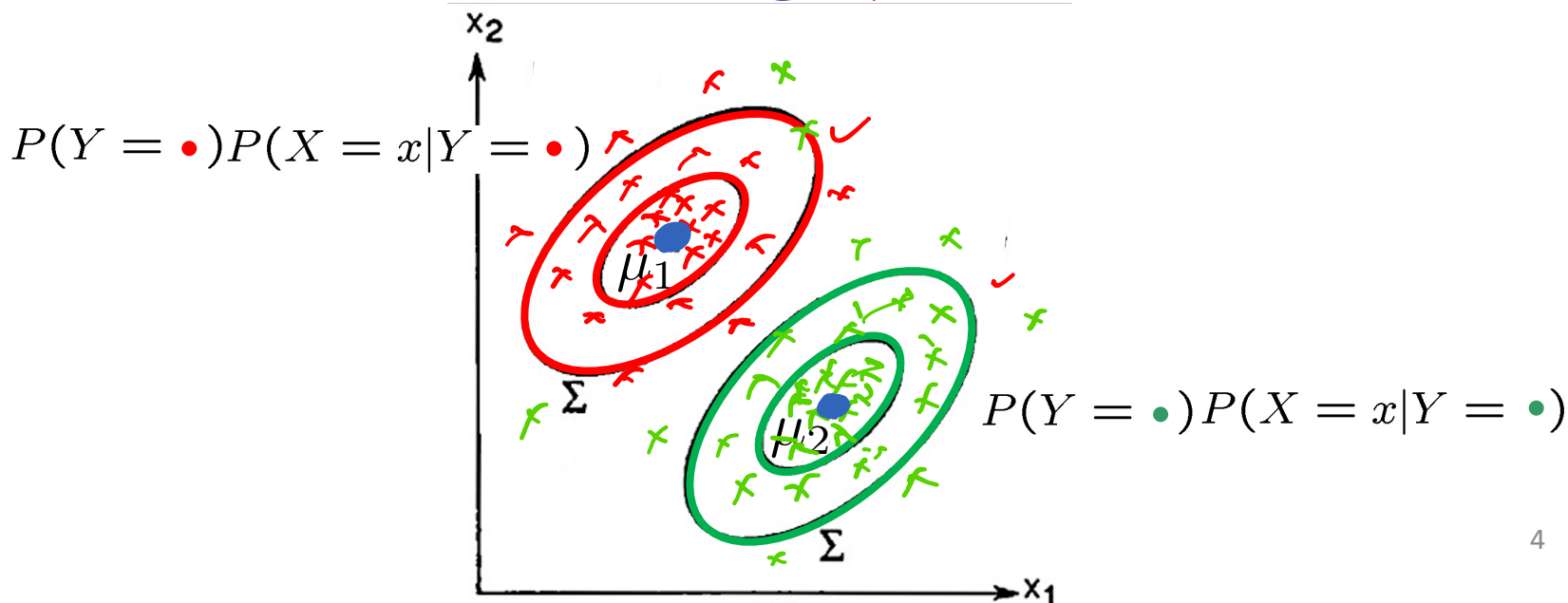
Learn parameters θ, μ_y, Σ_y from data

Class conditional
Distribution of inputs

Class distribution

Gaussian(μ_y, Σ_y)

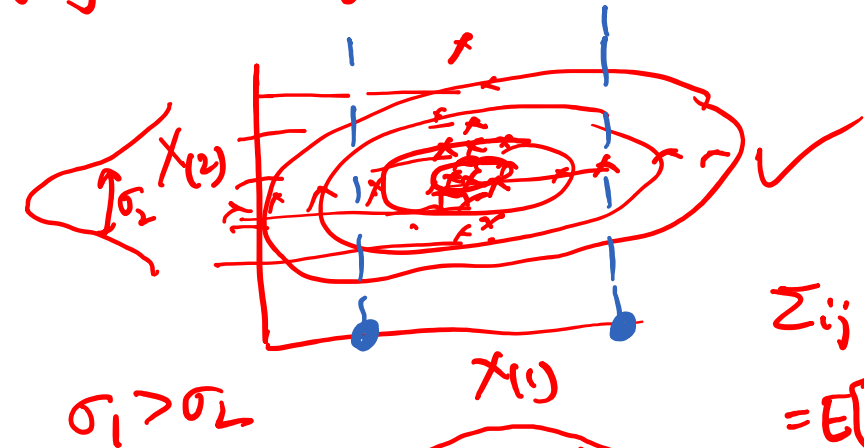
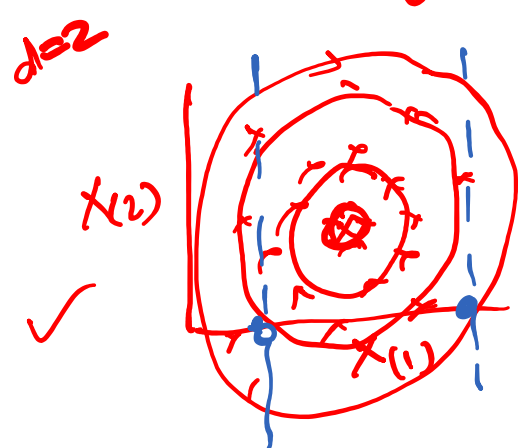
Bernoulli(θ)



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

$$\sigma_1 > \sigma_2 > \dots > \sigma_d$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots \\ 0 & & & \sigma_d^2 \end{bmatrix}$$

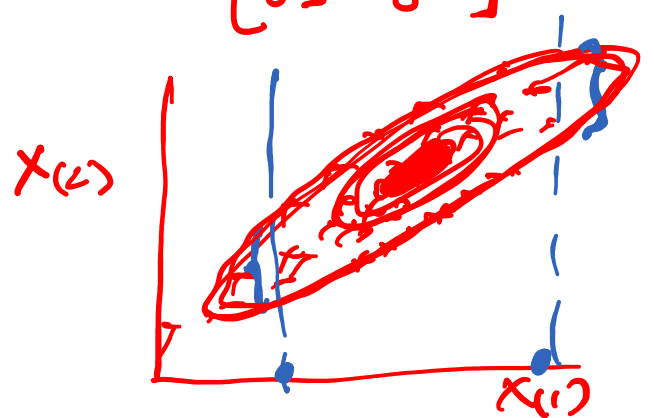


$$\Sigma_{ij} = \Sigma_{ji}$$

$$= E[(X_i - EX_i)(X_j - EX_j)]$$

$$\Sigma = \begin{bmatrix} \sigma^2 & 0.5 \\ 0.5 & \sigma^2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma^2 & -0.5 \\ -0.5 & \sigma^2 \end{bmatrix}$$



d-dim Gaussian Bayes classifier



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

Class conditional

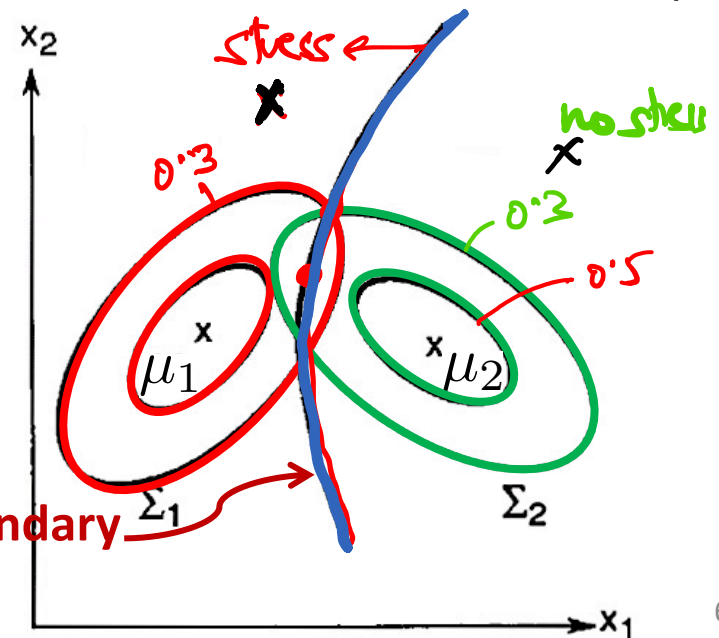
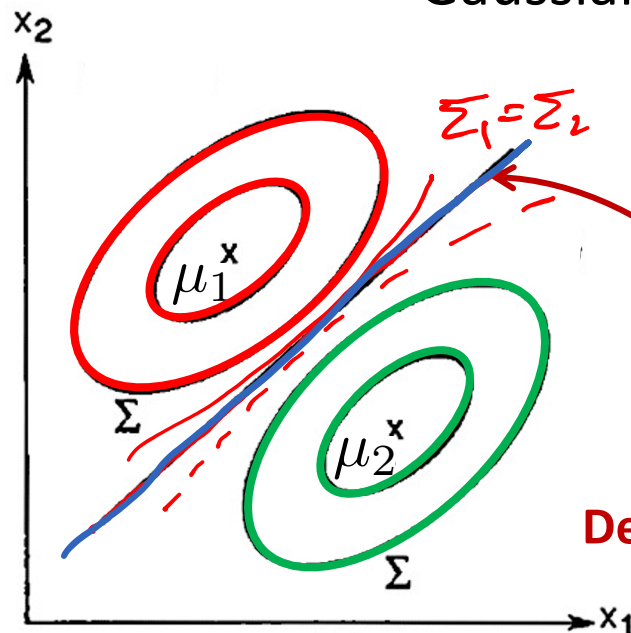
Class distribution

Distribution of inputs

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

Gaussian(μ_y, Σ_y)

Bernoulli(θ)



Decision Boundary of Gaussian Bayes

$$P(X|Y=\text{stress})P(Y=\text{stress}) = P(X|Y=\text{no s.})P(Y=\text{no s.})$$

stress no stress

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

Compute the ratio

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

$f(x) = 0$
 x^2 $x^T x$ $x^T A x$

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

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

Learn parameters θ, μ_y, Σ_y from data

Class conditional

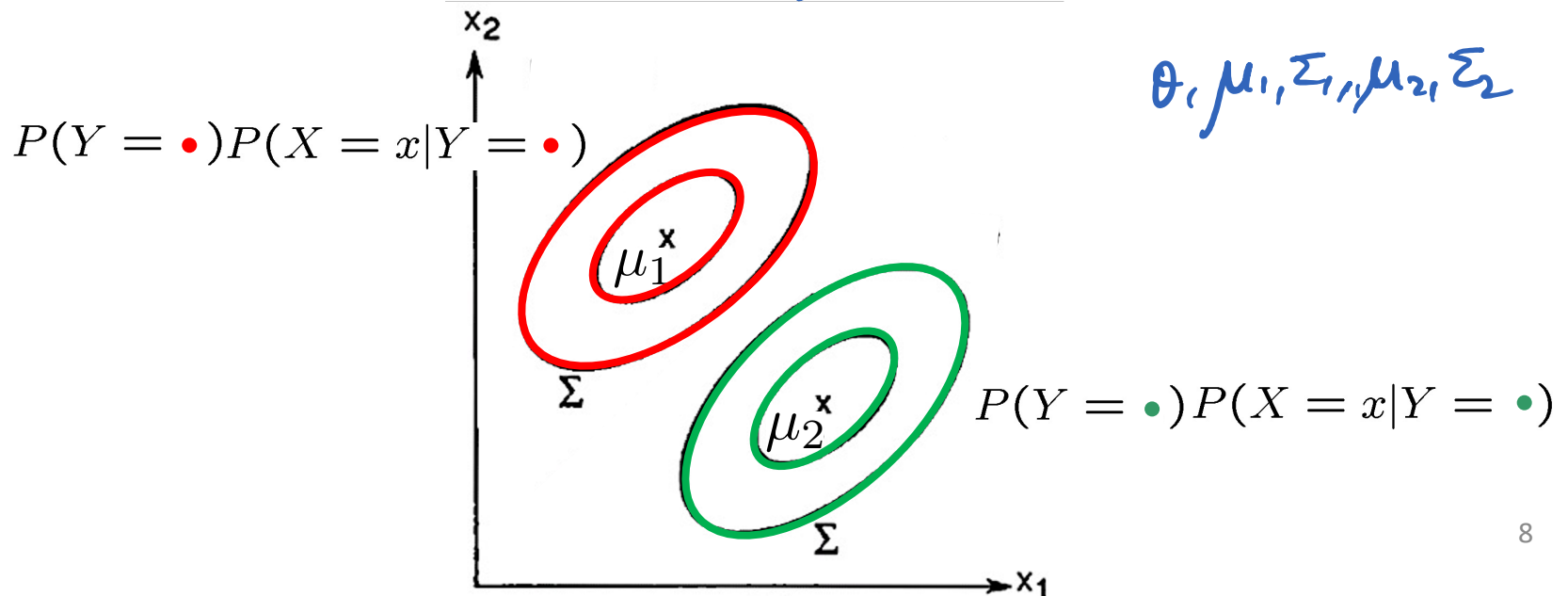
Distribution of inputs

Class distribution

Gaussian($\underline{\mu_y}, \underline{\Sigma_y}$)

Bernoulli(θ)

$\theta, \mu_1, \Sigma_1, \mu_2, \Sigma_2$

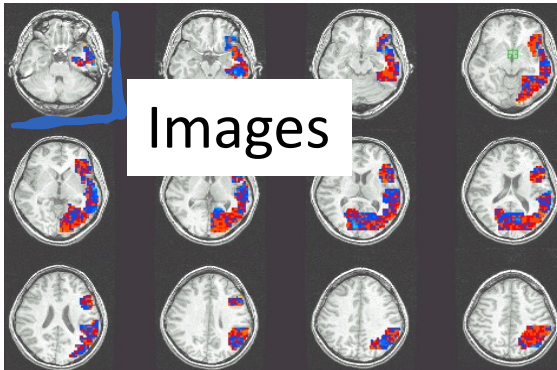


Notion of “Features aka Attributes”

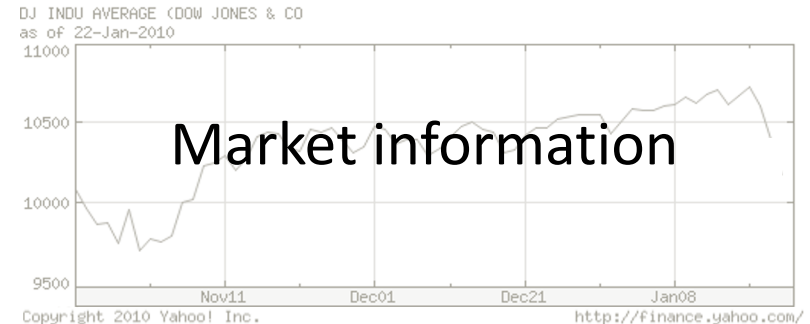
$$X = \begin{bmatrix} x_{00} \\ \vdots \\ x_{00} \\ \vdots \\ x_{00} \end{bmatrix}$$

Handwritten notes: 2.5 with an arrow pointing to the top of the vector, and -0.01 at the bottom.

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

$d = \text{size of vocabulary}$

$$X = \begin{bmatrix} 1^{\text{st}} \\ 2^{\text{nd}} \\ 10 \\ 5 \end{bmatrix} \begin{matrix} a \\ the \\ and \\ technology \\ book \end{matrix}$$

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!!

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

$X_{(1)}$	word1	5
$X_{(2)}$	word2	2
$X_{(3)}$	word3	10
	word4	20
	word5	12
	word6	5
	word7	8
	word8	4
	.	.
	.	.
	.	.



Model for input features

\downarrow

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

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

$$P(Y = \text{●}) = \theta$$

Head

$$P(Y = \text{●}) = 1 - \theta$$

Tail

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

Choose θ that maximizes the probability of observed data
aka Likelihood

Maximum Likelihood Estimation (MLE)

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

$$D = \{H, T, H, H, T, \dots\}$$

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

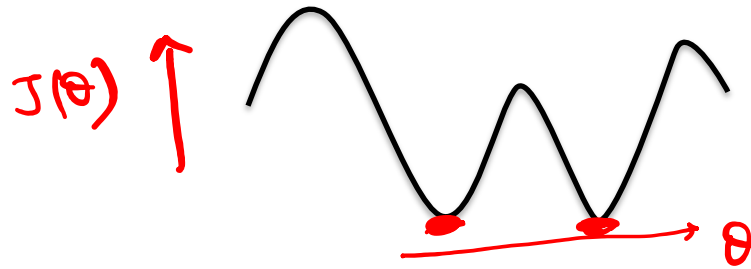
MLE of probability of head:

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

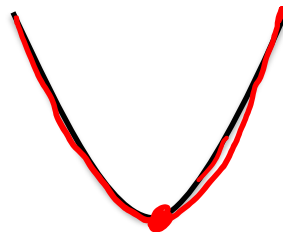
"Frequency of heads"

Short detour - Optimization

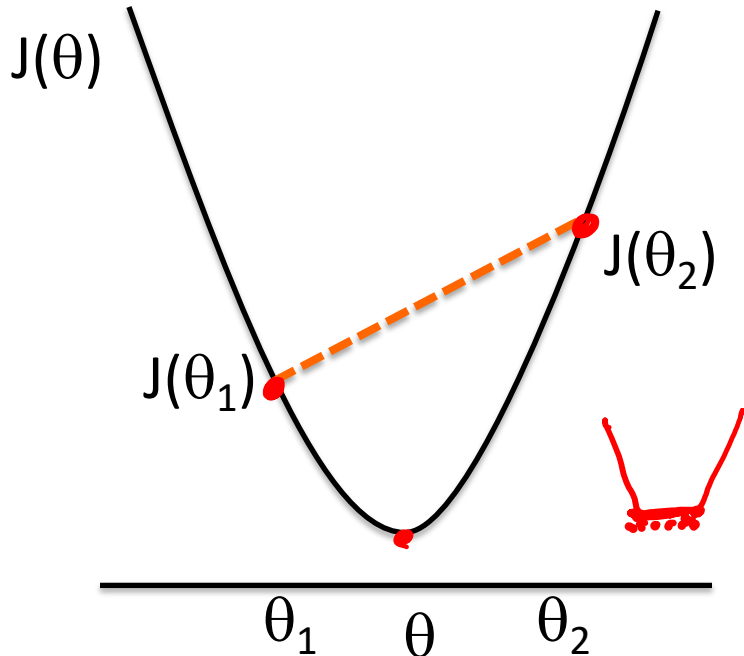
- Optimization objective $J(\theta)$
- Minimum value $J^* = \min_{\theta} J(\theta)$
- Minima (points at which minimum value is achieved) may not be unique



- If function is strictly convex, then minimum is unique

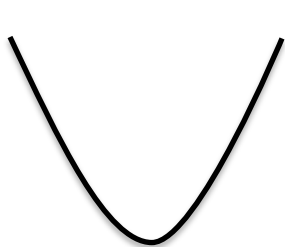


Convex functions

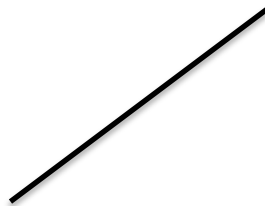


A function $J(\theta)$ is called **convex** if the line joining two points $J(\theta_1), J(\theta_2)$ on the function does not go below the function on the interval $[\theta_1, \theta_2]$

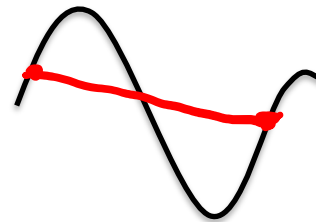
(Strictly) Convex functions have a unique minimum!



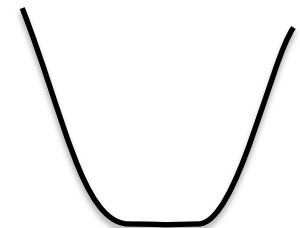
Convex



Both Concave & Convex



Neither



Convex but not strictly convex²⁰

Optimizing convex (concave) functions

- Derivative of a function

$$\frac{dJ(\theta)}{d\theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$



- Derivative is zero at minimum of a convex function



- Second derivative is positive at minimum of a convex function