


# Probability and Statistics

Nakul Gopalan  
Georgia Tech

# Outline

- Probability Distributions ←
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

# Probability

- A sample space  $S$  is the set of all possible outcomes of a conceptual or physical, repeatable experiment. ( $S$  can be finite or infinite.)
  - E.g.,  $S$  may be the set of all possible outcomes of a dice roll:  $S$   
(1 2 3 4 5 6)
  - E.g.,  $S$  may be the set of all possible nucleotides of a DNA site:  $S$   
(A C G T)  

  - E.g.,  $S$  may be the set of all possible time-space positions of an aircraft on a radar screen.
- An Event  $A$  is any subset of  $S$ 
  - Seeing "1" or "6" in a dice roll; observing a "G" at a site; UA007 in space-time interval

# Three Key Ingredients in Probability Theory

A **sample space** is a collection of all possible **outcomes**

Random variables  $X$  represents **outcomes** in sample space

Probability of a random variable to happen  $p(x) = p(X = \underline{x})$

$$\underline{p(x)} \geq 0$$

## Continuous variable

Continuous probability distribution

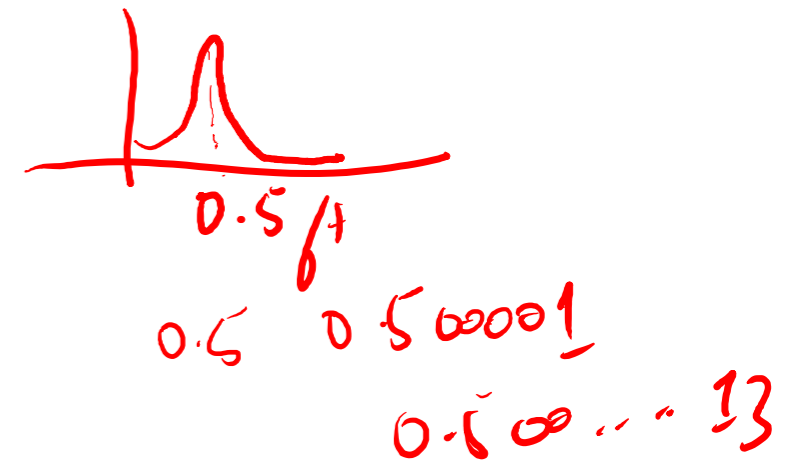
Probability density function

Density or likelihood value

Temperature (real number)

Gaussian Distribution

$$\int_x p(x) dx = 1 //$$



## Discrete variable

Discrete probability distribution

Probability mass function

Probability value

Coin flip (integer)

Bernoulli distribution

$$\sum_{x \in A} p(x) = 1 //$$

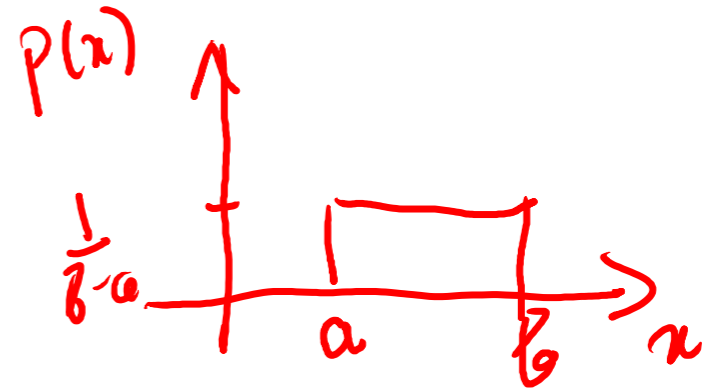


# Continuous Probability Functions

- Examples:

- Uniform Density Function:

$$f_x(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



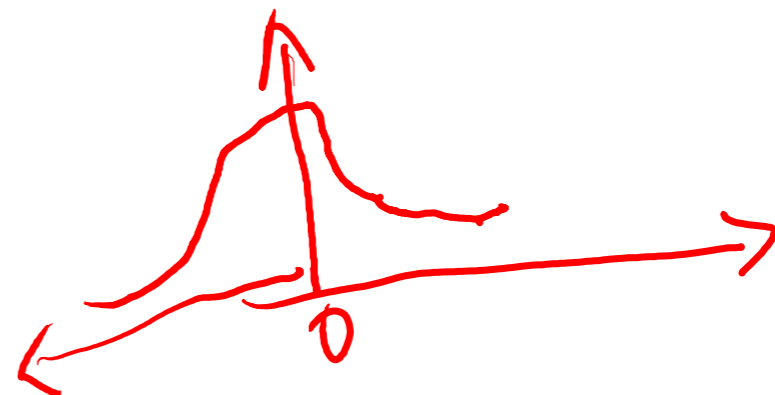
- Exponential Density Function:

$$f_x(x) = \frac{1}{\mu} e^{-\frac{x}{\mu}} \quad \text{for } x \geq 0$$

$$F_x(x) = 1 - e^{-\frac{x}{\mu}} \quad \text{for } x \geq 0$$

- Gaussian(Normal) Density Function

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



# Discrete Probability Functions

- Examples:

- Bernoulli Distribution:

- $$\begin{cases} 1 - p & \text{for } x = 0 \text{ (H)} \\ p & \text{for } x = 1 \text{ (T)} \end{cases}$$

$\downarrow$  0.5

In Bernoulli, just a **single** trial is conducted

- Binomial Distribution:

- $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$

$\rightarrow$  repeated coin tosses

**k** is number of successes

$P(X = k) =$

n  $\rightarrow$  total coin tosses

**n-k** is number of failures

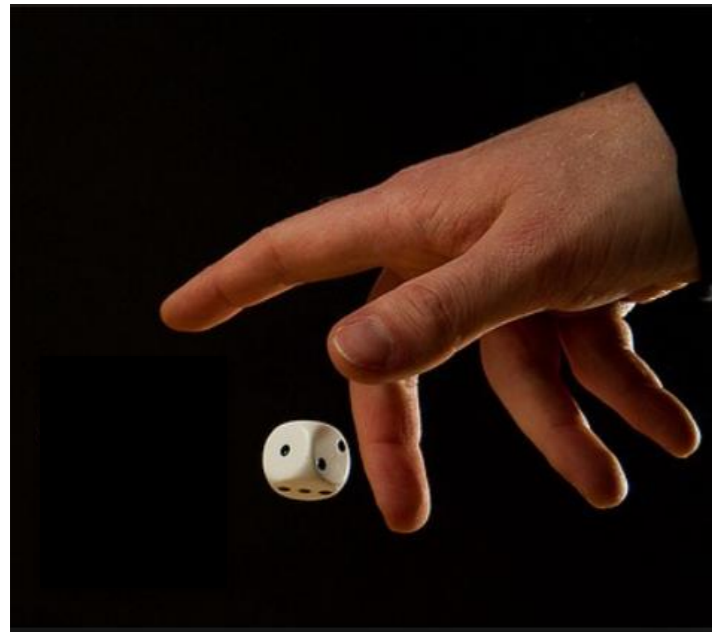
$\binom{n}{k}$  The total number of ways of selection **k** distinct combinations of **n** trials, **irrespective of order**.

# Outline

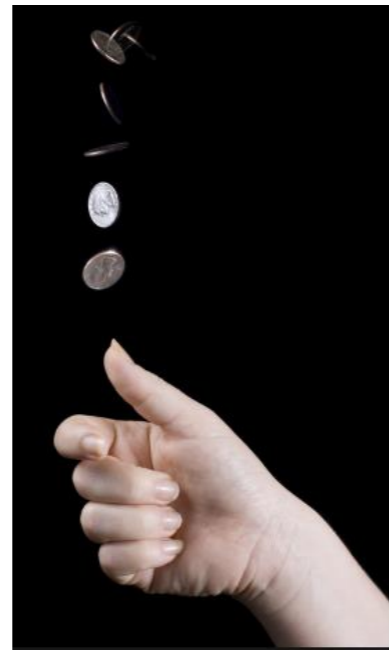
- Probability Distributions
- Joint and Conditional Probability Distributions ←
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation



# Example



X = Throw a dice



Y = Flip a coin

**X** and **Y** are random variables

**N** = total number of trials  
*die roll* (pointing to X)  
*coin toss* (pointing to Y)

$n_{ij}$  = Number of occurrences

		<b>X</b>						
		$x_{i=1} = 1$	$x_{i=2} = 2$	$x_{i=3} = 3$	$x_{i=4} = 4$	$x_{i=5} = 5$	$x_{i=6} = 6$	$C_j$
<b>Y</b>	$y_{j=2} = tail$	$n_{ij} = 3$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 5$	$n_{ij} = 1$	$n_{ij} = 5$	20
	$y_{j=1} = head$	$n_{ij} = 2$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 1$	15
$C_i$		5	6	6	7	5	6	N=35

		<b>X</b>						
		$x_{i=1} = 1$	$x_{i=2} = 2$	$x_{i=3} = 3$	$x_{i=4} = 4$	$x_{i=5} = 5$	$x_{i=6} = 6$	$C_j$
<b>Y</b>	$y_{j=2} = tail$	$n_{ij} = 3$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 5$	$n_{ij} = 1$	$n_{ij} = 5$	20
	$y_{j=1} = head$	$n_{ij} = 2$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 1$	15
$C_i$		5	6	6	7	5	6	N=35

$$\Pr(x=4, y=h) = \frac{2}{35}$$

$$\Pr(x=5) = \frac{\sum_y (n_{i=5})_y}{N} = \frac{(n_{i=5}, y=T) + (n_{i=5}, y=F)}{35} = \frac{5}{35}$$

*Now marginalization*

$$\Pr(y=t \mid x=3) = \frac{2}{2+4} = \frac{2}{6} = \frac{1}{3} = \frac{n_{ij}}{C_i} = \frac{P(y=t, x=3)}{P(x=3)} = \frac{2/35}{6/35} = \frac{2}{6}$$

Joint from conditional

$$\Pr(x=x_i, y=y_i) = P(y=y_i \mid x=x_i) P(x=x_i) = \frac{1}{3} \cdot \frac{2}{25} = \frac{2}{75}$$

**Probability:**

$$p(X = x_i) = \frac{c_i}{N}$$

**Joint probability:**

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

**Conditional probability:**

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

**Sum rule**

$$p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j) \Rightarrow p(X) = \sum_Y P(X, Y)$$

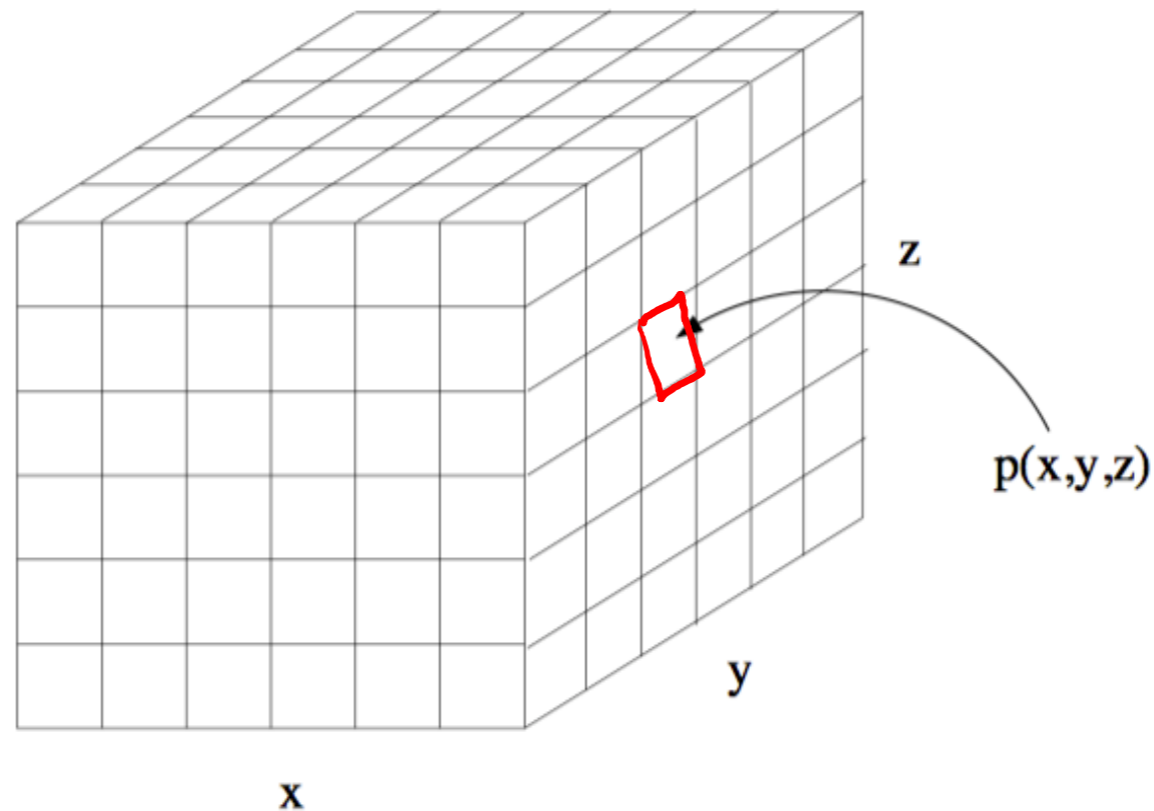
**Product rule**

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \frac{c_i}{N} = p(Y = y_j | X = x_i) p(X = x_i)$$

$$p(X, Y) = p(Y|X)p(X)$$

# Joint Distribution

- Key concept: two or more random variables may interact. Thus, the probability of one taking on a certain value depends on which value(s) the others are taking.
- We call this a joint ensemble and write
$$p(x, y) = \text{prob}(X = x \text{ and } Y = y)$$

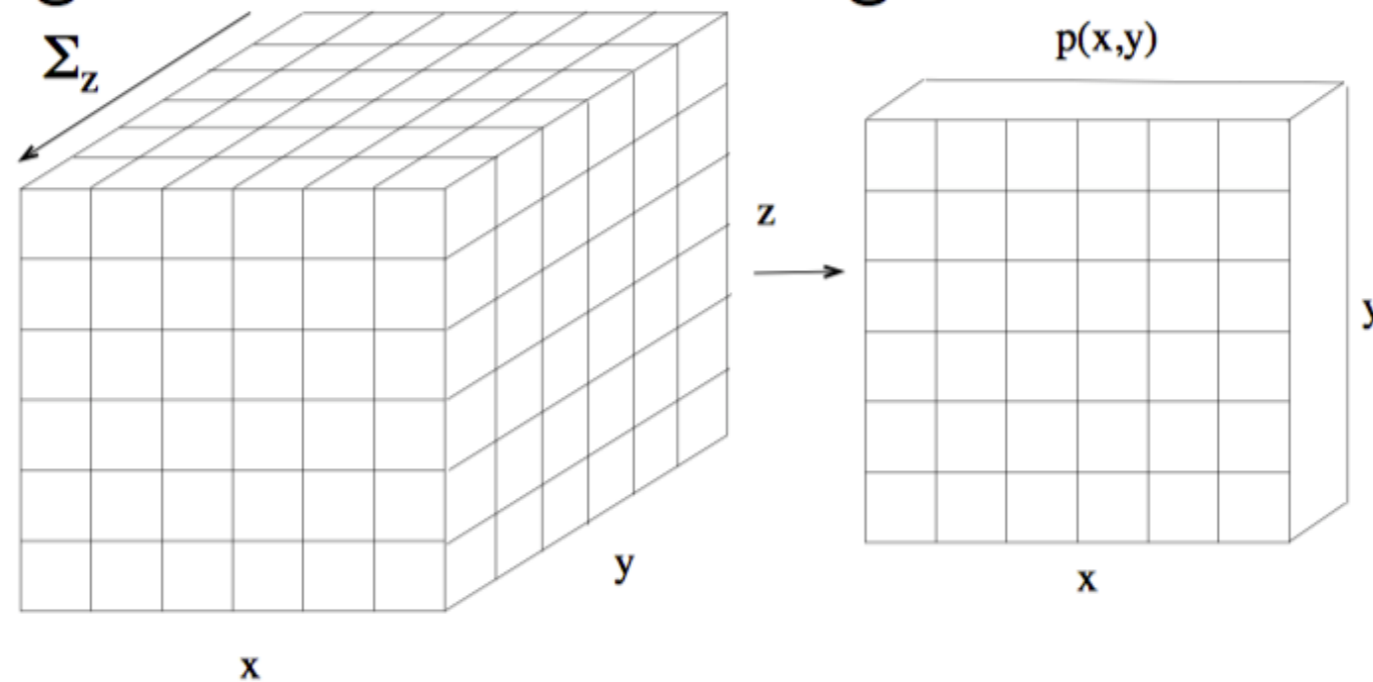


# Marginal Distribution

- We can "sum out" part of a joint distribution to get the *marginal distribution* of a subset of variables:

$$p(x) = \sum_y p(x, y)$$

- This is like adding slices of the table together.

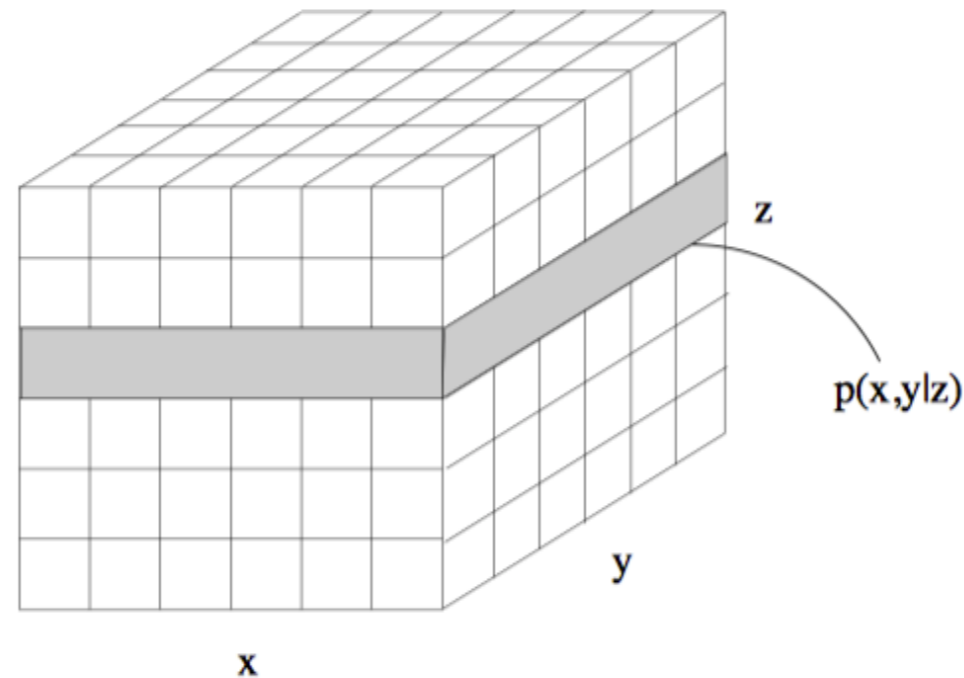


- Another equivalent definition:  $p(x) = \sum_y p(x|y)p(y)$ .

# Conditional Distribution

- If we know that some event has occurred, it changes our belief about the probability of other events.
- This is like taking a "slice" through the joint table.

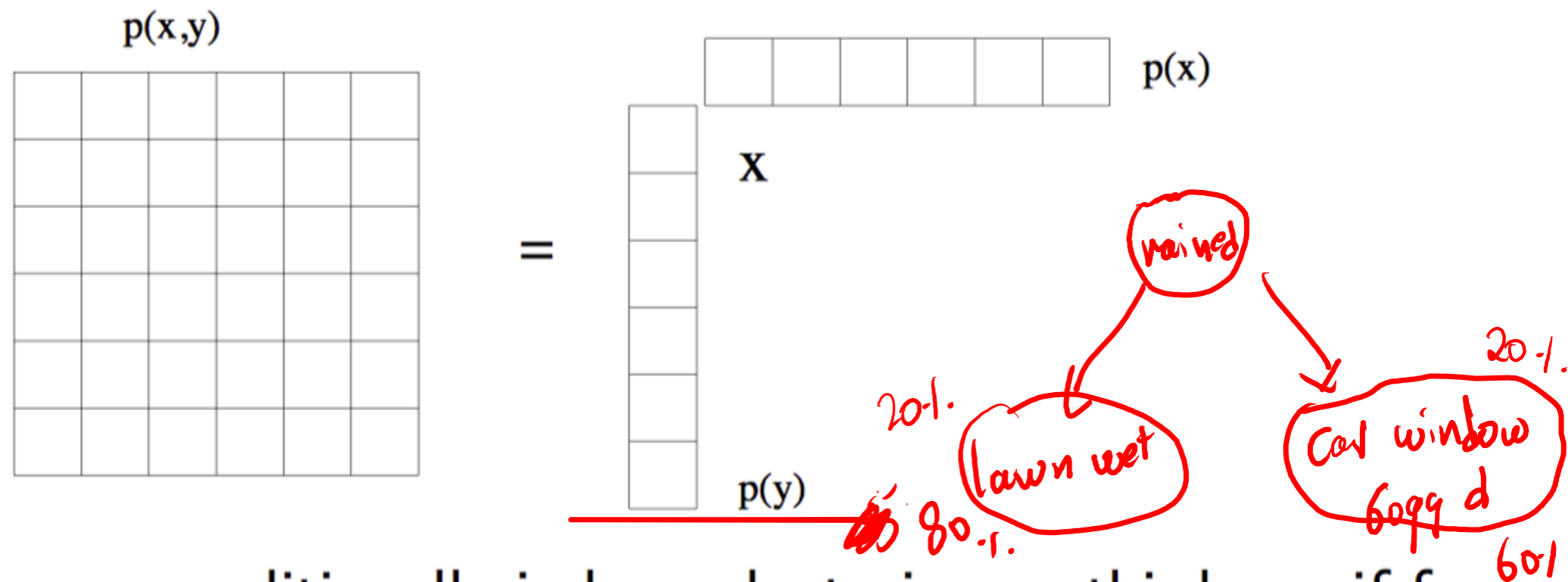
$$p(x|y) = p(x, y) / p(y)$$



# Independence & Conditional Independence

- Two variables are independent iff their joint factors:

$$p(x, y) = p(x)p(y)$$



- Two variables are conditionally independent given a third one if for all values of the conditioning variable, the resulting slice factors:

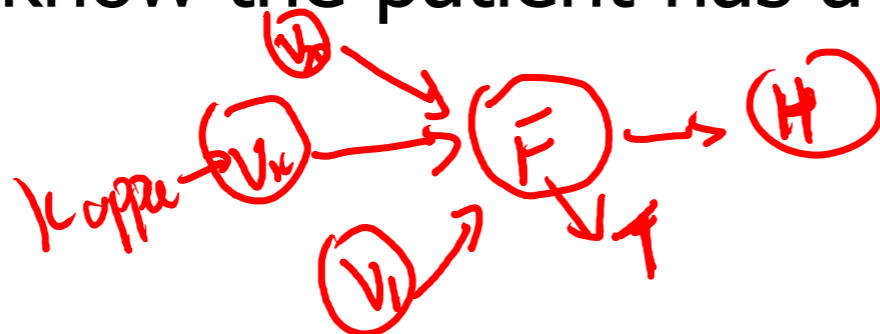
$$p(x, y|z) = p(x|z)p(y|z) \quad \forall z$$

# Poll

- A virus Kappa is known to cause a flu. And a flu is known to cause a headache sometimes. A flu can be caused by multiple reasons. A patient shows up with a diagnosis of a flu. We know the patient has flu. Does the probability of the patient having a headache depend on the virus Kappa now or not?

Options:

- 1) Yes, probability of a headache depends on the virus Kappa. 33%
- 2) No, probability of a headache is independent of the virus Kappa as we know the patient has a flu. 67%





# Conditional Independence

- Examples:

$$P(\text{Virus} \mid \text{Drink Beer}) = P(\text{Virus})$$

iff **Virus** is independent of **Drink Beer**

$$P(\text{Flu} \mid \text{Virus}; \text{Drink Beer}) = P(\text{Flu} \mid \text{Virus})$$

iff **Flu** is independent of **Drink Beer**, given **Virus**

$$P(\text{Headache} \mid \text{Flu}; \text{Virus}; \text{Drink Beer}) =$$

$$P(\text{Headache} \mid \text{Flu}; \text{Drink Beer})$$

iff **Headache** is independent of **Virus**, given **Flu** and **Drink Beer**

Assume the above independence, we obtain:

$$\begin{aligned} & P(\text{Headache}; \text{Flu}; \text{Virus}; \text{Drink Beer}) \\ &= P(\text{Headache} \mid \text{Flu}; \text{Virus}; \text{Drink Beer}) P(\text{Flu} \mid \text{Virus}; \text{Drink Beer}) \\ & \quad P(\text{Virus} \mid \text{Drink Beer}) P(\text{Drink Beer}) \\ &= P(\text{Headache} \mid \text{Flu}; \text{Drink Beer}) P(\text{Flu} \mid \text{Virus}) P(\text{Virus}) P(\text{Drink Beer}) \end{aligned}$$


↑ 3 D

↑ 2 D

↑ 1 D

↑ 1 D

# Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule ← 
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

# Bayes' Rule

- $P(X|Y)$  = Fraction of the worlds in which  $X$  is true given that  $Y$  is also true.
- For example:
  - $H$  = "Having a headache"
  - $F$  = "Coming down with flu"
  - $P(\text{Headache}|\text{Flu})$  = fraction of flu-inflicted worlds in which you have a headache. How to calculate?

- Definition:

$$P(X|Y) = \frac{P(X, Y)}{P(Y)} = \frac{P(Y|X)P(X)}{P(Y)}$$

Corollary:

$$P(X, Y) = P(Y|X)P(X)$$

This is called **Bayes Rule**



# Bayes' Rule

- $$P(\text{Headache}|\text{Flu}) = \frac{P(\text{Headache},\text{Flu})}{P(\text{Flu})}$$
$$= \frac{P(\text{Flu}|\text{Headache})P(\text{Headache})}{P(\text{Flu})}$$

Other cases:

- $$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X|Y)P(Y)+P(X|\neg Y)P(\neg Y)}$$
- $$P(Y = y_i|X) = \frac{P(X|Y)P(Y)}{\sum_{i \in S} P(X|Y = y_i)P(Y=y_i)}$$
- $$P(Y|X, Z) = \frac{P(X|Y, Z)P(Y, Z)}{P(X, Z)} =$$
$$\frac{P(X|Y, Z)P(Y, Z)}{P(X|Y, Z)P(Y, Z)+P(X|\neg Y, Z)P(\neg Y, Z)}$$

# Administrative business

- Office hours are live
- Live Q&A parallel to the lectures to TAs can help with answering questions. Use voting as well so I know what needs to be handled here.
- Chris is holding a python tutorial on Thursday at 6 pm *→ no mpy*
- Project questions answered on Thursday's lecture by me, come one, come all.

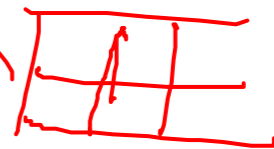
# Outline

- Probability Distributions

$$P(x)$$

- Joint and Conditional Probability Distributions

$$P(y)$$



- Bayes' Rule

$$P(x|y) = \frac{P(x,y)}{P(y)} = \frac{P(y|x)P(x)}{P(y)}$$

- Mean and Variance



- Properties of Gaussian Distribution

- Maximum Likelihood Estimation

# Mean and Variance

- Expectation: The mean value, center of mass, first moment:

$$E_X[g(X)] = \int_{-\infty}^{\infty} g(x) p_X(x) dx = \mu$$

$\int_{-\infty}^{\infty} x p(x) dx$       $h_1, h_2, \dots, h_n$       $\frac{\sum h_i}{n}$

- N-th moment:  $g(x) = x^n$

- N-th central moment:  $g(x) = (x - \mu)^n$

- Mean:  $E_X[X] = \int_{-\infty}^{\infty} x p_X(x) dx$

- $E[\alpha X] = \alpha E[X]$

- $E[\alpha + X] = \alpha + E[X]$

- Variance(Second central moment):  $Var(x) =$

$$E_X[(X - E_X[X])^2] = E_X[X^2] - E_X[X]^2$$

- $Var(\alpha X) = \alpha^2 Var(X)$

- $Var(\alpha + X) = Var(X)$

# For Joint Distributions

- Expectation and Covariance:

- $E[X + Y] = E[X] + E[Y]$

- $cov(X, Y) = E[(X - E_X[X])(Y - E_Y(Y))] = E[XY] - E[X]E[Y]$

- $Var(X + Y) = Var(X) + 2cov(X, Y) + Var(Y)$



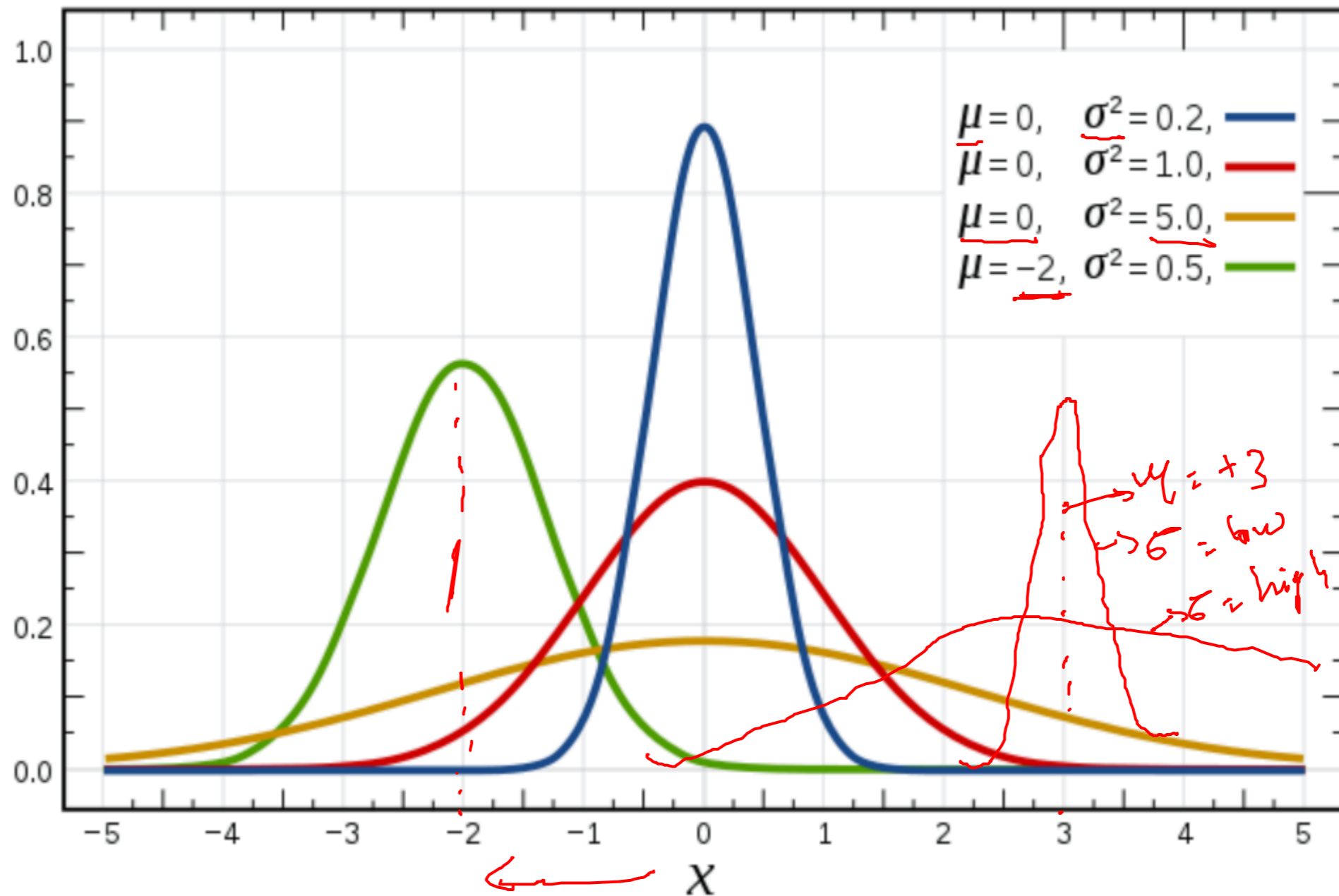
# Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution ←
- Maximum Likelihood Estimation

# Gaussian Distribution

- Gaussian Distribution:  $f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Probability density function //



# Multivariate Gaussian Distribution

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu)\right\}$$

$\begin{matrix} \mu_1 & \mu_1 & \downarrow & \mu_1, \mu_2 \\ \vdots & \vdots & & \\ \mu_n & \mu_n & & \end{matrix}$

- Moment Parameterization  $\mu = E(X)$

$$\Sigma = \text{Cov}(X) = E[(X - \mu)(X - \mu)^\top]$$

- Mahalanobis Distance  $\Delta^2 = (x - \mu)^\top \Sigma^{-1} (x - \mu)$

- Tons of applications (MoG, FA, PPCA, Kalman filter,...)

# Properties of Gaussian Distribution

- The **linear transform** of a Gaussian r.v. is a Gaussian. Remember that no matter how  $x$  is distributed

$$E(\underline{AX + b}) = \underline{AE(X) + b}$$

$$\underline{Cov(AX + b)} = ACov(X)A^T$$

this means that for Gaussian distributed quantities:

$$X \sim N(\mu, \Sigma) \rightarrow AX + b \sim N(A\mu + b, A\Sigma A^T)$$

- The **sum** of two independent Gaussian r.v. is a Gaussian

$$Y = X_1 + X_2, X_1 \perp X_2 \rightarrow \underline{\mu_y} = \mu_1 + \mu_2, \underline{\Sigma_y} = \Sigma_1 + \Sigma_2$$

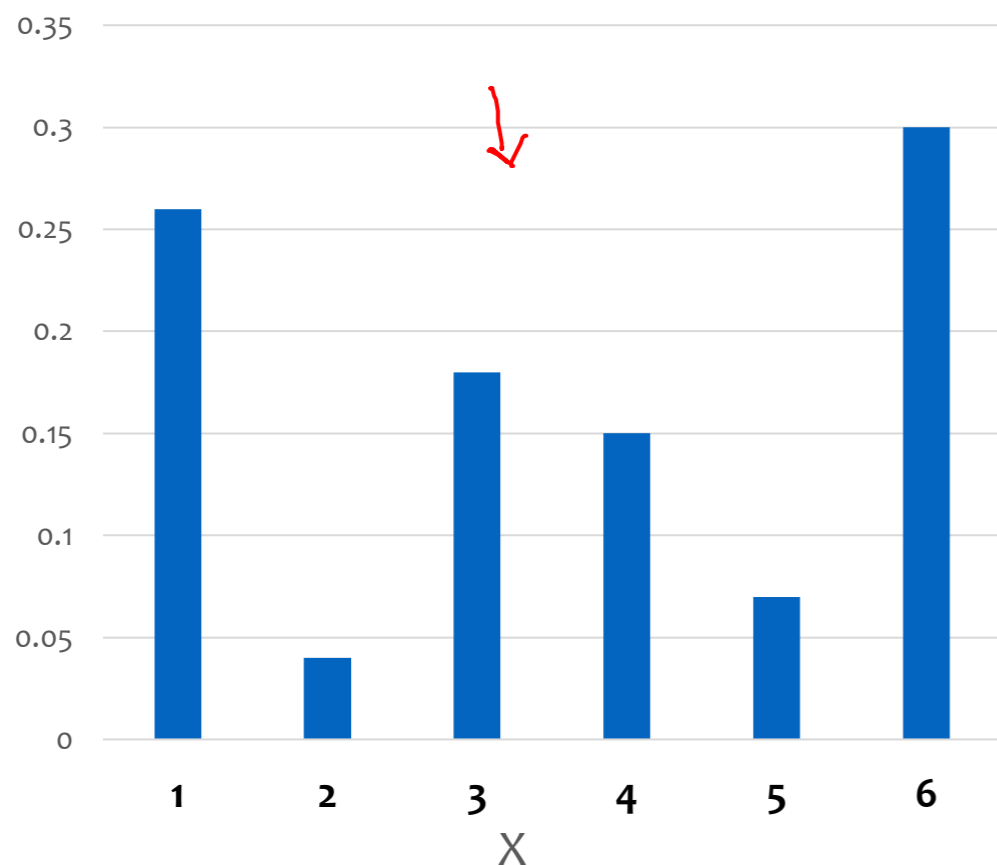
- The **multiplication** of two Gaussian functions is another Gaussian function (although no longer normalized)

$$N(a, A)N(b, B) \propto N(c, C),$$

$$\text{where } C = (A^{-1} + B^{-1})^{-1}, c = CA^{-1}a + CB^{-1}b$$

# Central Limit Theorem

Probability mass function of a **biased** dice



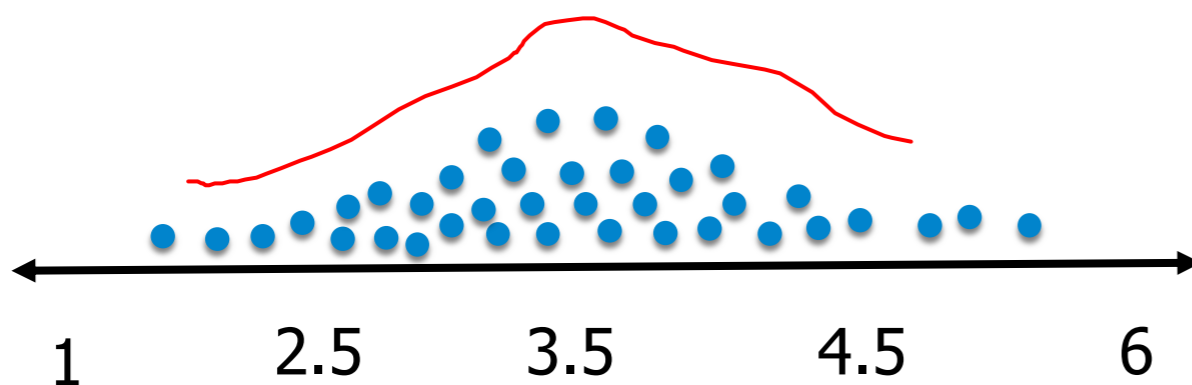
Let's say, I am going to get a sample from this pmf having a size of  **$n = 4$**

→  $S_1 = \{1,1,1,6\} \Rightarrow E(S_1) = 2.25$

→  $S_2 = \{1,1,3,6\} \Rightarrow E(S_2) = 2.75$

⋮

$S_m = \{1,4,6,6\} \Rightarrow E(S_m) = 4.25$



According to CLT, it will follow a bell curve distribution (normal distribution)

# CLT Definition

- Statement: The central limit theorem (due to Laplace) tells us that, subject to certain mild conditions, the sum of a set of random variables, which is of course itself a random variable, has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases (Walker, 1969).

# Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation ←

# Likelihood, what is it? Biased coin from a stranger

→ T

→ T

→ H

→ H

→ H

→ H

→ H

T

T

H

H

H

H

H

Head Biased

Tails biased

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

↑

↑

↑

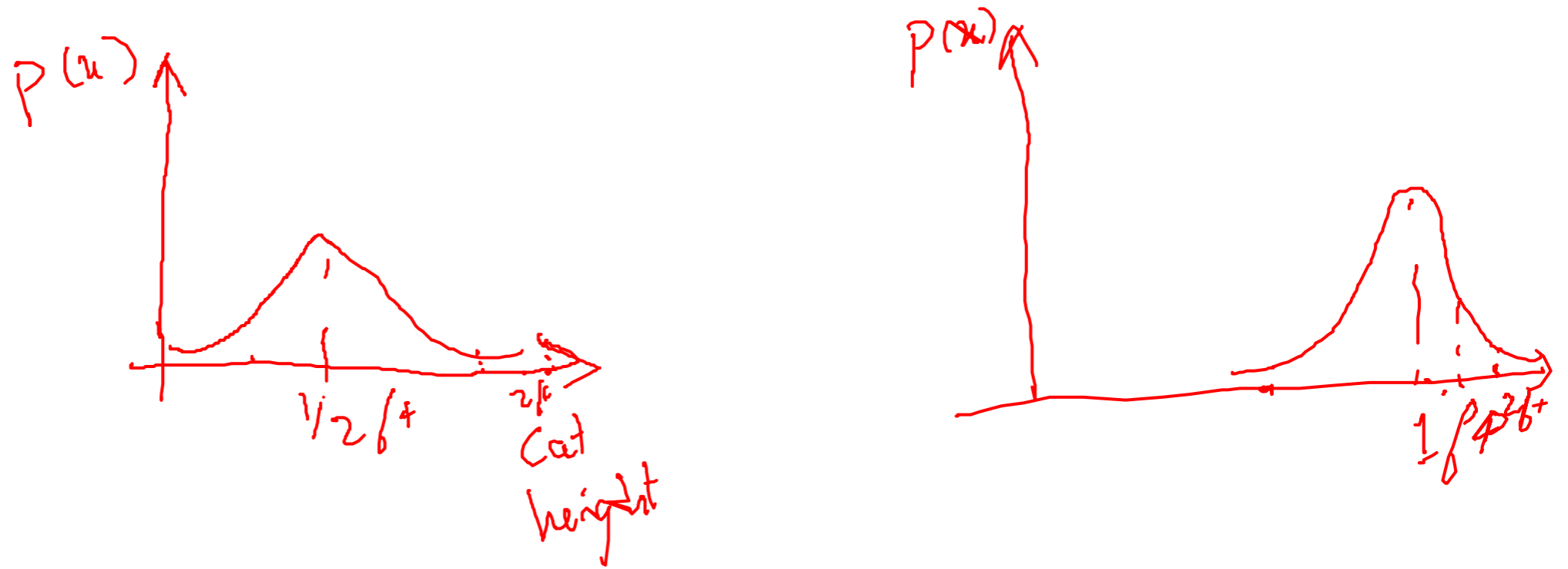
↑

↑



# Likelihood, what is it? or Cat or Dog? Do we know?

## Let's find out!



heard  $\rightarrow [2\sqrt{t}, 1-\sqrt{t}, 0.2\sqrt{t}, 2.2\sqrt{t} \dots]$

$$L(x_{1:n} | \text{cat}) = p(x_1 | \text{cat}) \cdot p(x_2 | \text{cat}) \cdot \dots \cdot p(x_n | \text{cat})$$

$$= \frac{0.001 \times 0.003 \times 0.05 \times 0.005}{C_{\text{cat}}}$$

$$L(x_{1:n} | \text{dog}) = p(x_1 | \text{dog}) \cdot p(x_2 | \text{dog}) \cdot \dots \cdot p(x_n | \text{dog})$$

$$= \frac{C_{\text{dog}}}{C_{\text{cat}}} < C_{\text{dog}}$$

# Maximum Likelihood Estimation

- Probability: inferring probabilistic quantities for data given fixed models (e.g. prob. of events, marginals, conditionals, etc).
- Statistics: inferring a model given fixed data observations (e.g. clustering, classification, regression).

Main assumption:

Independent and identically distributed random variables  
i.i.d

# Maximum Likelihood Estimation

For Bernoulli (i.e. flip a coin):

Objective function:  $f(x_i; p) = p^{x_i} (1 - p)^{1 - x_i}$   $x_i \in \{0, 1\}$  or  $\{head, tail\}$

$$L(p) = Pr(X = x_1, X = x_2, X = x_3, \dots, X = \underline{x_n})$$



We want to know what is the most “likely” value for the probability of success  $p$  given  $n$  observations???

# Maximum Likelihood Estimation

For Bernoulli (i.e. flip a coin):

Objective function:  $f(x_i; p) = p^{x_i}(1-p)^{1-x_i}$   $x_i \in \{0,1\}$  or {head, tail}

$$L(p) = Pr(X = x_1, X = x_2, X = x_3, \dots, X = x_n)$$

i.i.d assumption

$$= Pr(X = x_1)Pr(X = x_2) \dots Pr(X = x_n) = f(p; x_1)f(p; x_2) \dots f(p; x_n)$$

$$L(p) = \prod_{i=1}^n f(x_i; p) = \prod_{i=1}^n p^{x_i}(1-p)^{1-x_i}$$

$$L(p) = p^{x_1}(1-p)^{1-x_1} \times p^{x_2}(1-p)^{1-x_2} \dots \times p^{x_n}(1-p)^{1-x_n} =$$
$$= p^{\sum x_i}(1-p)^{\sum(1-x_i)}$$

We don't like multiplication, let's convert it into summation

What's the trick?

Take the log

$$L(p) = p^{\sum x_i} (1 - p)^{\sum (1 - x_i)}$$

$$\log L(p) = l(p) = \log(p) \sum_{i=1}^n x_i + \log(1 - p) \sum_{i=1}^n (1 - x_i)$$

How to optimize p?

$$\frac{\partial l(p)}{\partial p} = 0 \quad \frac{\sum_{i=1}^n x_i}{p} - \frac{\sum_{i=1}^n (1 - x_i)}{1 - p} = 0$$

$$p = \frac{1}{n} \sum_{i=1}^n x_i$$