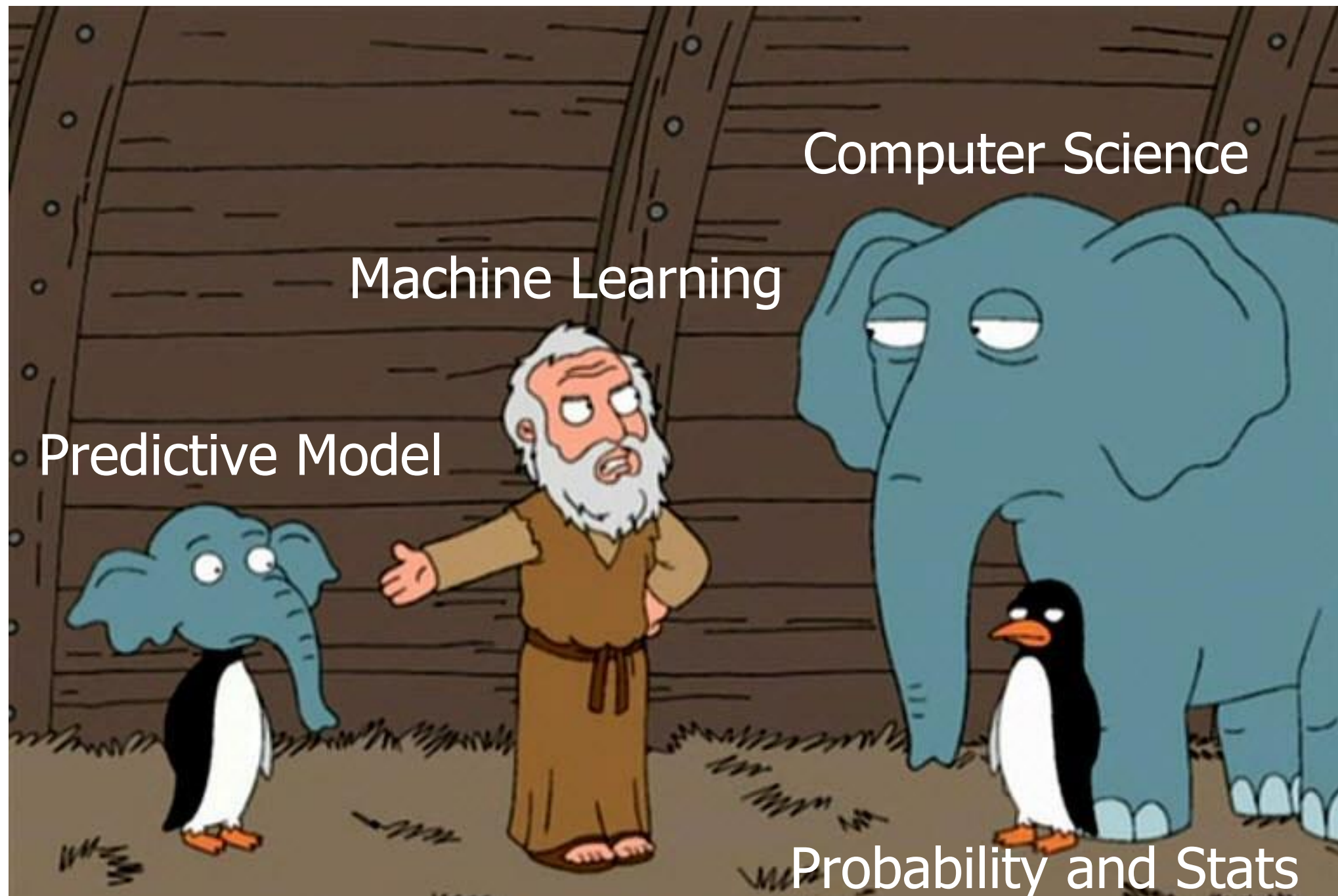




# Probability and Statistics

Mahdi Roozbahani  
Georgia Tech




Computer Science

Machine Learning

Predictive Model

Probability and Stats

# 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  $p(x) = p(X = x)$

$$p(x) \geq 0$$

## Continuous variable

Continuous probability distribution

Probability density function  $\leadsto$  pdf

Density or likelihood value

Temperature (real number)

Gaussian Distribution



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

$$P(X = x_t) = 0$$

Probability distribution function (pdf)

## Discrete variable

Discrete probability distribution

Probability mass function  $\leadsto$  pmf

Probability value

Coin flip (integer)

Bernoulli distribution

pmf  
Probability  
mass function

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

$$\frac{1}{6} + \frac{1}{6} + \dots + \frac{1}{6} = 1$$

# Continuous Probability Functions

Parameters  $\leadsto \Theta$

- Examples:

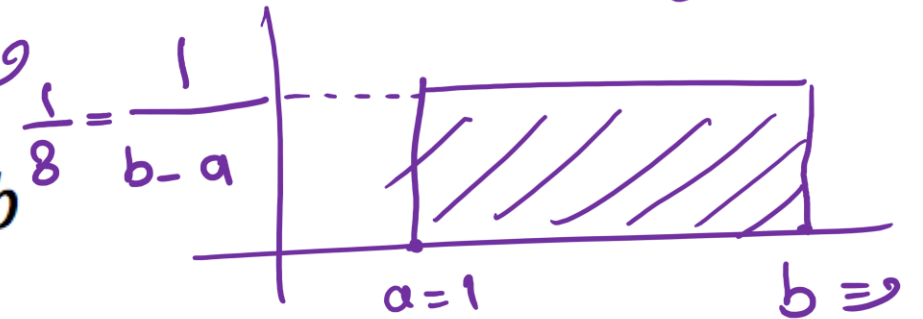
$$\{a, b\} \in \Theta$$

$$\frac{1}{8} \times \overbrace{(b-a)}^8 = 1$$

- Uniform Density Function:

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

$$a=1 \quad b=9$$



- Exponential Density Function:

$$\{\mu\} \in \Theta$$

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

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

$\mu$   $\rightarrow$  Arithmetic average  
 $\mu$   $\rightarrow$  Expected value or weighted avg

- Gaussian(Normal) Density Function

$$\{\mu, \sigma\} \in \Theta$$

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

$$\mu = \frac{1}{N} \sum x_i$$

$$\sigma^2 = \frac{1}{N} \sum (x_i - \mu)^2$$



# Discrete Probability Functions

- Examples:

- Bernoulli Distribution:

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

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

- Binomial Distribution:


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

**k** is number of successes

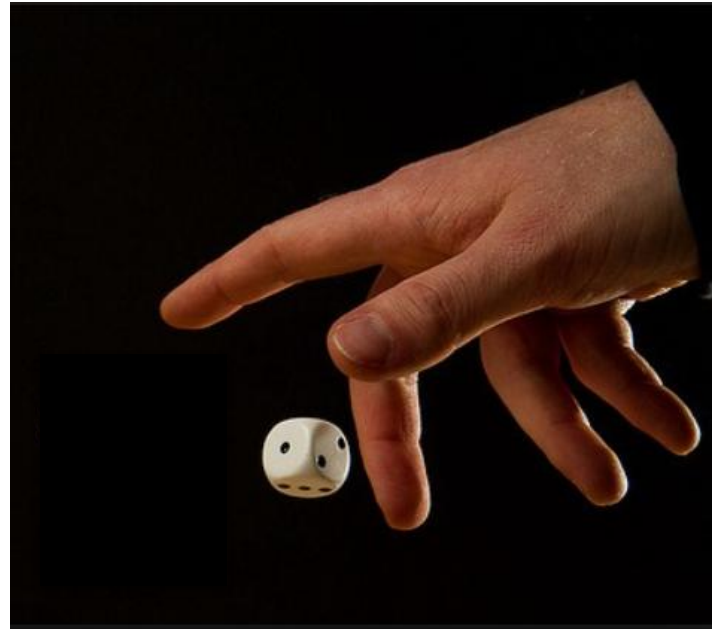
**n-k** is number of failures

$\binom{n}{k}$  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

$\mathbf{X}$  and  $\mathbf{Y}$  are random variables

$\mathbf{N}$  = total number of trials

$n_{ij}$  = Number of occurrence

Joint probability table

		$\mathbf{X}$						
		$x_{i=1} = 1$	$x_{i=2} = 2$	$x_{i=3} = 3$	$x_{i=4} = 4$	$x_{i=5} = 5$	$x_{i=6} = 6$	$C_j$
$\mathbf{Y}$	$y_{j=2} = tail$	$n_{ij} = 3$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 5$	$n_{ij} = 1$	$n_{ij} = 5$	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	$N=35$

**X**

$$x_{i=1} = 1 \quad x_{i=2} = 2 \quad x_{i=3} = 3 \quad x_{i=4} = 4 \quad x_{i=5} = 5 \quad x_{i=6} = 6$$

**(C<sub>j</sub>)**

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

$$P(Y = \text{tail}, X = 4) = \frac{5}{35} = \frac{n_{ij}}{N}$$

$$P(Y = \text{tail}) = \frac{20}{35} = \frac{C_j}{N} \quad P(X = 4) = \frac{7}{35} = \frac{C_i}{N}$$

$$P(Y = \text{tail} | X = 4) = \frac{5}{7} = \frac{n_{ij}}{C_i} \quad P(X = 4 | Y = \text{tail}) = \frac{5}{20} = \frac{n_{ij}}{C_j}$$

$$P(Y, X) = \frac{n_{ij}}{N} = \frac{n_{ij}}{C_j} \frac{C_j}{N} = P(X | Y) P(Y) \rightarrow \text{product rule}$$

$$= \frac{n_{ij}}{C_i} \frac{C_i}{N} = P(Y | X) P(X)$$

$$P(a, b, c) = P(a | \underline{b, c}) P(\underline{b, c})$$

**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**

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

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

# Conditional Independence

$$P(H, F, V, D) = P(H|F, V, D) \cdot P(F, V, D)$$

- Examples:

$$P(\text{Virus} | \text{DrinkBeer}) = P(\text{Virus})$$

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

$$P(\text{Flu} | \text{Virus}, \text{DrinkBeer}) = P(\text{Flu} | \text{Virus})$$

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

$$P(\text{Headache} | \text{Flu}, \text{Virus}, \text{DrinkBeer}) \\ = P(\text{Headache} | \text{Flu}, \text{DrinkBeer})$$

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

Assume the above independence, we obtain:

$$\begin{aligned} &P(\text{Headache}, \text{Flu}, \text{Virus}, \text{DrinkBeer}) \\ &= P(\text{Headache} | \text{Flu}, \text{Virus}, \text{DrinkBeer}) P(\text{Flu} | \text{Virus}, \text{DrinkBeer}) \\ &\quad P(\text{Virus} | \text{DrinkBeer}) P(\text{DrinkBeer}) \\ &= P(\text{Headache} | \text{Flu}, \text{DrinkBeer}) P(\text{Flu} | \text{Virus}) P(\text{Virus}) P(\text{DrinkBeer}) \end{aligned}$$


$$P(V, D) = \frac{P(V|D)P(D)}{P(V, D)}$$

$$P(H|F, D) \cdot P(F|V, D) \cdot P(V, D)$$

$$P(H|F, D) \cdot P(F|V) \cdot P(V|D) \cdot P(D)$$

$$P(H|F, D) \cdot P(F|V) \cdot P(V) \cdot P(D)$$

# 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.

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

- 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?

$$P(Y) = \sum_{x_i} P(Y, X=x_i) = \sum P(Y|X=x_i) P(X=x_i)$$

- 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

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

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

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

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

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

Other cases:

- $$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X|Y)P(Y) + P(X|\neg Y)P(\neg Y)}$$


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

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

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

# Outline

- ① product
- ② sum rule

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

$$P(Y) = \sum_x P(Y, X=x)$$

$$\begin{aligned} P(Y, X) &= P(Y|X) P(X) \\ &= P(\underline{X}|Y) P(Y) \end{aligned}$$

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

$E[\cdot]$

# Mean and Variance

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

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

$$E[g(x)] = \sum g(x)p(x)$$

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

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

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

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

$$\text{Var}(x) = E[(x - E[x])^2]$$

- Variance(Second central moment):  $\text{Var}(x) =$

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

$$\text{Var}(x) = E[x^2] - E[x]^2$$

- $\text{Var}(\alpha X) = \alpha^2 \text{Var}(X)$

- $\text{Var}(\alpha + X) = \text{Var}(X)$

## Mean and average $g(x) = x$

$$X = [1, 2, 3]$$

$$P(x) = \left[ \frac{1}{6}, \frac{2}{6}, \frac{3}{6} \right]$$

$$E[g(x)] = \sum_{x_i} g(x=i) P(x=i) =$$

$$g(x=1)P(x=1) + g(x=2)P(x=2) + g(x=3)P(x=3)$$

$$1 * \frac{1}{6} + 2 * \frac{2}{6} + 3 * \frac{3}{6}$$

$$= \frac{14}{6}$$

$$\text{avg}(x) \neq E[g(x)]$$

$$\text{avg}(x) = \frac{1+2+3}{3} = 2$$

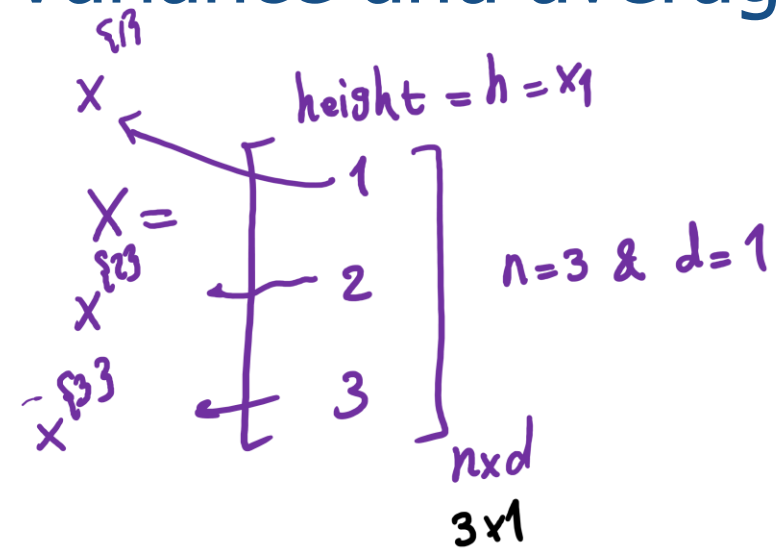
---

$$X = [1, 2, 2, 3, 3, 3] \Rightarrow$$

$$\text{avg}(x) = \frac{1+2+2+3+3+3}{6} = \frac{14}{6}$$

$$\text{avg}(x) = E[x]$$

# Variance and average:



$$\mu_h = \text{avg}_h = \frac{1+2+3}{3} = 2$$

$$\sigma_h^2 = \frac{1}{N} \sum_{i=1}^N (x^{\{i\}} - \mu_h)^2$$

For loop

$$\sigma_h^2 = E[(x^{\{i\}} - \mu_h)^2] = E[(x^{\{i\}} - E[x_h])^2]$$

$$\mu = E[\cdot]$$

$$\frac{1}{N} \sum_{i=1}^N (\cdot) = \text{avg} = E[\cdot]$$

Matrix operation

$$X \rightarrow \bar{X}$$

$$\begin{bmatrix} h \\ 1 \\ 2 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} h - \mu_h = \bar{h} \\ 1 - \mu_h = -1 \\ 2 - \mu_h = 0 \\ 3 - \mu_h = 1 \end{bmatrix}$$

$$\sigma_h^2 = \frac{1}{N} \underbrace{\bar{X}^T}_{1 \times 3} \underbrace{\bar{X}}_{3 \times 1} = \frac{1}{N} [1 - \mu_h \quad 2 - \mu_h \quad 3 - \mu_h] \begin{bmatrix} 1 - \mu_h \\ 2 - \mu_h \\ 3 - \mu_h \end{bmatrix} = \frac{1}{N} [(1 - \mu_h)^2 + (2 - \mu_h)^2 + (3 - \mu_h)^2]$$

$$= \frac{1}{N} \sum_{i=1}^N (x^{\{i\}} - \mu_h)^2$$

# Covariance:

$$X = \begin{matrix} & \begin{matrix} h & \text{weight} = w \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} \end{matrix}$$

$3 \times 2 = n \times d$   
 $2 \times 2 = d \times d$

$$X \rightarrow \bar{X} \Rightarrow \bar{X} = \begin{matrix} \begin{matrix} \bar{h} & \bar{w} \end{matrix} \\ \begin{bmatrix} 1 - \mu_h & 4 - \mu_w \\ 2 - \mu_h & 5 - \mu_w \\ 3 - \mu_h & 6 - \mu_w \end{bmatrix} \end{matrix}$$

$$COV = \frac{1}{N} \begin{matrix} \bar{X}^T \\ \begin{matrix} 2 \times 3 & 3 \times 2 \end{matrix} \end{matrix} \bar{X} = \frac{1}{N} \begin{bmatrix} 1 - \mu_h & 2 - \mu_h & 3 - \mu_h \\ 4 - \mu_w & 5 - \mu_w & 6 - \mu_w \end{bmatrix} \begin{bmatrix} 1 - \mu_h & 4 - \mu_w \\ 2 - \mu_h & 5 - \mu_w \\ 3 - \mu_h & 6 - \mu_w \end{bmatrix} = \begin{matrix} h & w \\ \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \end{matrix}$$

①  $\frac{1}{N} [(1 - \mu_h)^2 + (2 - \mu_h)^2 + (3 - \mu_h)^2] = \sigma_h^2$

②  $\frac{1}{N} [(1 - \mu_h)(4 - \mu_w) + (2 - \mu_h)(5 - \mu_w) + (3 - \mu_h)(6 - \mu_w)] = \sigma_{hw}$

③  $= \textcircled{2} = \sigma_{wh}$

④  $\sigma_w^2$

$$COV = \begin{matrix} & \begin{matrix} h & w \end{matrix} \\ \begin{matrix} h & w \end{matrix} & \begin{bmatrix} \sigma_h^2 & \sigma_{hw} \\ \sigma_{wh} & \sigma_w^2 \end{bmatrix} \end{matrix}$$

Symmetrical

Correlation: EDA

$$X \rightarrow \bar{X} \rightarrow \bar{X}^*$$

$$\bar{X} = \begin{bmatrix} \bar{h} & \bar{w} \\ 1-\mu_h & 4-\mu_w \\ 2-\mu_h & 5-\mu_w \\ 3-\mu_h & 6-\mu_w \end{bmatrix}$$

$$\bar{X}^*$$

$$= \begin{bmatrix} \frac{\bar{h}^*}{\sigma_h} & \frac{\bar{w}^*}{\sigma_w} \\ \frac{1-\mu_h}{\sigma_h} & \frac{4-\mu_w}{\sigma_w} \\ \frac{2-\mu_h}{\sigma_h} & \frac{5-\mu_w}{\sigma_w} \\ \frac{3-\mu_h}{\sigma_h} & \frac{6-\mu_w}{\sigma_w} \end{bmatrix}$$

$$\bar{X}^{*T}$$

$$Cor = \frac{1}{N} \bar{X}^{*T} \bar{X}^* = \begin{bmatrix} \textcircled{1} & \textcircled{2} \\ \textcircled{3} & \textcircled{4} \end{bmatrix}_{d \times d}$$

Standardization

$$\textcircled{1} \frac{1}{N} \left( \frac{1-\mu_h}{\sigma_h} \right)^2 + \left( \frac{2-\mu_h}{\sigma_h} \right)^2 + \left( \frac{3-\mu_h}{\sigma_h} \right)^2 = \frac{(1-\mu_h)^2 + (2-\mu_h)^2 + (3-\mu_h)^2}{N \sigma_h^2} = \frac{\sigma_h^2}{\sigma_h^2} = 1$$

$$Cor = \begin{matrix} & \begin{matrix} h & w \end{matrix} \\ \begin{matrix} h \\ w \end{matrix} & \begin{bmatrix} 1 & -1 \leq \sigma_{hw} \leq 1 \\ -1 \leq \sigma_{wh} \leq 1 & 1 \end{bmatrix} \end{matrix}$$



# For Joint Distributions

$$\text{cov}(X, Y) = E[XY] - E[X]E[Y] = \underbrace{E[z^3]}_0 - \underbrace{E[z]}_0 \underbrace{E[z^2]}_{=1} = 0$$

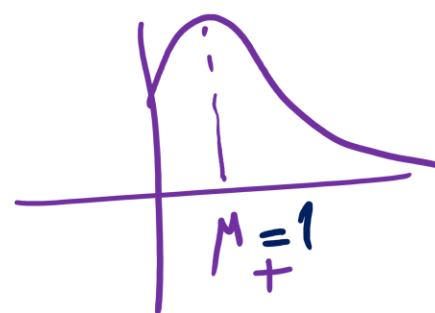
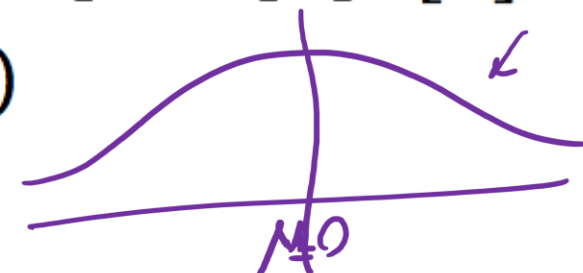
## Expectation and Covariance:

- $E[X + Y] = E[X] + E[Y]$       $E[a+b+c] = E[a] + E[b+c]$
- $\text{cov}(X, Y) = E[(X - E_X[X])(Y - E_Y[Y])] = E[XY] - E[X]E[Y]$
- $\text{Var}(X + Y) = \text{Var}(X) + 2\text{cov}(X, Y) + \text{Var}(Y)$

$X = z \rightsquigarrow$  Standard Gaussian distribution

$Y = z^2$  chi-squared

$Y = z^3$




$$\text{Var}(X) = 1 = E[X^2] - \underbrace{(E[X])^2}_{=0}$$

$$1 = E[z^2] - 0^2 \Rightarrow E[z^2] = 1$$



# Outline

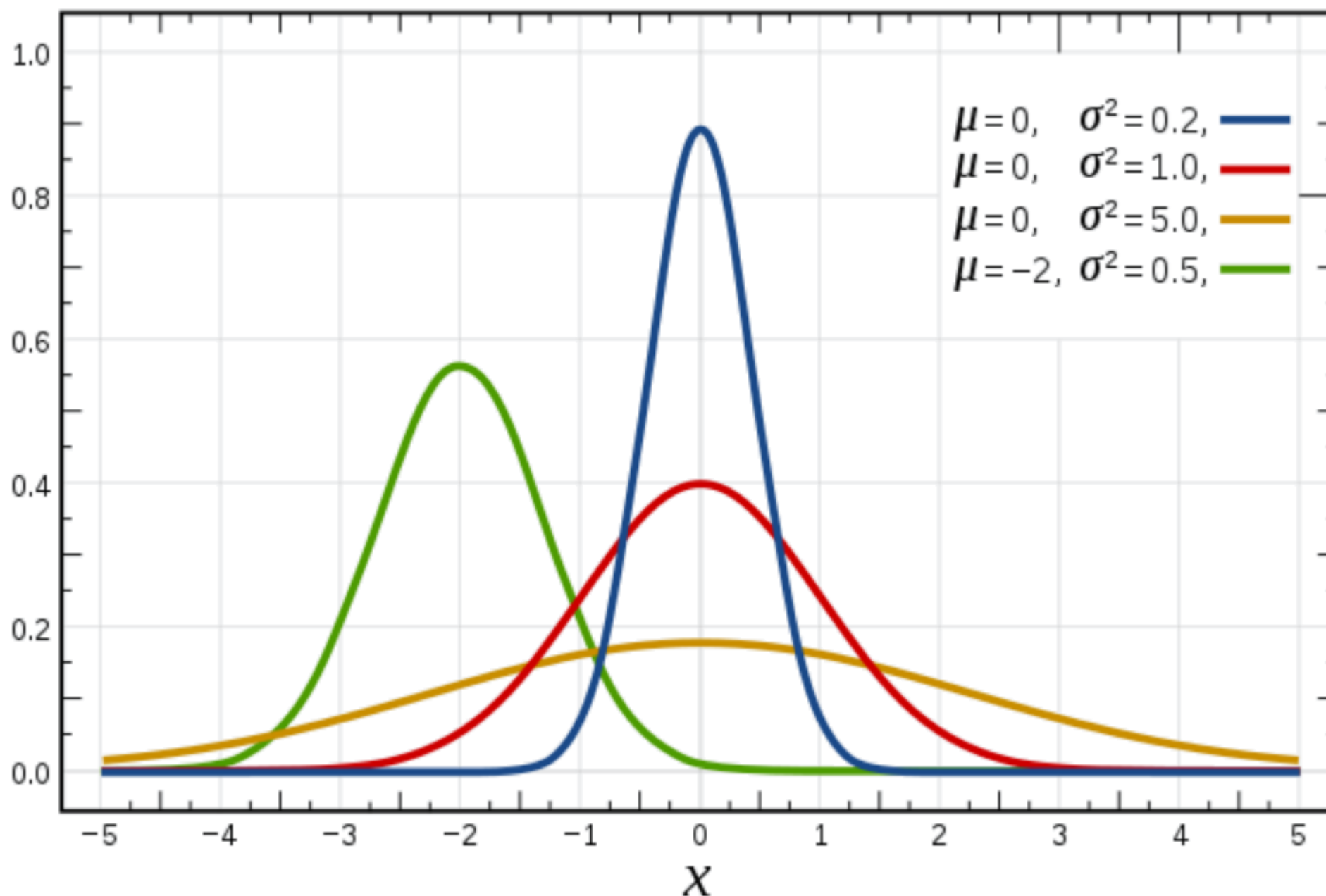
- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution ← 
- Maximum Likelihood Estimation

# Gaussian Distribution

$$\{\mu, \sigma^2\} \in \theta$$

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

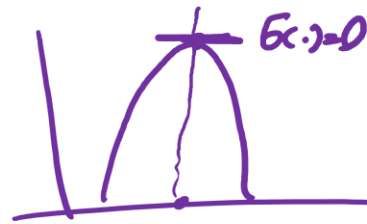
Probability density function



Probability versus likelihood

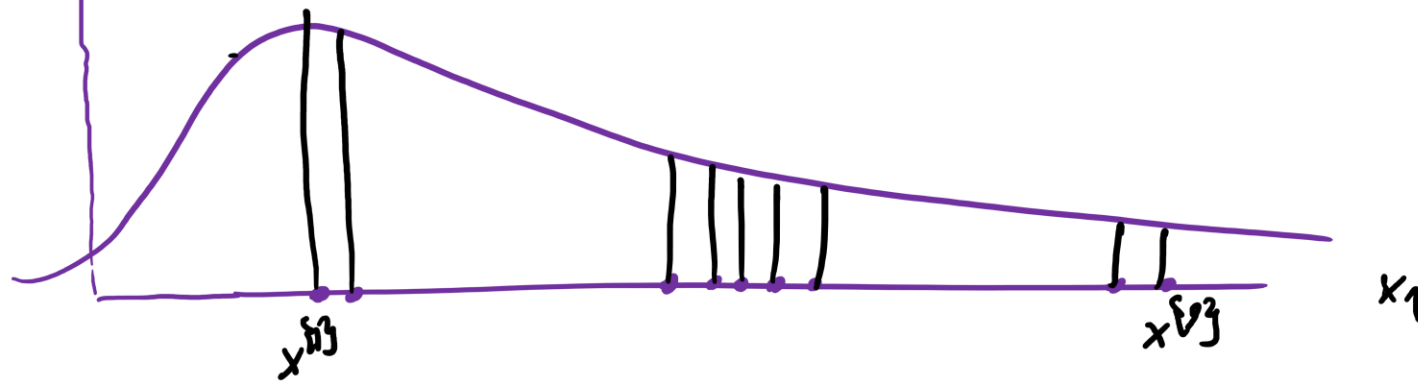
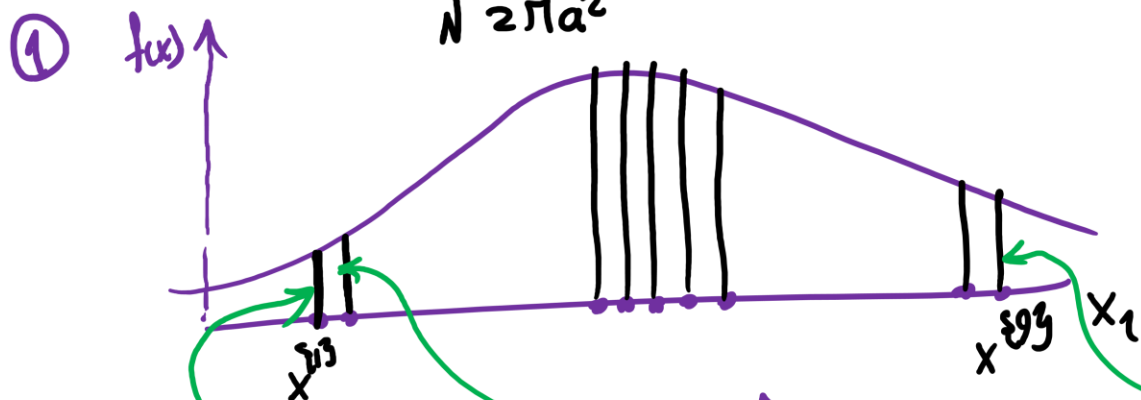
Prob vs Likelihood ①  
 coin  $\rightarrow P(X=T) = \frac{1}{2}$

②  $\Rightarrow [T, T, T, H] \Rightarrow P(X=T) = \frac{3}{4}$



$$f(x|a,b) = \frac{1}{\sqrt{2\pi a^2}} \exp\left(-\frac{(x-b)^2}{2a^2}\right)$$

$f(x)$  density or likelihood



Objective function

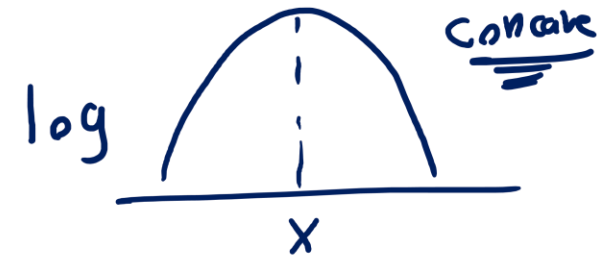
$\text{Max } L(a,b|x) \rightarrow a=1, b=1$   
 $\rightarrow a=1, b=-1$

$$f(x|a,b) = f(x^{(1)}, x^{(2)}, \dots, x^{(n)}|a,b)$$

$$f(x^{(1)}, x^{(2)}, \dots, x^{(n)}|a,b) = f(x^{(1)}|a,b) f(x^{(2)}|a,b) \dots f(x^{(n)}|a,b)$$

# Prob vs Likelihood

$$\log(a \cdot b) = \log a + \log b$$



$$\text{Max } L(a, b | x) \rightsquigarrow \text{Max } \log L(a, b | x)$$

$$f(x^{(1)} | a, b) \dots f(x^{(N)} | a, b) \rightarrow \log(f(x^{(1)} | a, b)) + \dots + \log(f(x^{(N)} | a, b))$$

$$\prod_{i=1}^N f(x^{(i)} | a, b)$$

$$\sum_{i=1}^N \log f(x^{(i)} | a, b)$$

$$\text{Max } \ell(a, b | x) = \sum_{i=1}^N \log f(x^{(i)} | a, b)$$

$$\frac{\partial \ell(a, b | x)}{\partial a} = 0 \rightsquigarrow a = \sigma^2 \quad \frac{\partial \ell(a, b | x)}{\partial b} = 0 \Rightarrow b = \mu$$

# Multivariate Gaussian Distribution

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

- 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

$$E(AX + b) = AE(X) + b$$

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

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

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

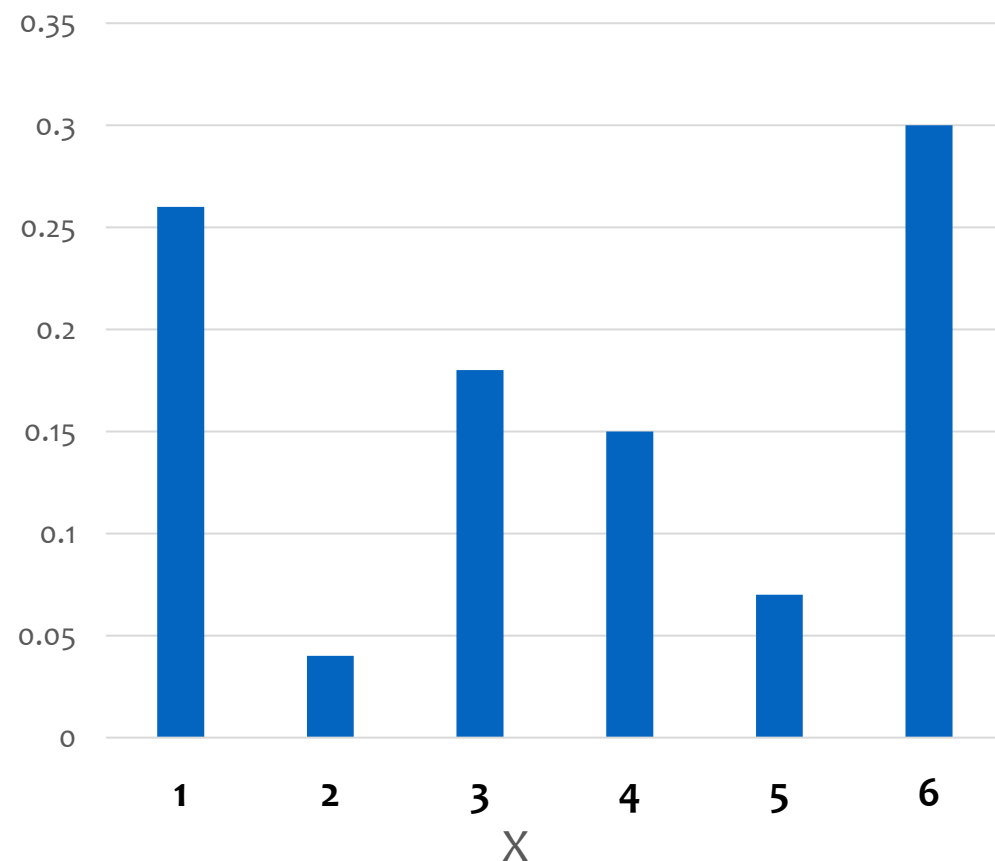
- The **multiplication** of two Gaussian functions is another Gaussian function (although no longer normalized)

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

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

# Central Limit Theorem

Probability mass function of a **biased** dice



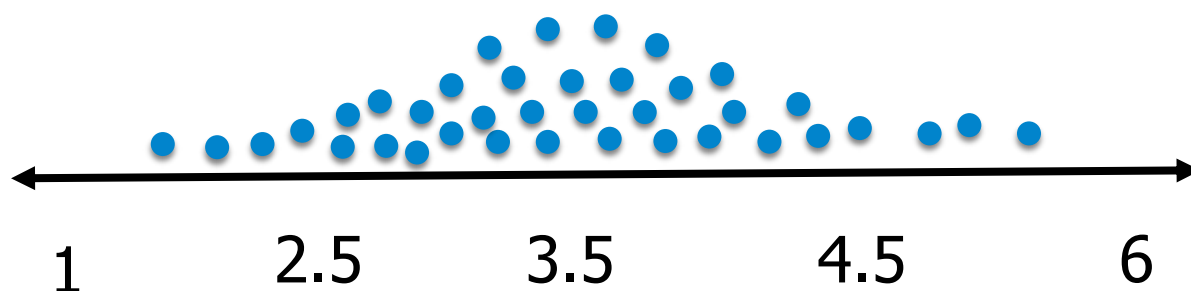
Let's say, I am going to get a sample from this pmf having a size of  **$n = 4$**

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

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


$\vdots$

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



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

# Outline

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



# Maximum Likelihood Estimation

$M$   $L$   $E$  ←

- 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:  $L(\theta|X) = \prod_{i=1}^N p(x_i|\theta)$   $\ell(\theta|x) = \sum_{i=1}^N \log p(x_i|\theta)$

$P(x_i|\theta) = \theta^{x_i}(1-\theta)^{1-x_i}$   $x_i \in \{0,1\}$  or  $\{head, tail\}$

$$L(\theta|X) = L(\theta|X = x_1, X = x_2, X = x_3, \dots, X = x_n)$$

i.i.d assumption

$$L(\theta|X) = \prod_{i=1}^n P(x_i|\theta)$$

Objective function

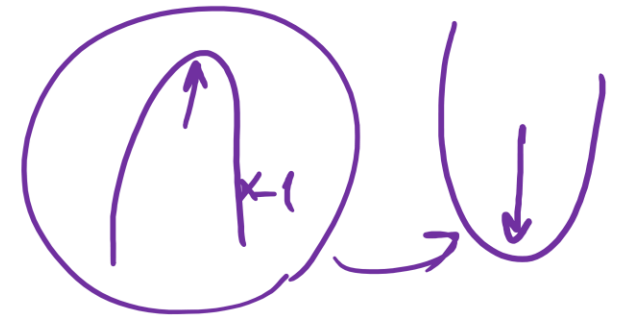
$$L(\theta|X) = \prod_{i=1}^n P(x_i|\theta) = \prod_{i=1}^n \theta^{x_i}(1-\theta)^{1-x_i}$$

$$\begin{aligned} L(\theta|X) &= \theta^{x_1}(1-\theta)^{1-x_1} \times \theta^{x_2}(1-\theta)^{1-x_2} \dots \times \theta^{x_n}(1-\theta)^{1-x_n} = \\ &= \theta^{\sum x_i} (1-\theta)^{\sum (1-x_i)} \end{aligned}$$

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

What's the trick?

Take the log



$$L(\theta|X) = \theta^{\sum x_i} (1 - \theta)^{\sum (1-x_i)}$$

Max

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

How to optimize  $\theta$ ?

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

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

$x=1$        $70 \rightarrow H$   
 $x=0$        $30 \rightarrow T$

$$\Rightarrow \Theta = \frac{1+1+1+\dots+0+0}{100} = 0.7$$