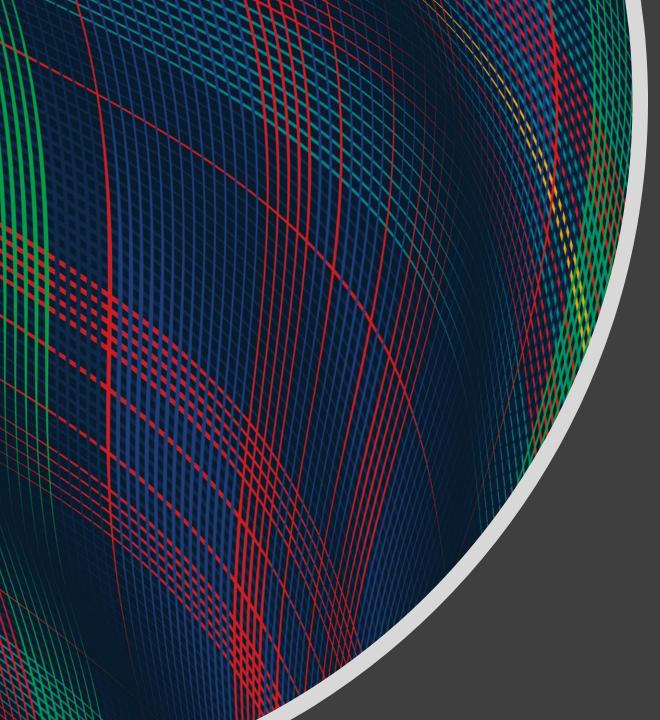
### Plan

### Today

- Wrap-up regularization (for now)
- MLE
  - Probability / likelihood
  - Maximum likelihood estimation
  - Probabilistic formulation of linear and logistic regression

### Wrap up Neural Nets

Switch to regression slides



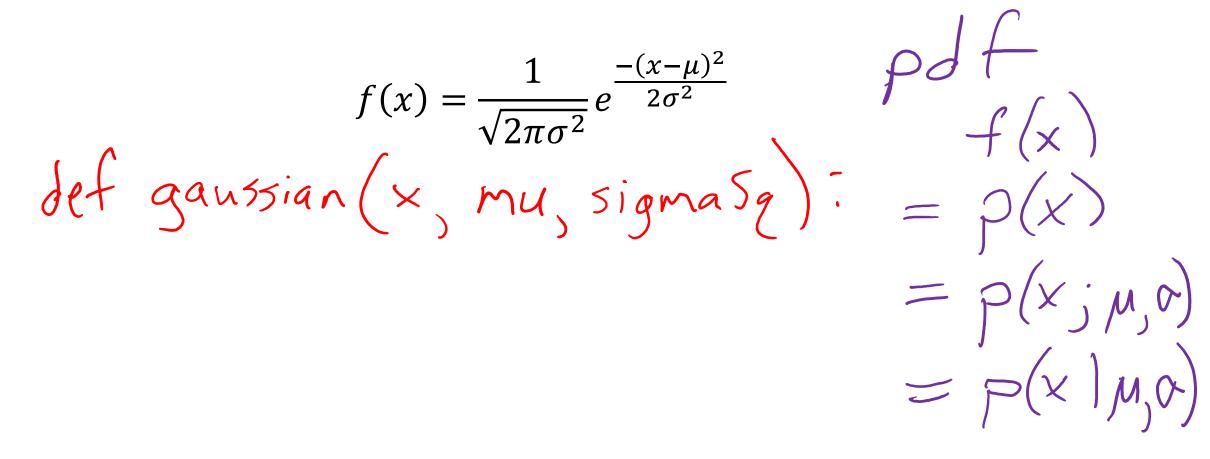
# 10-315 Introduction to ML

## MLE and Probabilistic Formulation of Machine Learning

Instructor: Pat Virtue

### Poll 1: Exercise

Implement a function in Python for the pdf of a Gaussian distribution. Python numpy or math packages are fine, no scipy, etc.



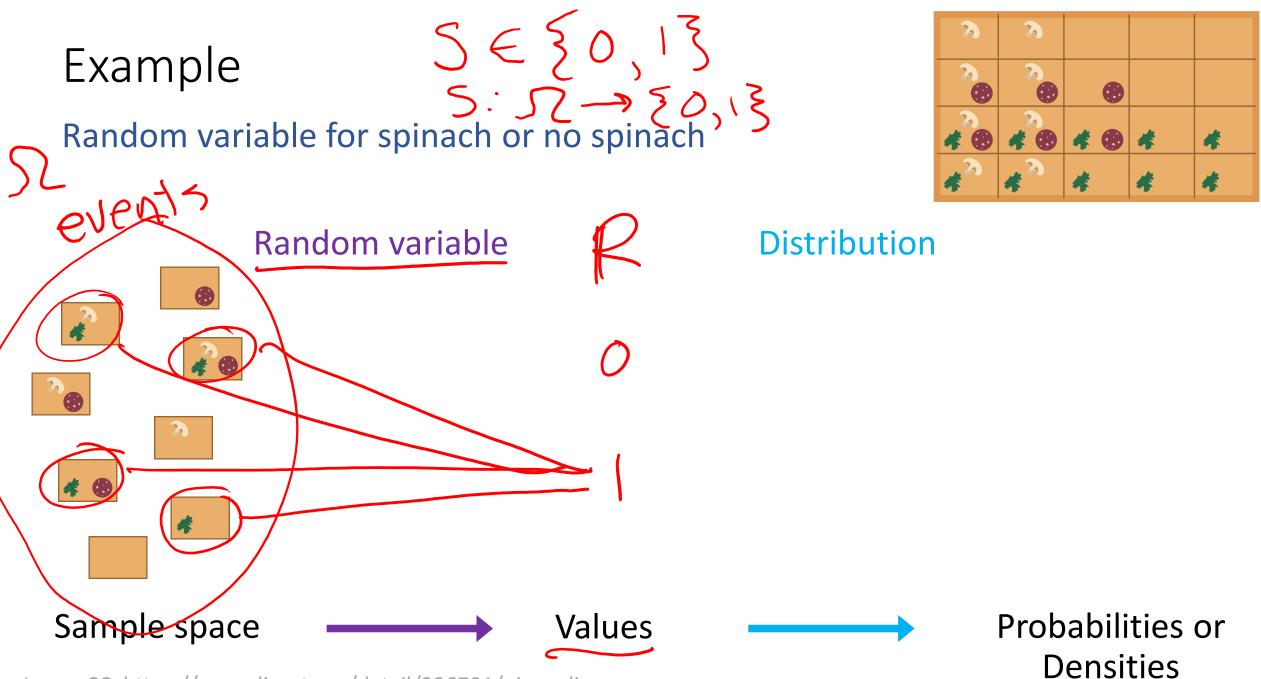
### Exercises

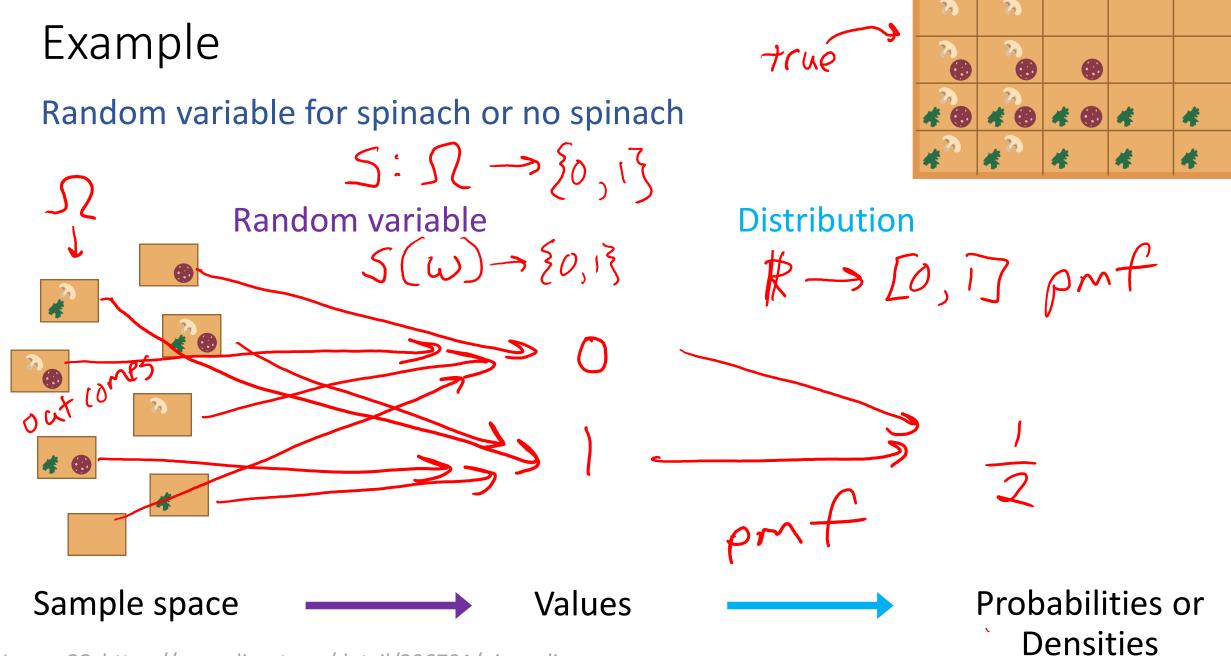
Calculate the probability of these event sequences happening

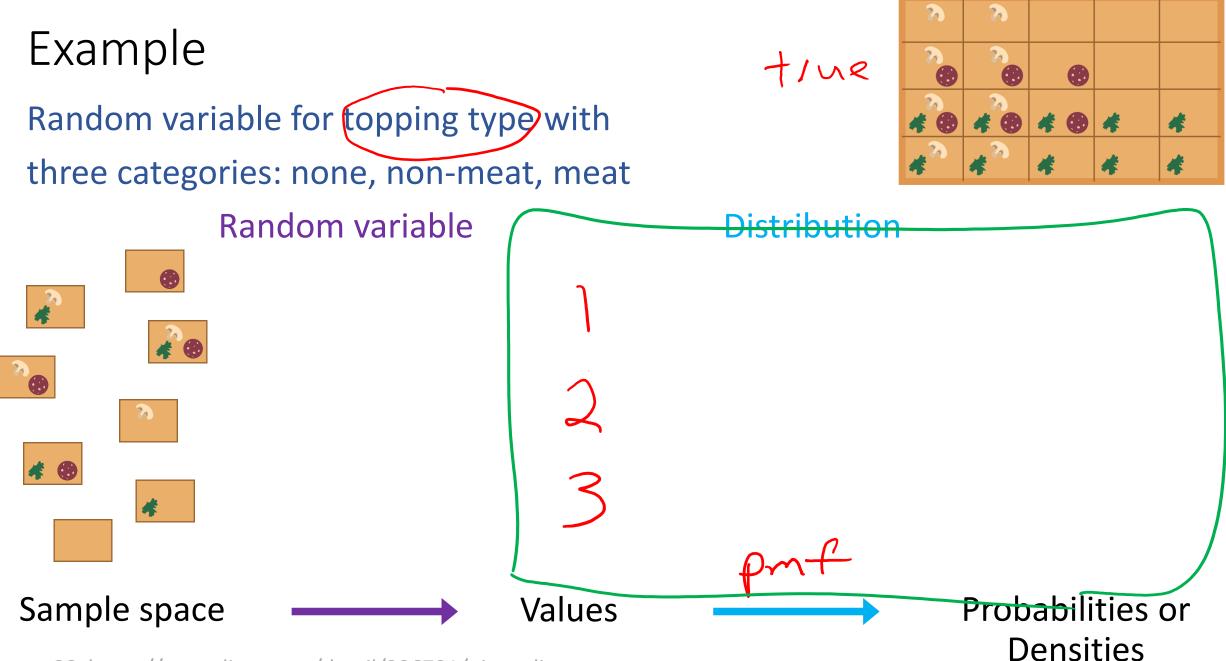
Coin  $\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{1}$ <u> </u>Fair: H, H, T, H **b** Biased,  $\phi = \frac{3/4}{4}$  heads: H, H, T, H  $\frac{3}{4} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} = \frac{7}{256}$ 4-sided die with sides: A, B, C, D 44-444 **A** Fair: A, B, D, D, A • Weighted,  $[\phi_A, \phi_B, \phi_C, \phi_D] = [1/10, 2/10, 3/10, 4/10]$ A, B, D, D, A  $\frac{1}{10} \cdot \frac{2}{10} \cdot \frac{4}{10} \cdot \frac{4}{10} \cdot \frac{4}{10} \cdot \frac{4}{10}$ 

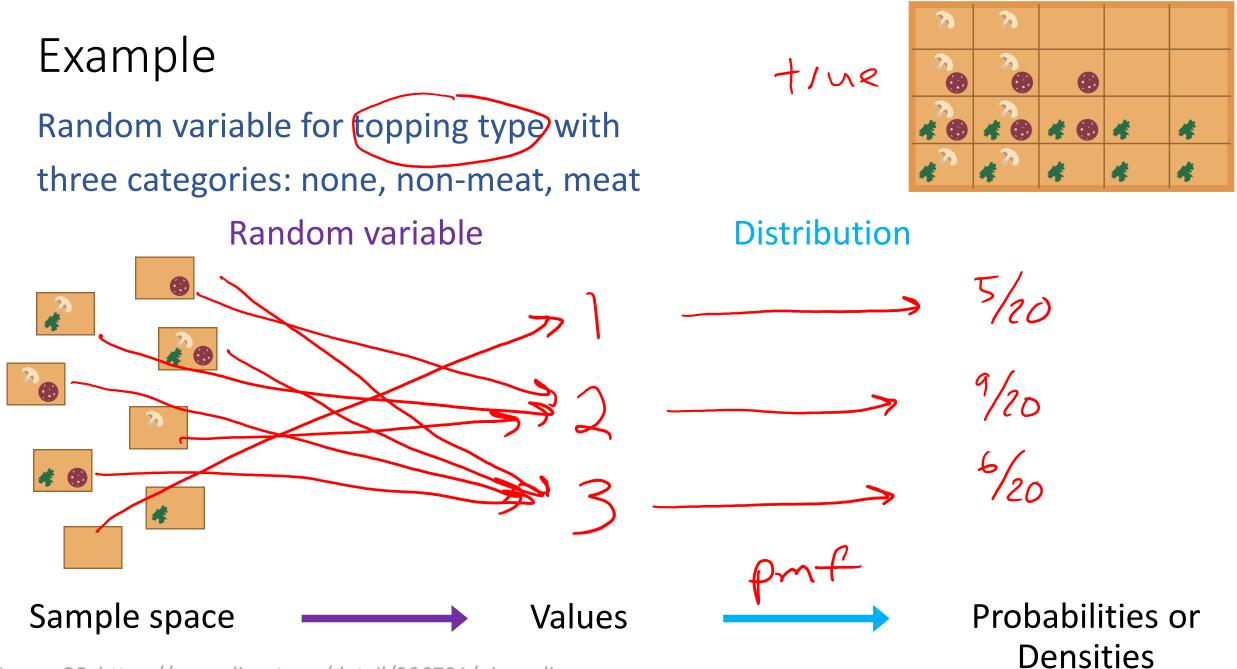
# Probability

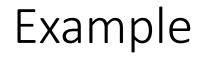
- Probability Vocab
- Outcomes
- Sample space
- Events
- Probability
- Random variable
- Discrete random variable
- Continuous random variable
- **Probability mass function**
- Probability density function
- Parameters







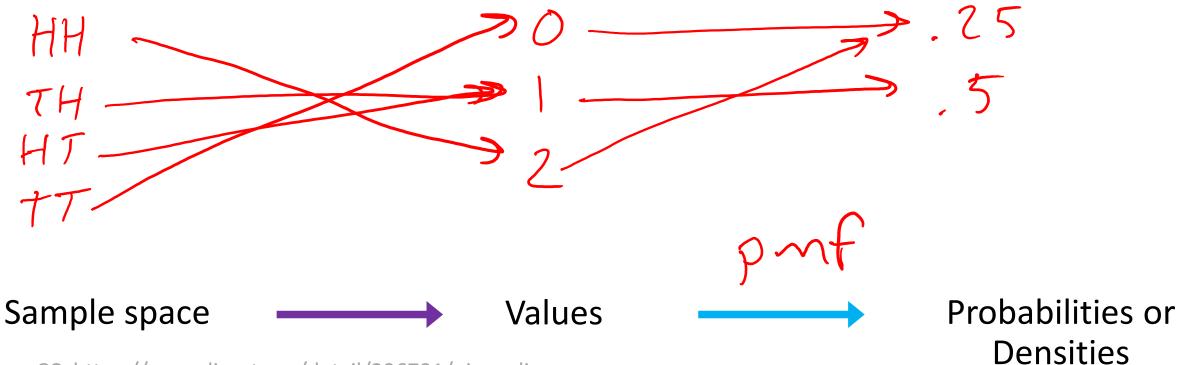




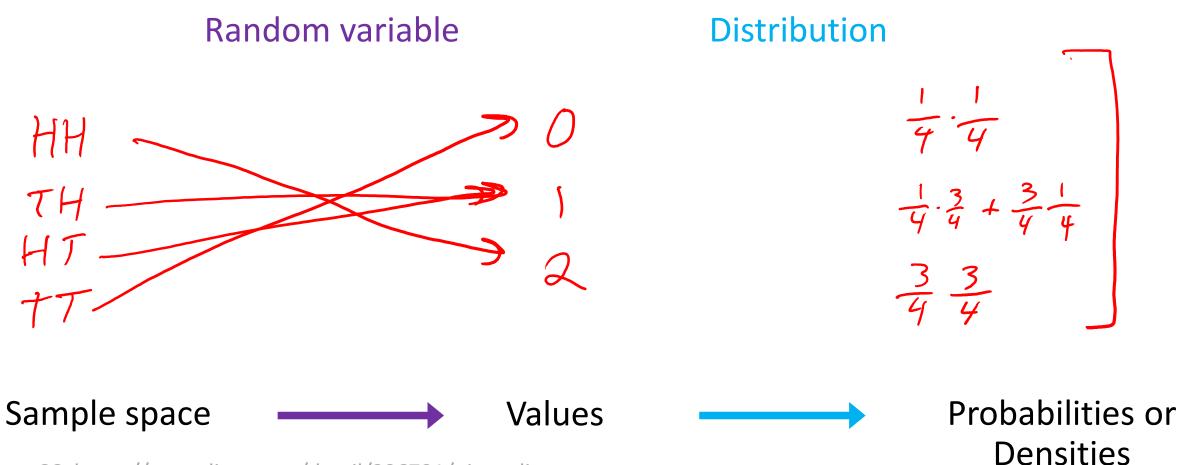
Random variable for number of heads after two flips of a fair coin

Random variable

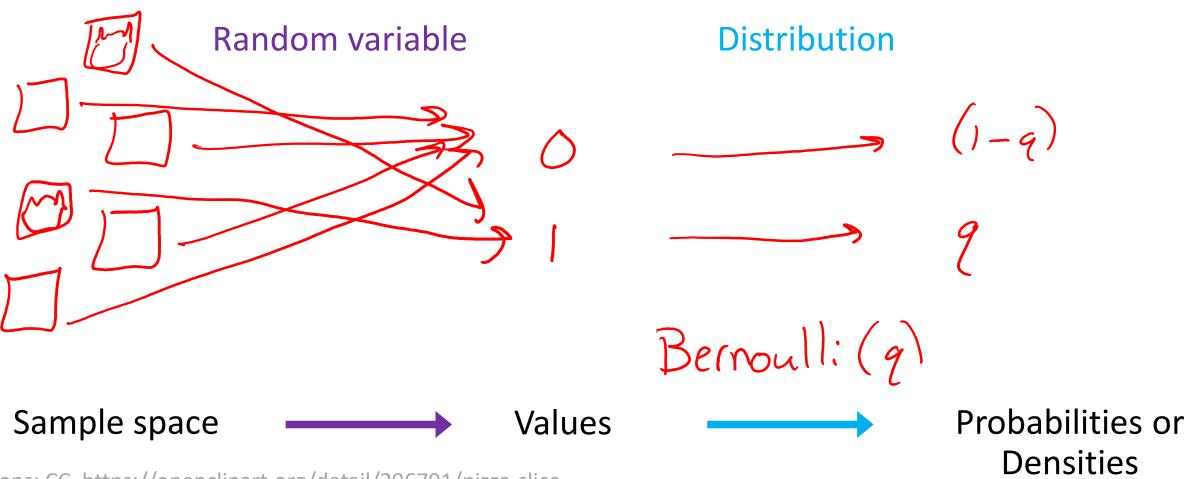
Distribution

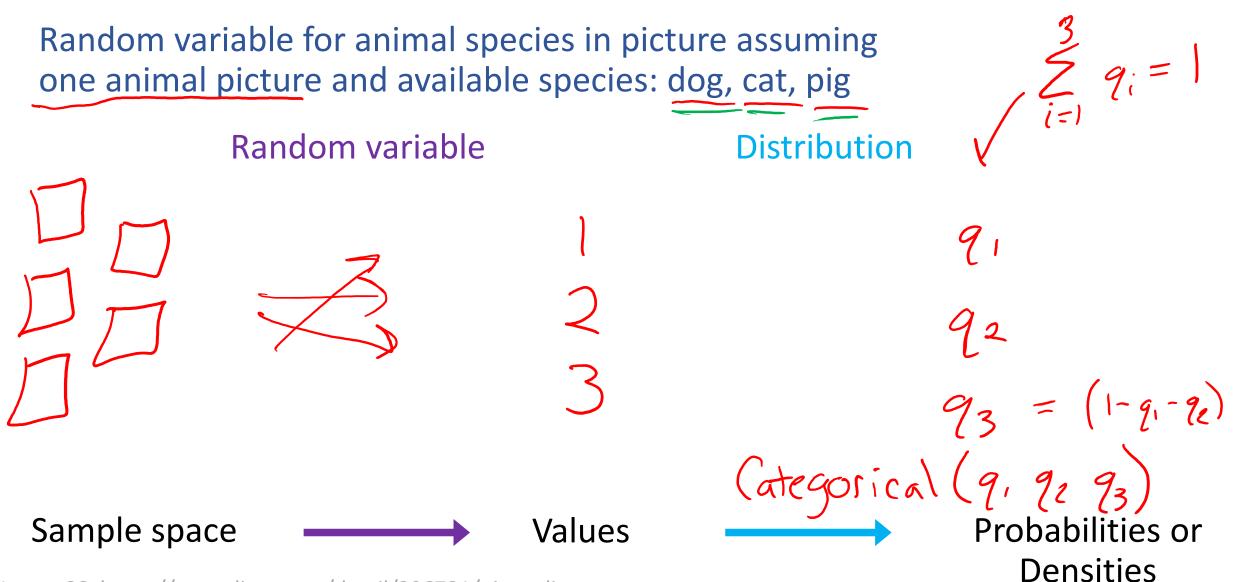


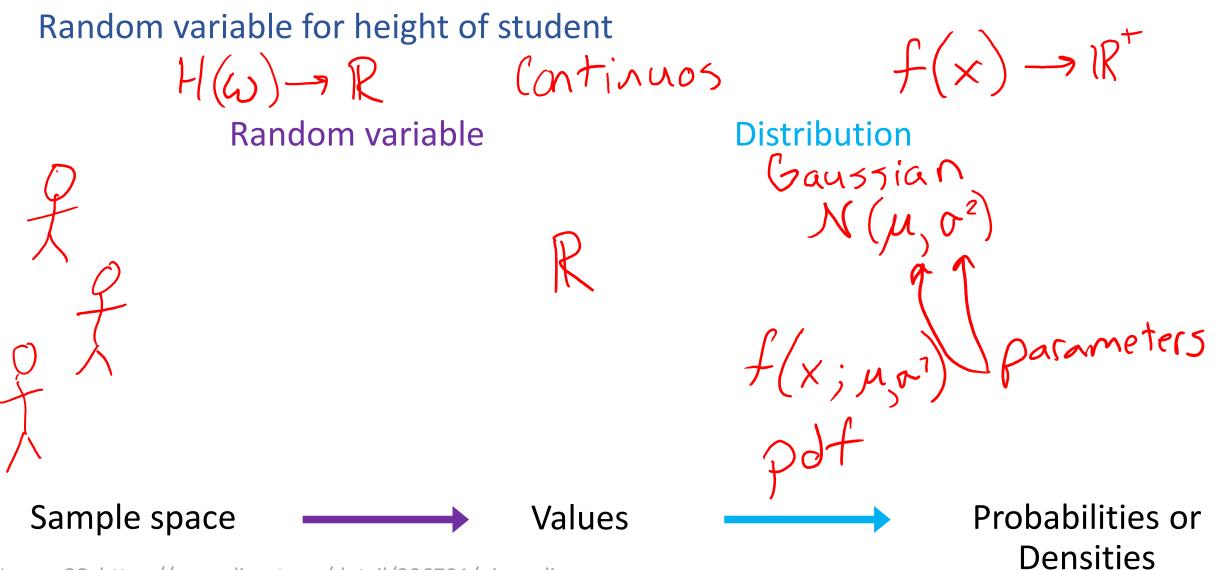
Random variable for number of heads after two flips of a *biased* coin that lands heads 75%



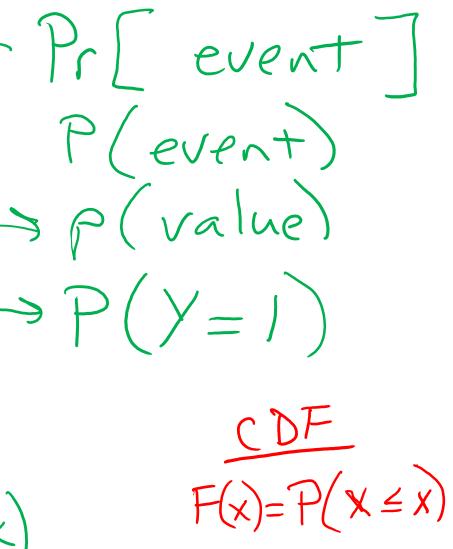
Random variable for cat in picture or not



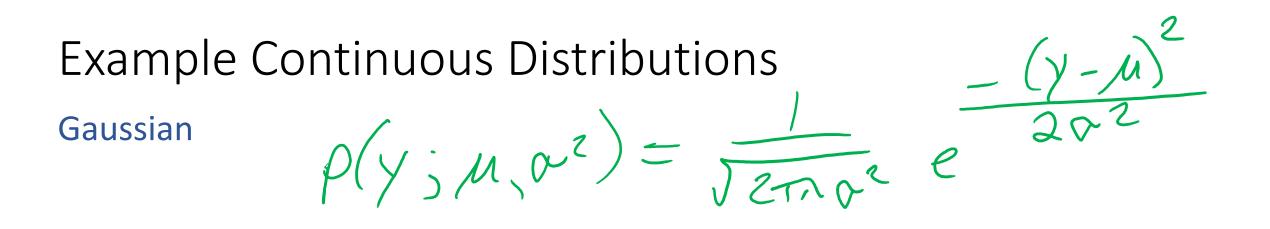




Probability Vocab Outcomes Sample space **Events Probability** Random variable **Discrete random variable** Continuous random variable Probability mass function P(x) = P(x = x)Probability density function f(x)**Parameters** 



Example Discrete Distributions Bern  $P(Y=1) = \phi$ Categorial Multinom Bernoulli  $Y \in \{0, 1\}$   $P(Y=0) = 1 - \phi$ Categorical  $P(Y_2 = 1) =$  $Y_{\chi} \in \mathbb{Z}0, \mathbb{Z}$  $\phi_z = \sum_{k=1}^{\infty}$ **Binomial Multinomial**  $P(Y_{\nu}=)=$ Jniform



Beta

Laplace

 $p(x) = \sum_{y} p(x, y)$  Marginalizing Probability Vocab Marginal p(x, y, z, w)p(x,y)Joint Conditional p(x|y) p(x, v|z, y)

### Notation

#### Dataset

#### Parameters, generically $\theta$

 $p(\mathcal{D} \mid \theta), p(\mathcal{D} ; \theta)$ Random variables

Capital Values

lower case

#### Random variable: function that maps events to values

P(Y=y)

Y is rand variable that maps the event of a coin toss being heads to value one and the event of a coin toss being tails to zero

$$P(Y = 1 | \phi) = 3/4$$
, where  $\phi = 3/4$   
 $P(Y = 1) = 3/4$   
Sometimes even

$$P(Y = heads) = 3/4$$



# Probability Toolbox

- Algebra
- Three axioms of probability
- Theorem of total probability
- Definition of conditional probability
- Product rule
- Bayes' theorem
- Chain rule
- Independence
- Conditional independence

Probability Tools Summary

Adding to the toolbox

1. Definition of conditional probability

3. Bayes' theorem

4. Chain Rule...

 $P(A|B) = \frac{P(A,B)}{P(B)}$  $P(A,B) = P(A \mid B)P(B)$  $P(B|A) = \frac{P(A | B)P(B)}{P(A)}$  $P(A_1, \dots, A_N) = P(A_1) \sum_{i=1}^{N} P(A_i \mid A_{i-1})$ 

# Likelihood

### Likelihood

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .

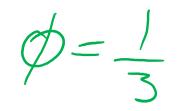
### Likelihood

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .

Grades  

$$\begin{aligned}
\mathcal{J} &= \left\{ \begin{array}{c} \mathcal{Y} & (1) & (2) & (3) \\ \mathcal{Y} & \mathcal{Y} & \mathcal{Y} & \mathcal{Y} \\ \end{array} \right\} \\
&= \left\{ \begin{array}{c} \mathcal{R} & \mathcal{Q} & \mathcal{Q} & \mathcal{R} \\ \mathcal{Q} & \mathcal{Q} & \mathcal{Q} & \mathcal{R} \\ \mathcal{Q} & \mathcal{Q} & \mathcal{Q} & \mathcal{R} \\ \end{array} \right\} \\
&= \left\{ \begin{array}{c} \mathcal{R} & \mathcal{Q} & \mathcal{Q} & \mathcal{R} \\ \end{array} \right\} \\
&= \left\{ \begin{array}{c} \mathcal{R} & \mathcal{R} & \mathcal{R} \\ \mathcal{Q} & \mathcal{Q} \\ \mathcal{Q} & \mathcal{R} \\ \mathcal{Q} & \mathcal{Q} \\ \mathcal{Q} \\$$





### Trick coin: comes up heads only 1/3 of the time

1 flip: H probability:  $\frac{1}{3}$ 

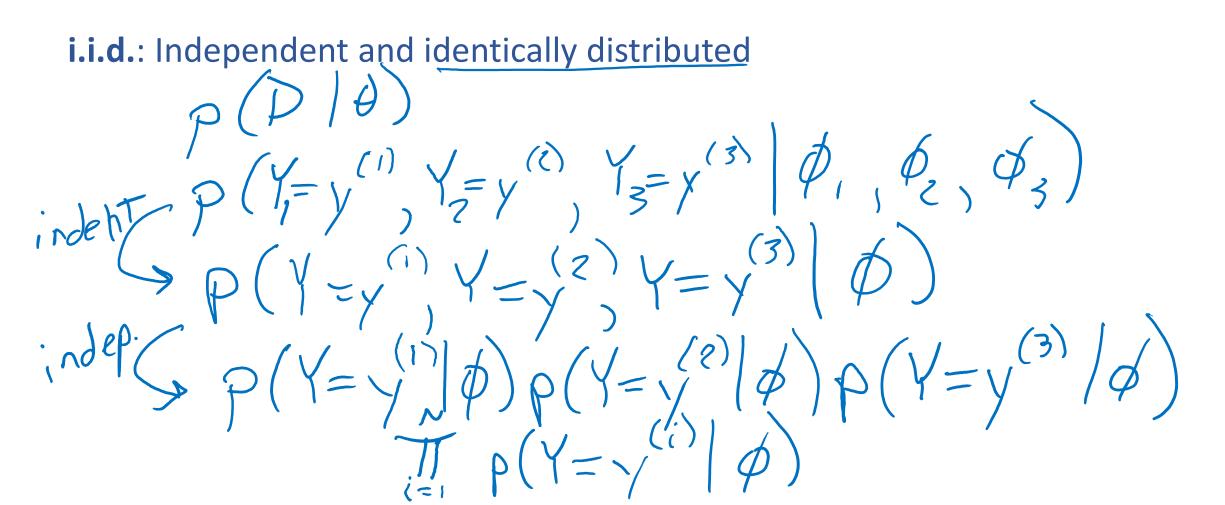
2 flips: H,H probability:  $\frac{1}{3} \cdot \frac{1}{3}$ 

3 flips: H,H,T probability:  $\frac{1}{3} \cdot \frac{1}{3} \cdot \left(1 - \frac{1}{3}\right)$ 

But why can we just multiply these?

### Likelihood and i.i.d

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .



### Bernoulli Likelihood

Bernoulli distribution:

$$Y \sim Bern(\phi) \qquad p(y \mid \phi) = \begin{cases} \phi, & y = 1\\ 1 - \phi, & y = 0 \end{cases}$$

What is the likelihood for three i.i.d. samples, given parameter  $\phi$ :

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\} \qquad (\gamma^{(i)} = 1)$$
  

$$\prod_{i=1}^{N} p(Y = y^{(i)} | \phi)$$
  

$$= \underline{\phi} \cdot \phi \cdot (1 - \phi) \qquad (1 - \phi)^{\circ} \qquad (1 - \phi)^{\circ} \qquad (1 - \phi)^{\circ}$$

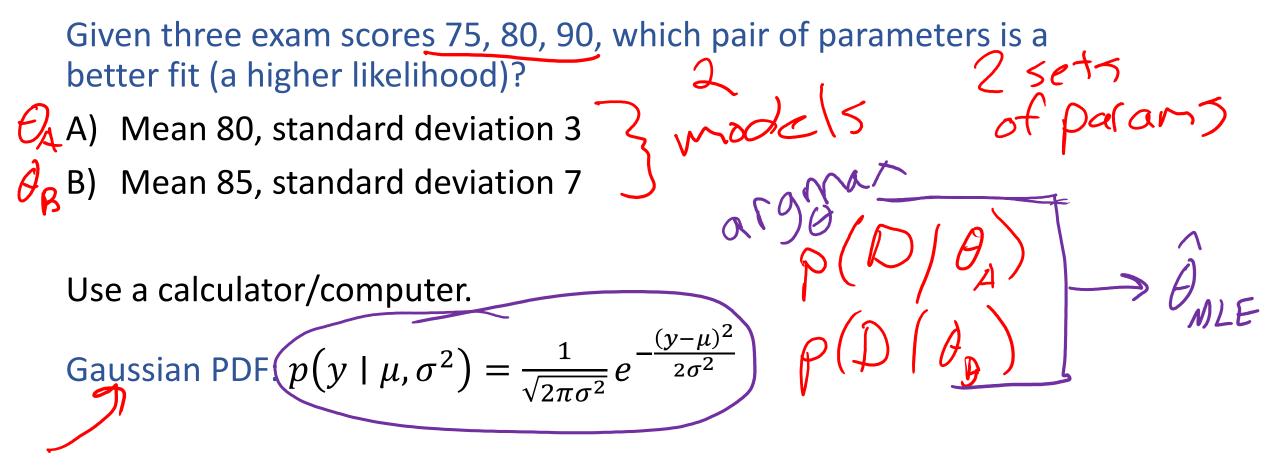
### MLE Maximum likelihood estimation

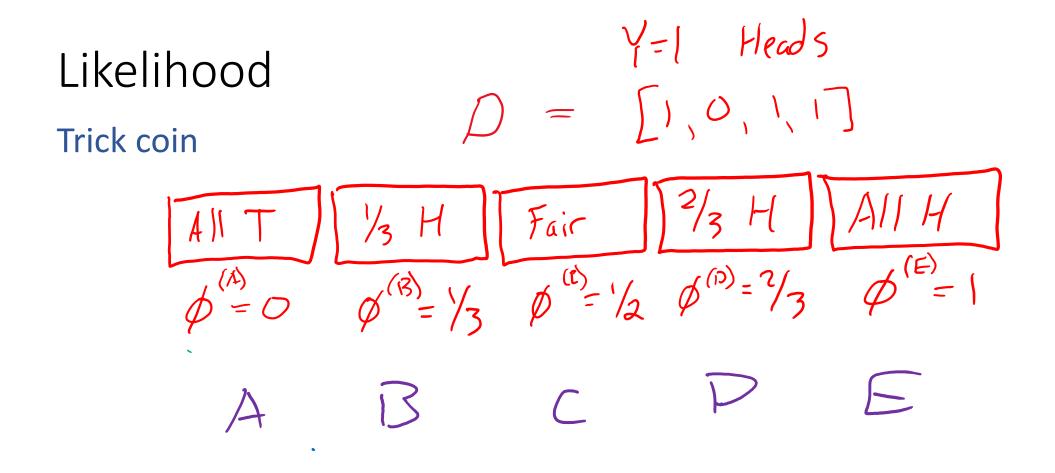
From Probability to Statistics

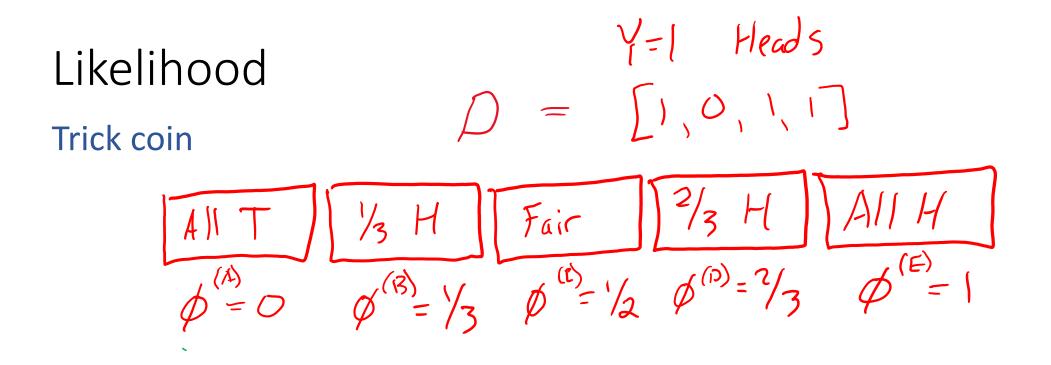
 $\phi \rightarrow \rho$   $D \rightarrow ($ \$ predict

### Poll 2

Assume that exam scores are drawn independently from the same Gaussian (Normal) distribution.







### Estimating Parameters with Likelihood

We model the outcome of a single mysterious weighted-coin flip as a Bernoulli random variable:

$$p(y \mid \phi) = \begin{cases} \varphi \sim Bern(\phi) \\ \phi, & y = 1 \ (Heads) \\ 1 - \phi, & y = 0 \ (Tails) \end{cases}$$

Given the ordered sequence of coin flip outcomes: [1, 0, 1, 1]

What is the estimate of parameter  $\widehat{\phi}$ ?

### Estimating Parameters with Likelihood

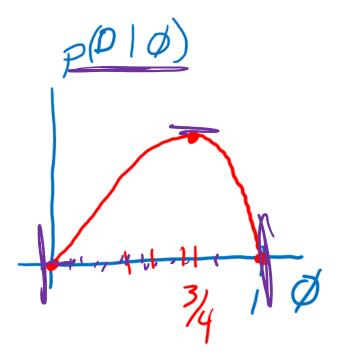
We model the outcome of a single mysterious weighted-coin flip as a Bernoulli random variable:

$$p(y \mid \phi) = \begin{cases} \phi, & y = 1 \ (Heads) \\ 1 - \phi, & y = 0 \ (Tails) \end{cases}$$

Given the ordered sequence of coin flip outcomes: [1, 0, 1, 1]

What is the estimate of parameter  $\widehat{\phi}$ ?

$$p(D \mid \phi) = \phi \cdot \phi \cdot (1 - \phi) \cdot \phi$$
$$= \phi^{3}(1 - \phi)^{1}$$
$$\mathcal{N}_{\tau} = 300 \qquad \mathcal{N}_{\tau} = (00)$$



Likelihood and Maximum Likelihood Estimation

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .

Likelihood function: The value of likelihood as we change theta

(same as likelihood, but conceptually we are considering many different values of the parameters)

$$\mathcal{I}(\mathcal{J}; \mathcal{D}) = \mathcal{P}(\mathcal{D}|\mathcal{A})$$

**Maximum Likelihood Estimation (MLE)**: Find the parameter value that maximizes the likelihood.

### MLE as Data Increases

Given the ordered sequence of coin flip outcomes:

$$p(\mathcal{D} \mid \phi) = \prod_{i}^{N} p(y^{(i)} \mid \phi) = \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}}$$

What happens as we flip more coins?

## MLE for Gaussian

#### Gaussian distribution:

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$
$$p(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

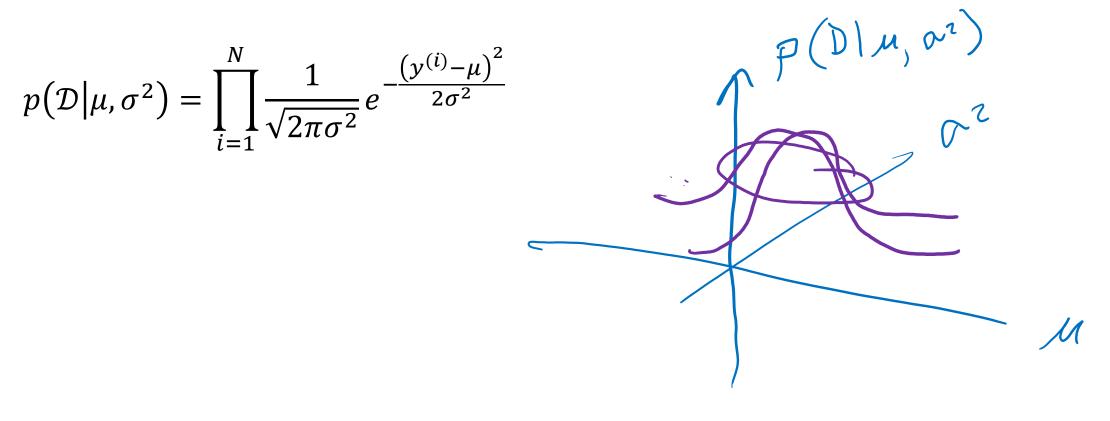
What is the log likelihood for three i.i.d. samples, given parameters  $\mu$ ,  $\sigma^2$ ?  $\mathcal{D} = \{y^{(1)} = 65, y^{(2)} = 95, y^{(3)} = 85\}$ 

$$L(\mu,\sigma^{2}) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{(y^{(i)}-\mu)^{2}}{2\sigma^{2}}} \qquad \qquad \hat{\theta}_{MLE} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i}^{N} p(y^{(i)} \mid \boldsymbol{\theta})$$

### MLE for Gaussian

Assume that exam scores are drawn independently from the same Gaussian (Normal) distribution.

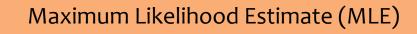
Given three exam scores 75, 80, 90, which pair of parameters is the best fit (the highest likelihood)?

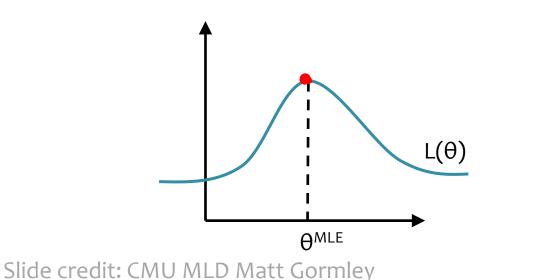


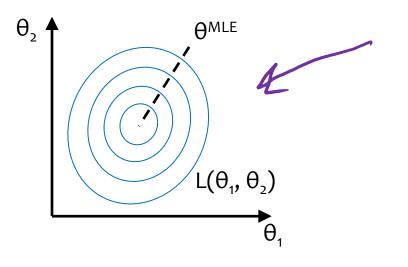
# MLE

Suppose we have data  $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$ 

Principle of Maximum Likelihood Estimation: Choose the parameters that maximize the likelihood of the data.  $\theta^{\mathsf{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$ 







## Likelihood and Log Likelihood

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .

 $\log XZ = \log x + \log z$ 

**Likelihood function**: The value of likelihood as we change theta (same as likelihood, but conceptually we are considering many different values of the parameters)

## Likelihood and Log Likelihood

**Likelihood**: The probability (or density) of random variable Y taking on value y given the distribution parameters,  $\theta$ .

$$P(D|\theta) = \pi p(y^{(i)}|\theta)$$

 $\log XZ = \log x + \log z$ 

Likelihood function: The value of likelihood as we change theta

(same as likelihood, but conceptually we are considering many different values of the parameters)

$$likelihood \mathcal{L}(\Theta; D) = P(D|\Theta) = TP(y^{(i)}|\Theta)$$

$$log likelihood \mathcal{L}(\Theta; D) = log P(D|\Theta) = Z log P(y^{(i)}|\Theta)$$

### Maximum Likelihood Estimation

MLE of parameter  $\theta$  for i.i.d. dataset  $\mathcal{D} = \{y^{(i)}\}_{i=1}^{N}$  $\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} p(\mathcal{D} \mid \theta)$ 

.

# Recipe for Estimation

#### MLE

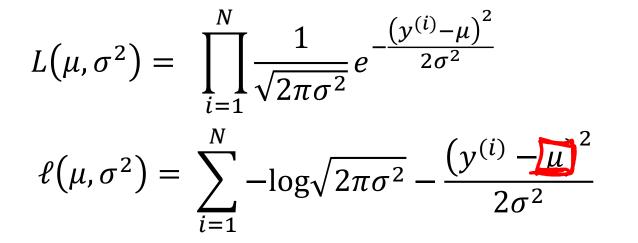
- 1. Formulate the likelihood,  $p(\mathcal{D} \mid \theta)$
- 2. Set objective  $J(\theta)$  equal to negative log of likelihood  $J(\theta) = -\log p(\mathcal{D} \mid \theta)$
- 3. Compute derivative of objective,  $\partial J/\partial \theta$
- 4. Find  $\hat{\theta}$ , either
  - a. Set derivate equal to zero and solve for  $\boldsymbol{\theta}$
  - b. Use (stochastic) gradient descent to step towards better  $\theta$

### MLE for Gaussian

Gaussian distribution:

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$
$$p(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

What is the log likelihood for three i.i.d. samples, given parameters  $\mu$ ,  $\sigma^2$ ?  $\mathcal{D} = \{y^{(1)} = 65, y^{(2)} = 95, y^{(3)} = 85\}$ 



$$\hat{\theta}_{MLE} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i}^{N} p(y^{(i)} \mid \boldsymbol{\theta})$$
$$\hat{\theta}_{MLE} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i}^{N} \log p(y^{(i)} \mid \boldsymbol{\theta})$$

Probabilistic Formulation for ML MLE for Linear and Logistic Regression

### Using Statistics for Machine Learning

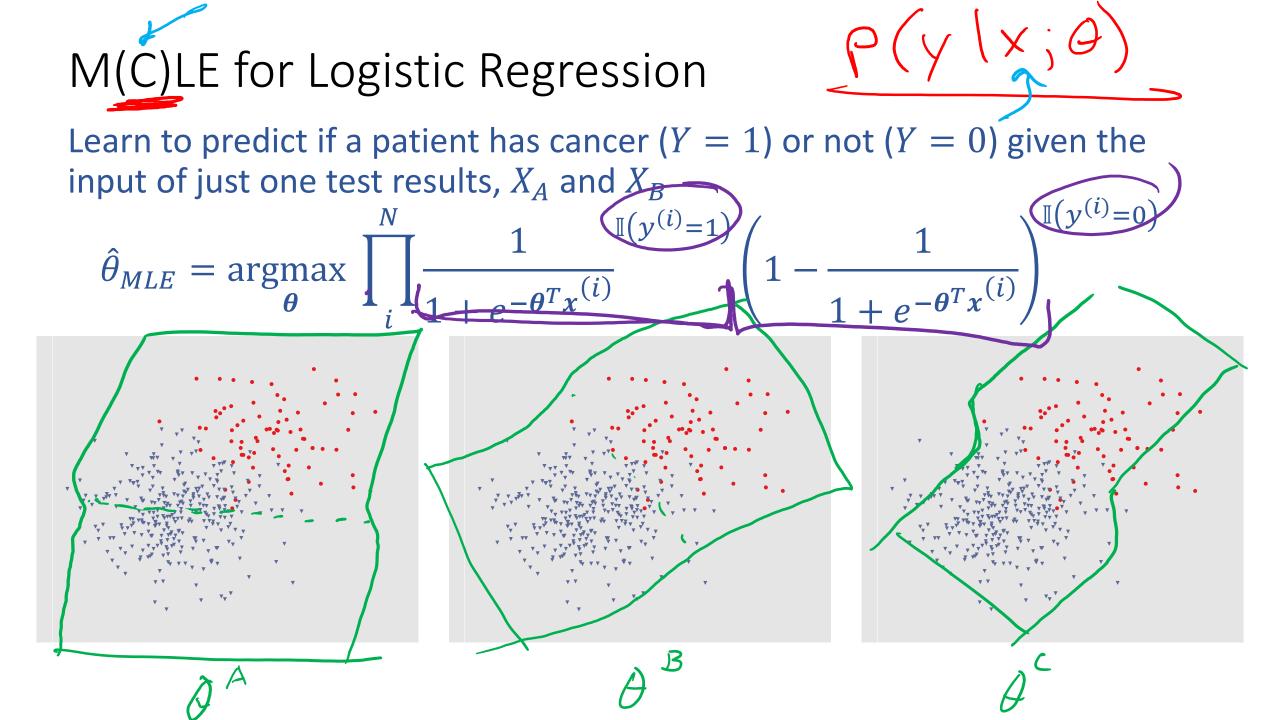
likelihood -> p(D(D)  $TP(Y^{(i)}|\theta)$ 

conditional likelihood  $TTP(Y^{(i)}|X, \theta)$  $P(D|\theta)$ 

# Recipe for Estimation

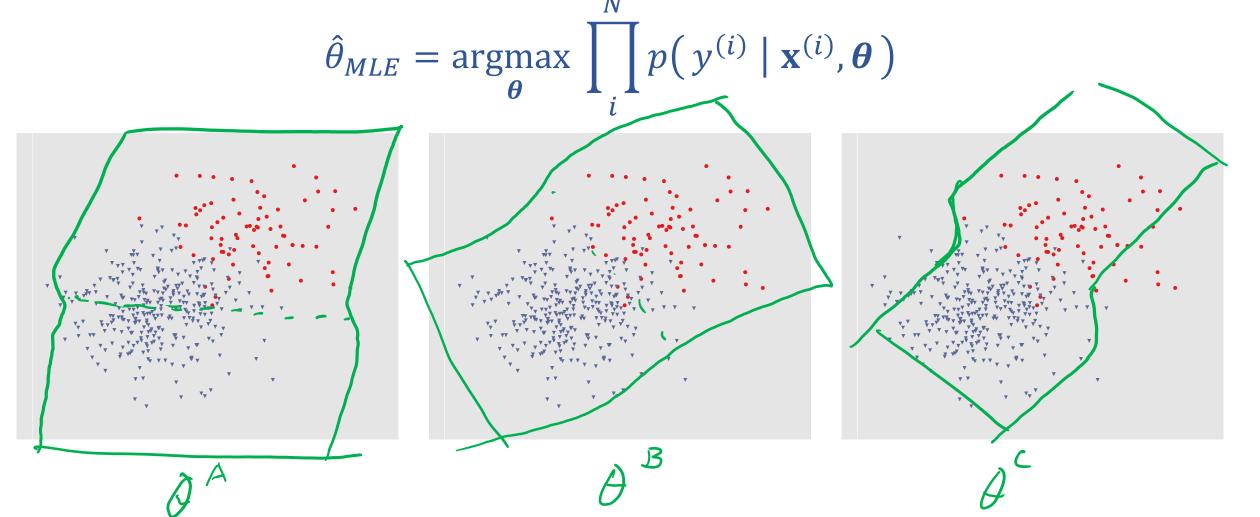
#### MLE

- 1. Formulate the likelihood,  $p(\mathcal{D} \mid \theta)$
- 2. Set objective  $J(\theta)$  equal to negative log of likelihood  $J(\theta) = -\log p(\mathcal{D} \mid \theta)$
- 3. Compute derivative of objective,  $\partial J/\partial \theta$
- 4. Find  $\hat{\theta}$ , either
  - a. Set derivate equal to zero and solve for  $\boldsymbol{\theta}$
  - b. Use (stochastic) gradient descent to step towards better  $\theta$



# M(C)LE for Logistic Regression

Learn to predict if a patient has cancer (Y = 1) or not (Y = 0) given the input of just one test results,  $X_A$  and  $X_B$ 



# M(C)LE for Multi-class Logistic Regression

Learn to predict if probability of output belonging to class k,  $Y_k$ , given input X,  $P(Y_k = 1 | X, \theta_1, ..., \theta_K)$ 

$$\widehat{\Theta}_{MLE} = \underset{\Theta}{\operatorname{argmax}} \prod_{i}^{N} \prod_{k}^{K} \frac{e^{\theta_{k}^{T} x^{(i)}}}{\sum_{l=1}^{K} e^{\theta_{l}^{T} x^{(i)}}} \underbrace{\mathbb{I}(y_{k}^{(i)}=1)}_{\sum_{l=1}^{K} e^{\theta_{l}^{T} x^{(i)}}}$$

M(C)LE for Multi-class Logistic Regression

 $L(\Theta; \mathcal{D})$ 

Learn to predict if probability of output belonging to class k,  $Y_k$ , given input X,  $P(Y_k = 1 \mid X, \theta_1, ..., \theta_K)$  $M = \sum_{k=1}^{N} \sum_{k=1}^{K} \sum_{k=1}^{N} \theta_k^T x^{(i)} = 1$ 

 $\sum_{l=1}^{K} e^{\theta_{l}^{T} x^{(i)}}$ 

 $= \sum_{k} \sum_{k} I(Y_{k}^{(i)}=1) \log_{k}$ 

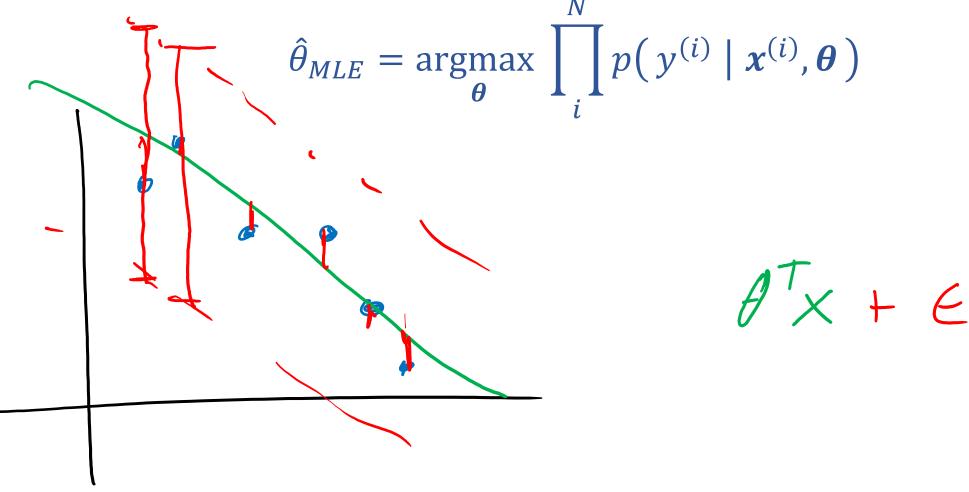
5M

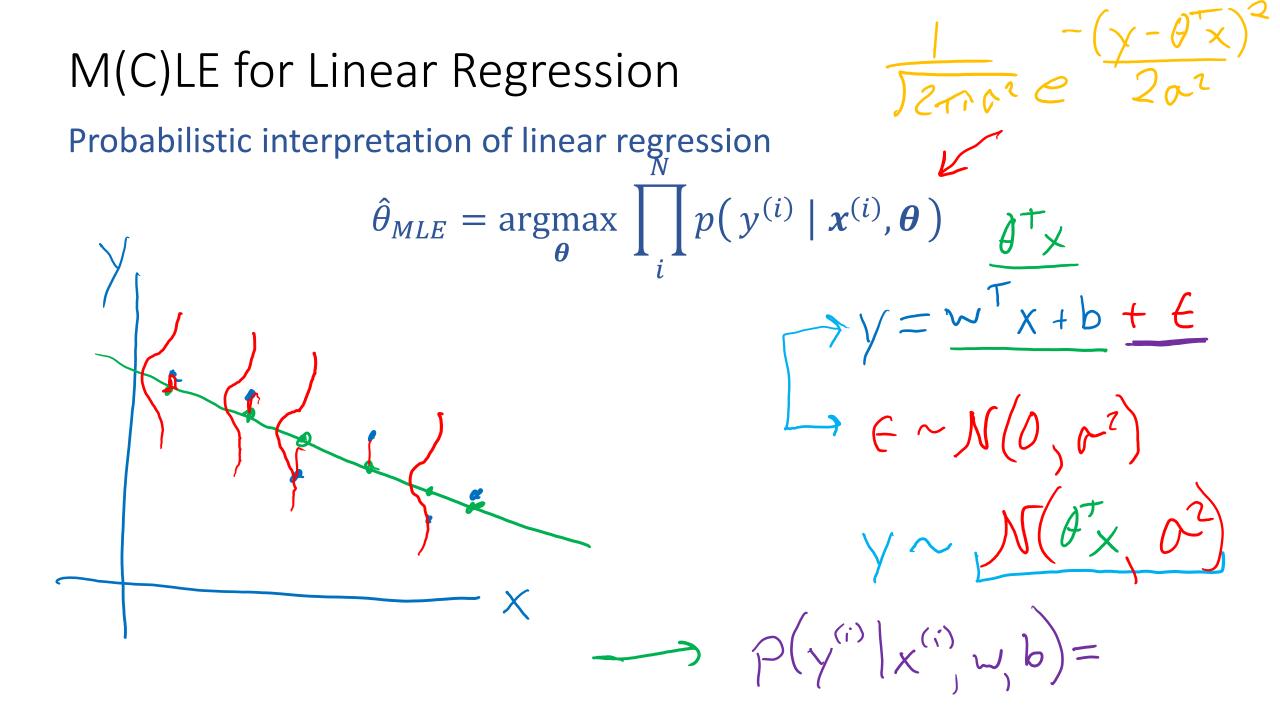
log Ví.

 $y_{k}^{(i)} = 1$ 

## M(C)LE for Linear Regression

Probabilistic interpretation of linear regression





### M(C)LE for Linear Regression

Probabilistic interpretation of linear regression

$$L(\theta; \mathcal{D}) = \prod_{i}^{N} p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

