

10-315 Introduction to ML

MLE and Probabilistic Formulation of Machine Learning

Instructor: Pat Virtue

#### Logistics

#### Exam

Wed in class, see Piazza for details

#### Coming this week

- Probability Primer
- MAP Pre-reading
- Mini project info

#### **Today**

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

#### Course feedback link on Piazza

Are you finished with the course feedback form?

- A. Yes
- B. No
- C. Still working on it

#### Exercises

#### Calculate the probability of these event sequences happening

#### 1. Coin

- a) Fair: {H, H, T, H}
- b) Biased,  $\phi = 3/4$  heads {H, H, T, H}
- 2. 4-sided die with sides: A, B, C, D
  - a) Fair: {A, B, D, D, A}
  - b) Weighted,  $[\phi_A, \phi_B, \phi_C, \phi_D] = [1/10, 2/10, 3/10, 4/10]$  {A, B, D, D, A}

Implement a function in Python for the pdf of a Gaussian distribution.

Python numpy or math packages are fine, no scipy, etc.

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

def gaussian(x, mu, sigmaSq):

What is gaussian(3.3, 2.2, 1.1)?

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 a better fit (a higher likelihood)?

- A) Mean 80, standard deviation 3
- B) Mean 85, standard deviation 7
- C) I don't know

Use a calculator/computer.

Gaussian PDF: 
$$p(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

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 what property allows us to just multiply these?

$$\mathcal{D} = \left\{ y^{(1)}, y^{(2)}, y^{(3)} \right\}$$

Which of these is the correct labeling of these properties?

$$Y \sim Bern(\phi)$$

$$p(\mathcal{D} \mid \theta) =$$

Pre-reading

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

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

Grades

Gaussian PDF: 
$$p(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

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,  $\theta$ .

i.i.d.: Independent and identically distributed

#### 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  $\phi$ :

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

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

$$= \phi \cdot \phi \cdot (1 - \phi)$$

## MLE

Maximum likelihood estimation

## From Probability to Statistics

## Estimating Parameters with Likelihood

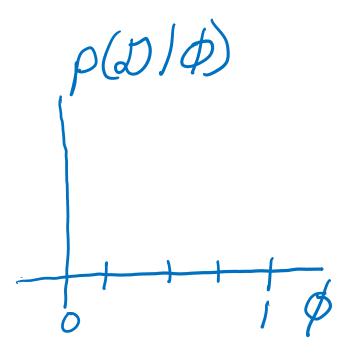
We model the outcome of a single mysterious weighted-coin flip as a Bernoulli random variable:

$$p(y \mid \phi) = \begin{cases} \phi, & y = 1 \text{ (heads)} \\ 1 - \phi, & y = 0 \text{ (tails)} \end{cases}$$

Given the ordered sequence of coin flip outcomes:

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

$$p(D \mid \phi) = \phi \cdot \phi \cdot (1 - \phi) \cdot \phi$$
$$= \phi^{3} (1 - \phi)^{1}$$



https://www.desmos.com/calculator/kr7m2m6cf7

#### 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)

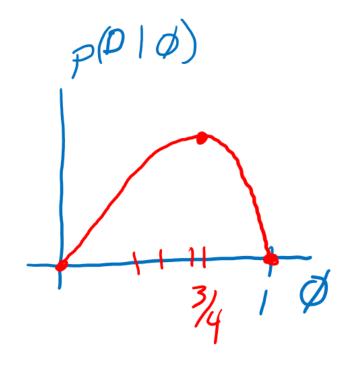
Maximum Likelihood Estimation (MLE): Find the parameter value that maximizes the likelihood.

#### 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:

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

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

$$\mathcal{D} = \{y^{(1)} = 65, y^{(2)} = 95, y^{(3)} = 85\}$$

Formulate the likelihood for three i.i.d. samples, given parameters  $\mu$ ,  $\sigma^2$ ?

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

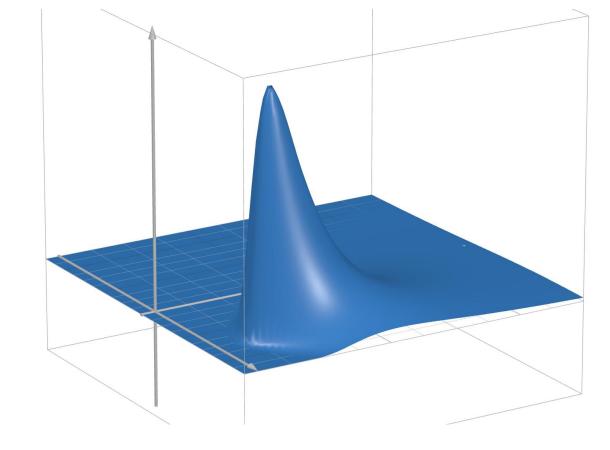
#### MLE for Gaussian

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

Given three exam scores 2, 3, 4, which pair of parameters is the best fit

(the highest likelihood)?

$$p(\mathcal{D}|\mu,\sigma^2) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)}-\mu)^2}{2\sigma^2}}$$



#### MLE

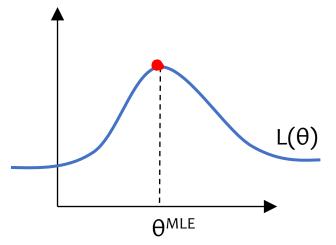
Suppose we have data  $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$ 

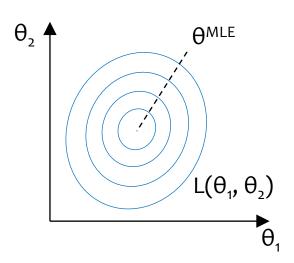
#### Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the likelihood of the data. N

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^{n} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)





## Likelihood and Log Likelihood

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

$$p(\mathcal{D} \mid \theta) = \prod_{i=1}^{N} p(y^{(i)} \mid \theta)$$

Likelihood function: The value of likelihood as we change theta

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

### MLE Objective

Updating our objective: Minimize neg. log likelihood

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^{N} p(y^{(i)} | \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(y^{(i)} | \theta)$$

Log is monotonic:

If 
$$a < b$$

then  $\log a < \log b$ 

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^{N} \log p(y^{(i)} \mid \theta)$$

Minimize 
$$J(\theta) = -\log \mathcal{L}(\theta; \mathcal{D}) = -\sum_{i=1}^{n} \log p(y^{(i)} \mid \theta)$$

#### 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  $\mu$ ,  $\sigma^2$ ?

$$\mathcal{D} = \{y^{(1)} = 75, y^{(2)} = 80, y^{(3)} = 90\}$$

$$L(\mu, \sigma^{2}) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{(y^{(i)} - \mu)^{2}}{2\sigma^{2}}}$$

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

$$\ell(\mu, \sigma^{2}) = \sum_{i=1}^{N} -\log\sqrt{2\pi\sigma^{2}} - \frac{(y^{(i)} - \mu)^{2}}{2\sigma^{2}}$$

$$\hat{\theta}_{MLE} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(y^{(i)} \mid \boldsymbol{\theta})$$

## 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  $\theta$
  - b. Use (stochastic) gradient descent to step towards better  $\theta$

# Probabilistic Formulation for ML MLE for Linear and Logistic Regression

## Using Statistics for Machine Learning

Likelihood vs conditional likelihood

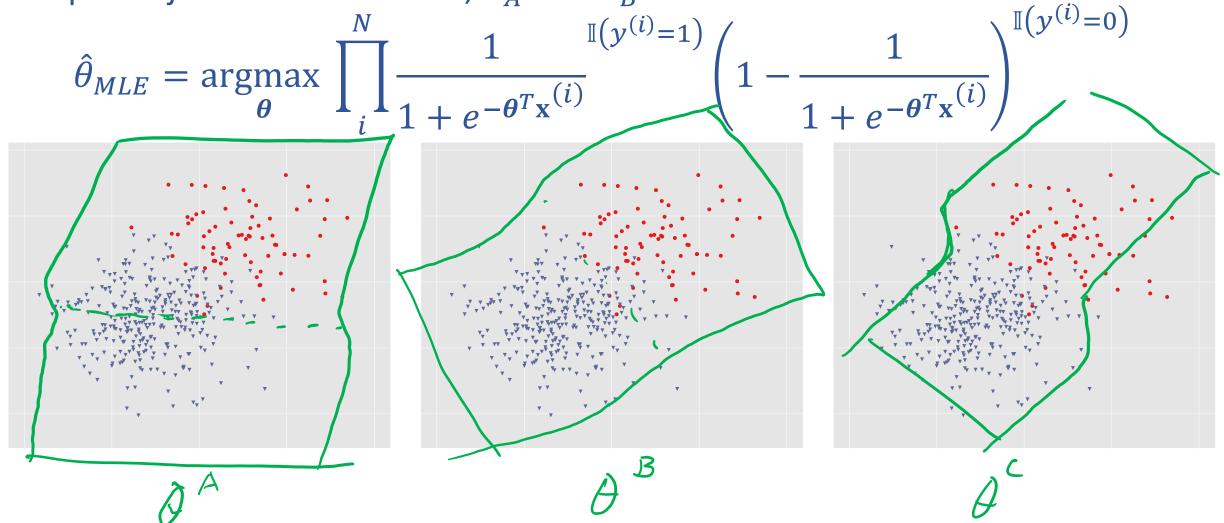
## Recipe for Estimation

#### **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  $\theta$
  - b. Use (stochastic) gradient descent to step towards better  $\theta$

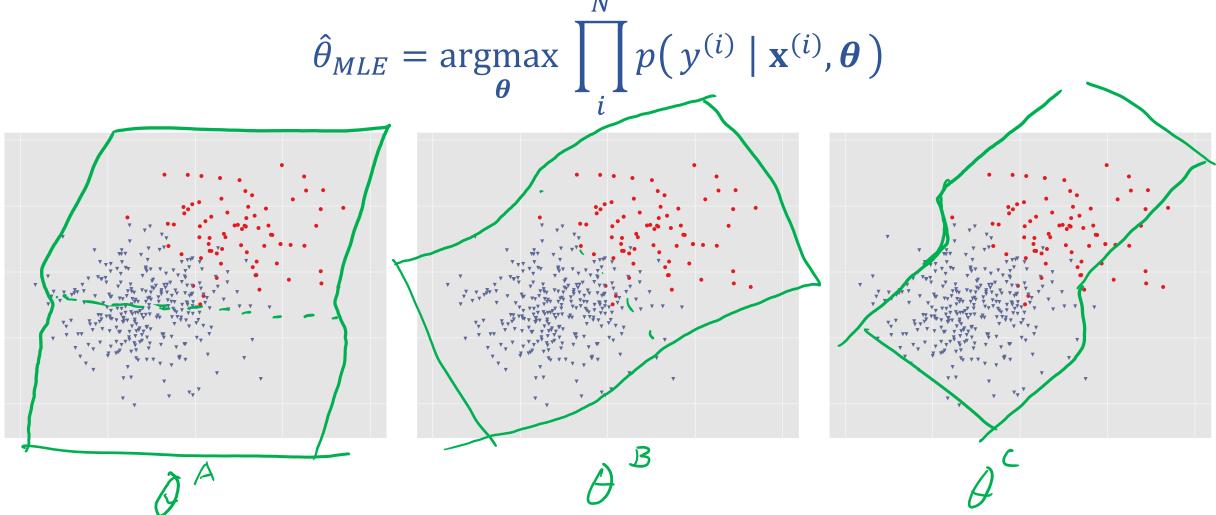
## M(C)LE for Logistic Regression

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 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,  $P(Y_k = 1 \mid X, \boldsymbol{\theta}_1, ..., \boldsymbol{\theta}_K)$ 

$$\widehat{\Theta}_{MLE} = \underset{\mathbf{\Theta}}{\operatorname{argmax}} \prod_{i}^{N} \prod_{k}^{K} \frac{e^{\boldsymbol{\theta}_{k}^{T} \mathbf{x}^{(i)}}}{\sum_{l=1}^{K} e^{\boldsymbol{\theta}_{l}^{T} \mathbf{x}^{(i)}}}$$

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

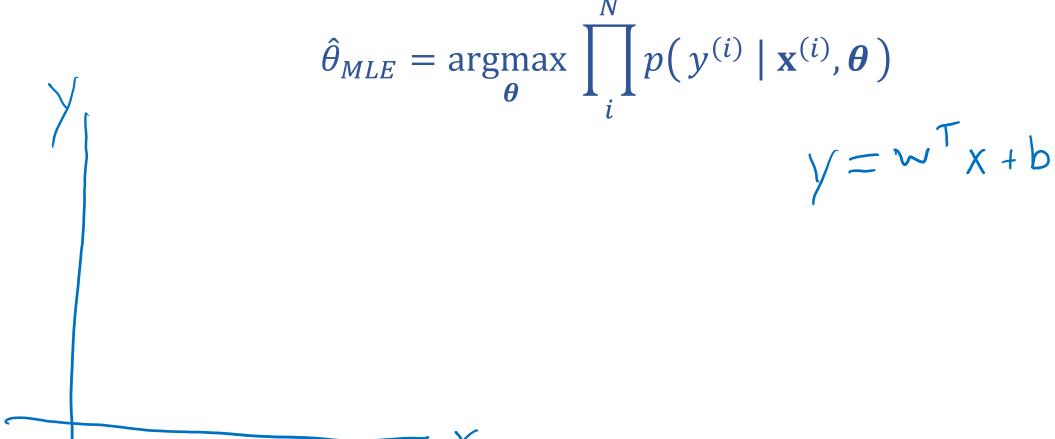
Learn to predict if probability of output belonging to class k,  $Y_k$ , given input X,  $P(Y_k = 1 \mid X, \boldsymbol{\theta}_1, ..., \boldsymbol{\theta}_K)$ 

$$\mathcal{L}(\Theta; \mathcal{D}) = \prod_{i}^{N} \prod_{k}^{K} \frac{e^{\theta_{k}^{T} \mathbf{x}^{(i)}}}{\sum_{l=1}^{K} e^{\theta_{l}^{T} \mathbf{x}^{(i)}}} \mathbb{I}(y_{k}^{(i)} = 1)$$

## M(C)LE for Linear Regression

$$f(z; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

Probabilistic interpretation of linear regression



## M(C)LE for Linear Regression

Probabilistic interpretation of linear regression

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

$$f(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(z-\mu)^2}{2\sigma^2}}$$

$$\prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z^{(i)} - \mu)^2}{2\sigma^2}}$$

$$\sum_{i=1}^{N} -\log\sqrt{2\pi\sigma^2} - \frac{\left(z^{(i)} - \mu\right)^2}{2\sigma^2}$$

## Additional Slides

# Probability Primer

### Probability Vocab

Outcomes

Sample space

**Events** 

**Probability** 

Random variable

Discrete random variable

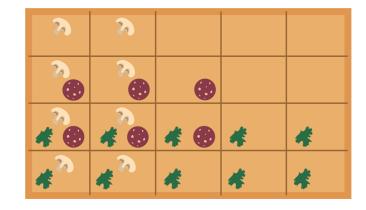
Continuous random variable

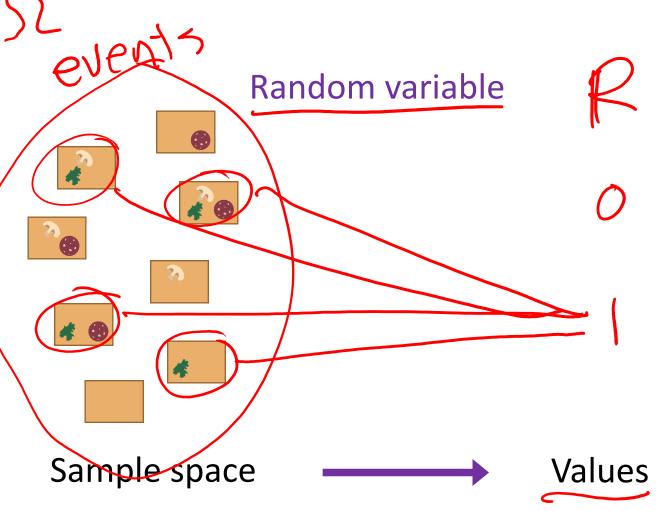
Probability mass function

Probability density function

**Parameters** 

Example  $S \in \{0, 1\}$ Random variable for spinach or no spinach

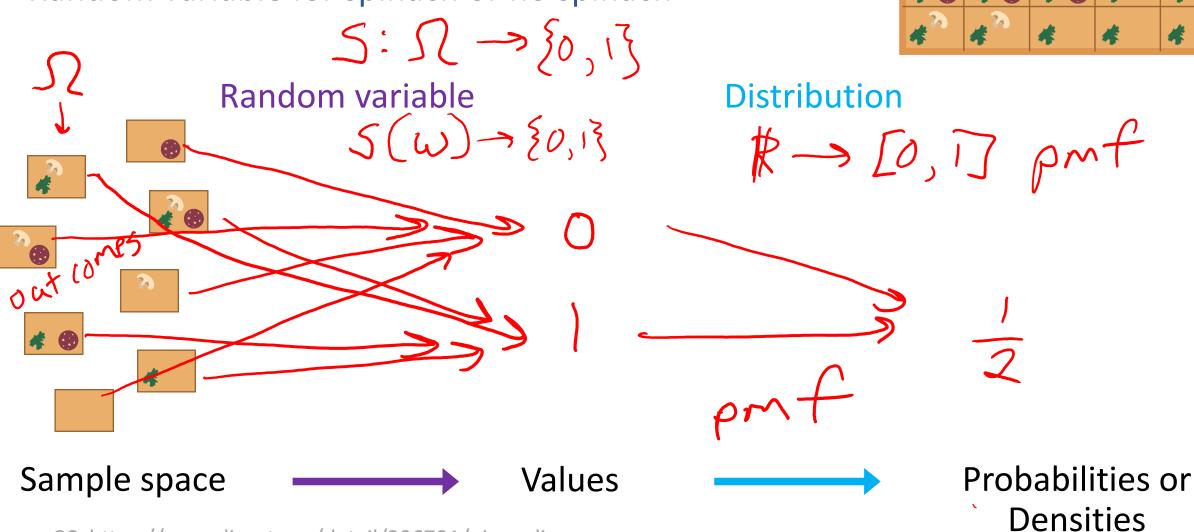




Distribution

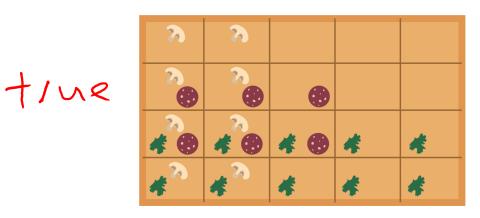
Probabilities or **Densities** 

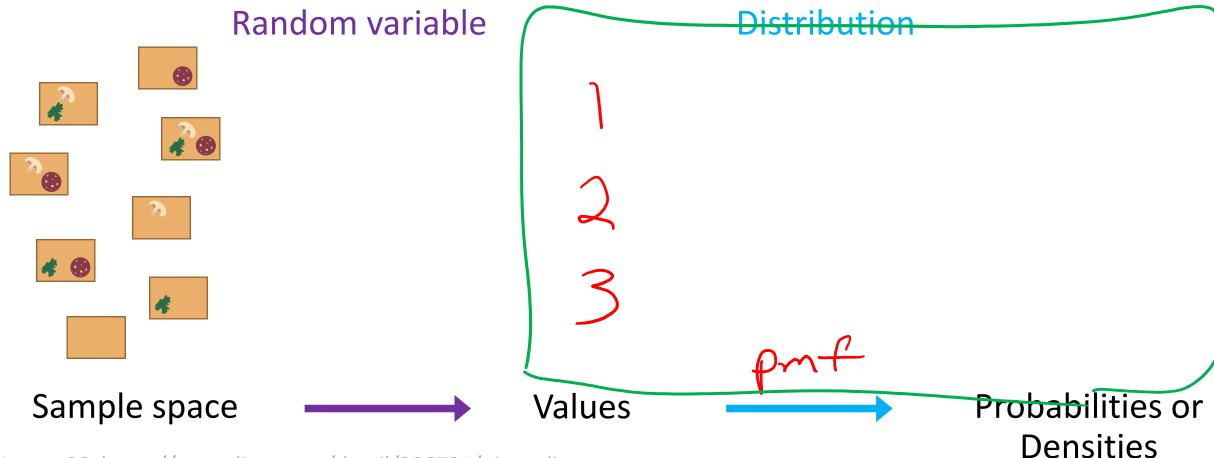
Random variable for spinach or no spinach



Random variable for topping type with

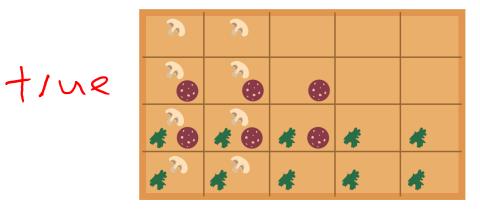
three categories: none, non-meat, meat





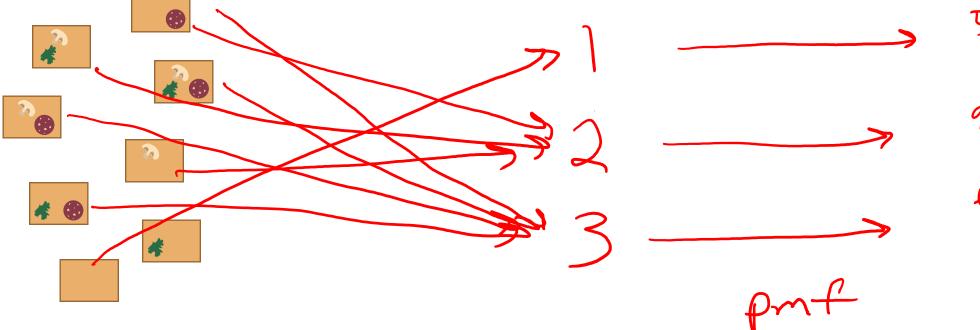
Random variable for topping type with

three categories: none, non-meat, meat



Distribution

#### Random variable



Sample space

**Values** 

Probabilities or **Densities** 

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

Random variable Distribution Probabilities or Sample space Values **Densities** Icons: CC, https://openclipart.org/detail/296791/pizza-slice

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

#### Random variable

## HH TH HT 2

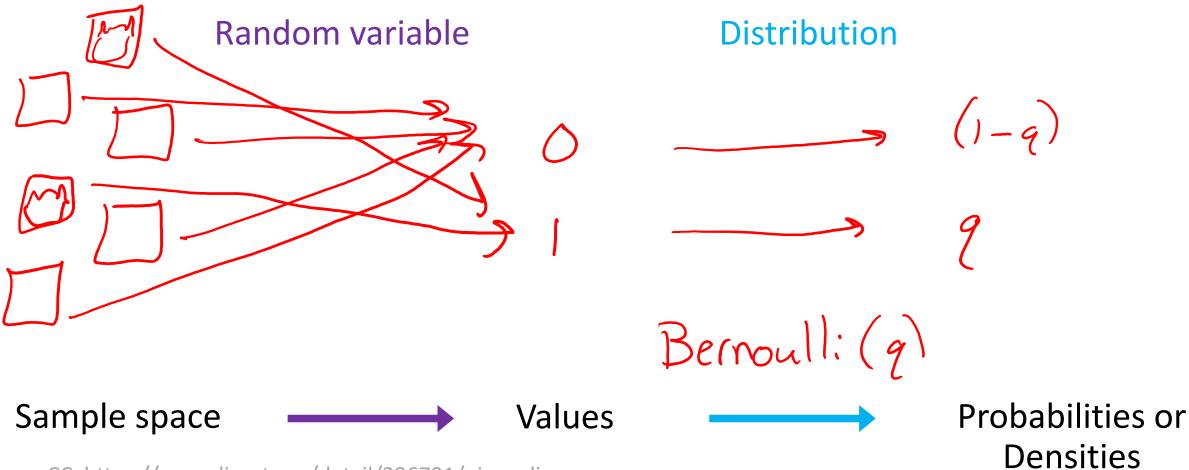
#### Distribution

Sample space

Values

Probabilities or Densities

Random variable for cat in picture or not



Random variable for animal species in picture assuming one animal picture and available species: dog, cat, pig

Distribution Random variable Probabilities or Sample space Values

**Densities** 

Random variable for height of student

Continuos

$$f(x) \rightarrow \mathbb{R}^{+}$$

Random variable

**Distribution** 

Sample space

**Values** 

Probabilities or **Densities** 

## **Probability Vocab**

**Outcomes** 

Sample space

**Events** 

**Probability** 

Random variable

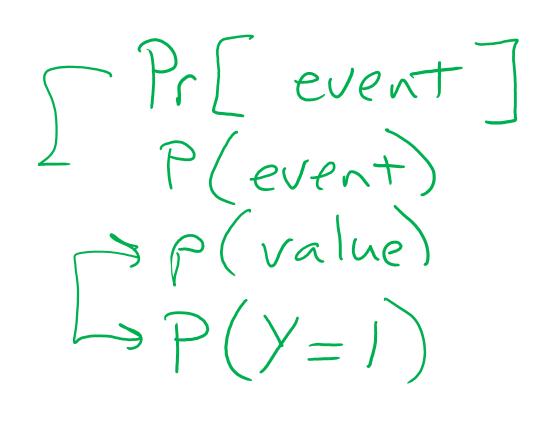
Discrete random variable

Continuous random variable

Probability mass function P(x) = P(X = x)

Probability density function f(x)

**Parameters** 



$$\frac{CDF}{F(x)=P(X \leq x)}$$

## Example Discrete Distributions

Bernoulli

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_K \end{bmatrix}$$

Multinomial

$$\langle -1 \rangle = \overline{p}$$

$$P(Y_2=1)=\phi_2$$

$$P(Y_K = 1) = P_K$$

$$=$$
  $\leq$   $\neq_k$ 

Example Continuous Distributions

Gaussian

 $\rho(y; \mu, \alpha^2) = \frac{1}{\sqrt{2\pi n}\alpha^2}$ 

 $\frac{-(\gamma-\mu)}{2\sigma^2}$ 

Beta

Laplace

#### Probability Vocab

$$\rho(x) = 0$$

Vocab
$$\rho(x) = \sum_{Y} \rho(x, y) \quad Marginalizing$$

$$p(x, y, z, \omega)$$

$$p(x|y)$$
  $p(x, y|z, y)$ 

#### Notation

#### **Dataset**

Parameters, generically  $\theta$ 

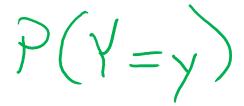
$$p(\mathcal{D} \mid \theta), p(\mathcal{D} ; \theta)$$

Random variables

Capital

Values

lower case





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

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

$$P(A|B) = \frac{P(A,B)}{P(B)}$$

2. Chain Rule Product rule

$$\underline{P(A,B)} = P(A \mid B)P(B)$$

3. Bayes' theorem

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

4. Chain Rule...

$$P(A_1, ..., A_N) = P(A_1) \sum_{i=2}^{N} P(A_i \mid A_{i-1})$$