

Probability and Statistics

Nimisha Roy

Lecturer, SCI, College of Computing, Georgia Tech

Director, Online Undergraduate Initiatives

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

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

random variable

specific outcome

$$p(x) \geq 0$$

Continuous variable

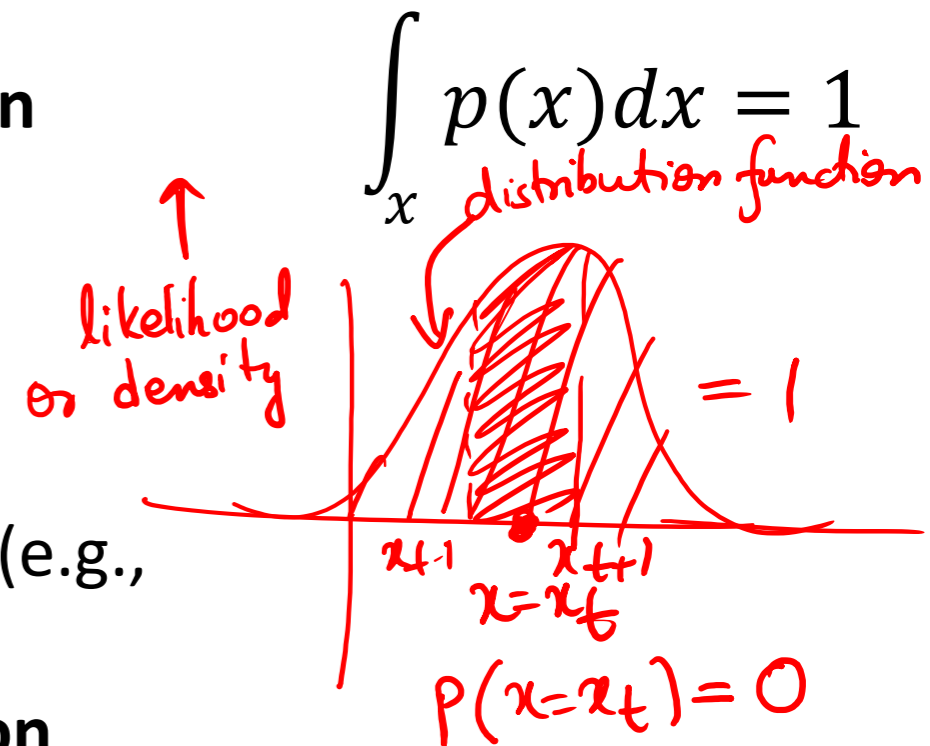
Definition: Takes values from a continuous range (e.g., any real number within an interval).

Distribution: Governed by a **probability density function** (PDF).

Key Concepts:

- The **density** represents likelihood, but not actual probability at a specific point.
- Example: **Temperature**, which can be any real value (e.g., 72.3°F).
- Common distribution: **Gaussian (Normal) Distribution**.

probability distribution function



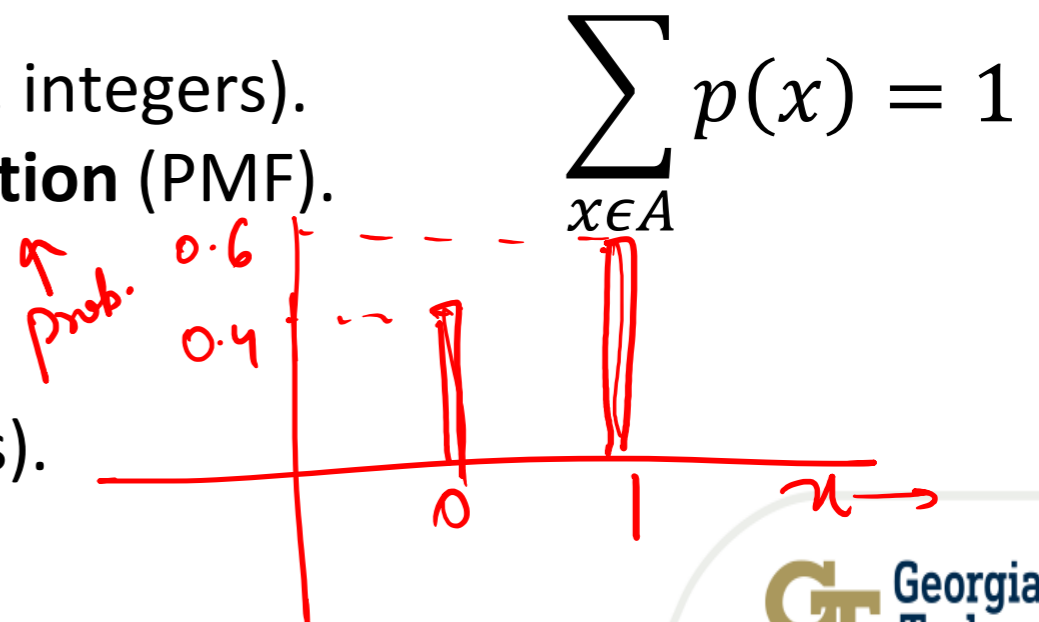
Discrete variable

Definition: Takes values from a countable set (e.g., integers).

Distribution: Governed by a **probability mass function** (PMF).

Key Concepts:

- The function directly gives **probability values**.
- Example: **A coin flip** (e.g., 0 for tails, 1 for heads).
- Common distribution: **Bernoulli Distribution**.



$f(x) \rightarrow$ objective function

$\theta \rightarrow$ parameters wt.

Continuous Probability Functions

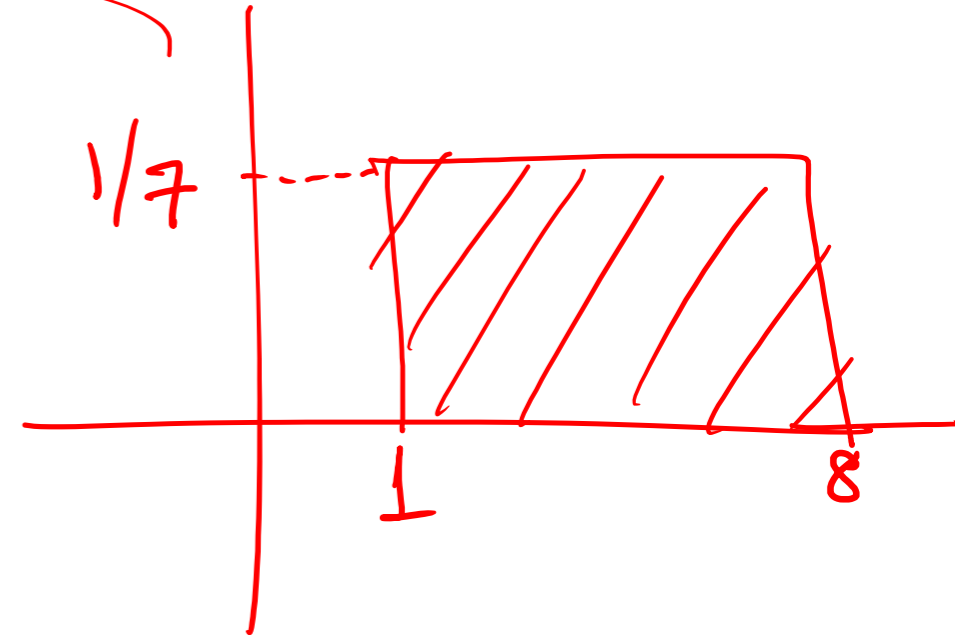
$$\frac{1}{7}(8-1) = 1$$

Examples:

Uniform Density Function:

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

$\theta \in (a, b)$
min \leftarrow a , b \leftarrow max



Exponential Density Function:

$$f_x(x) = \frac{1}{\mu} e^{-\frac{x}{\mu}}$$

$\theta \in a$

for $x \geq 0$

$$f(x) = \frac{1}{a} e^{-x/a}$$

$a = \mu$

Gaussian(Normal) Density Function

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

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**.

E.g. - Given a biased coin with $\mu = 0.3$ for tail ($x = 1$), what is the probability of getting 4 heads given we flip the coin 10 times?

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

n_{ij} = Number of occurrence

		X						C_j
		$x_{i=1} = 1$	$x_{i=2} = 2$	$x_{i=3} = 3$	$x_{i=4} = 4$	$x_{i=5} = 5$	$x_{i=6} = 6$	
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
 $x_{i=1} = 1$ $x_{i=2} = 2$ $x_{i=3} = 3$ $x_{i=4} = 4$ $x_{i=5} = 5$ $x_{i=6} = 6$ C_j
Y $y_{j=2} = \text{tail}$ $y_{j=1} = \text{head}$ C_i

$n_{ij} = 3$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 5$	$n_{ij} = 1$	$n_{ij} = 5$	20
$n_{ij} = 2$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 1$	15
5	6	6	7	5	6	<u>N=35</u>

Joint probabilities

$$P(x=2, y=t) = 4/35 = n_{ij}/N$$

Marginal Prob.

$$P(x=2) = 6/35 = C_i/N$$

$$P(y=t) = 20/35 = C_j/N$$

Sum rule \Rightarrow

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

$$P(y) = \sum_x P(x,y)$$

$$P(y) = \sum_{x,z,w} P(x,y,z,w)$$

Cond. Prob.

$$P(x=2 | y=t) = 4/20 = \frac{n_{ij}}{C_j}$$

$$P(y=t | x=2) = 4/6 = \frac{n_{ij}}{C_i}$$

Shrunk sample space

Product rule =

$$P(x,y) = P(x|y) \cdot P(y)$$

$$= P(y|x) \cdot P(x)$$

$$P(x,y,z) = P(x|y,z) P(y,z)$$

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

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

If $P(x|y) = P(x) \Rightarrow x$ is conditionally ind. of y

Conditional Independence

- 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**

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

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

$$\begin{aligned} & P(h, f, d, v) \\ &= P(h | f, v, d) P(f, v | d) \\ &= P(h | f, d) P(f | v, d) P(v | d) \\ &= P(h | f, d) P(f | v) P(v | d) \\ &= P(h | f, d) P(f | v) P(v | d) P(d) \\ &\leftarrow P(h | f, d) P(f | v) P(v) P(d) \end{aligned}$$

Assume the above independence, we obtain:

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

Recap

• Sample Space and Event

all possible outcomes

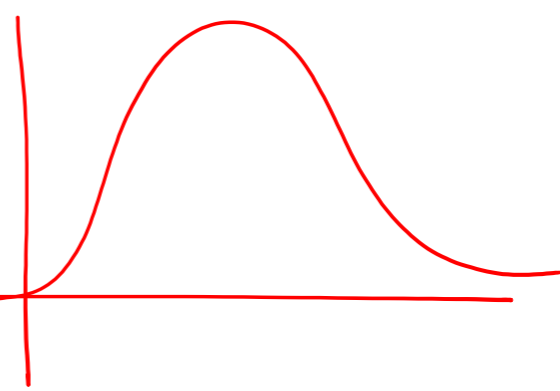
1 possible outcome

• Probability Distribution Function

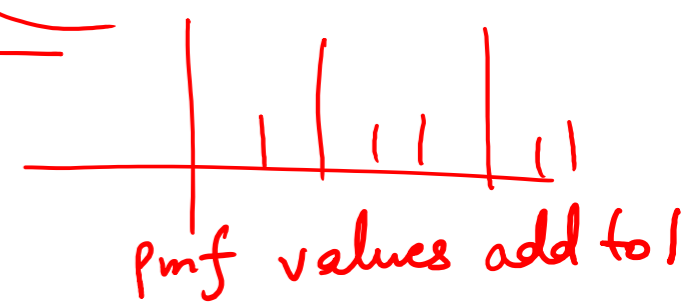
continuous
prob. density function

discrete

prob. mass function



area under the curve is 1
OR



• Objective Function

$f(x)$ ← optimize it → θ parameters.

• Joint Probability, Marginal Probability & Conditional Probability

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

$p(x)$

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

• 2 rules

SUM RULE →

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


$$p(x) = \sum_{y, z, w} p(x, y, z, w)$$

PRODUCT RULE
 $p(x|y) = p(x)$
 x & y are conditionally independence

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

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

Outline

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

Bayes' Rule

Mix of sum & product rule

$$P(x, y) = P(x|y)P(y) \\ = P(y|x)P(x)$$

- $P(X|Y)$ = 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?

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

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

- 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

$$P(x, z) = \sum_y P(x, z, y)$$

$$= \sum_y P(x|z, y) P(z, y) = P(x|z, y) P(z, y) + P(x|z, \bar{y}) P(z, \bar{y})$$

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

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

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

binary ←

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

Other cases:

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

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

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

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

binary ←

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

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


$$P(y|x, z) P(x, z) = P(y, x, z)$$

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

$$P(x, z) = \sum_y P(x, z, y)$$

binary ←

Outline

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

$$E[g(x)] = \int_{-\infty}^{\infty} g(x) p(x) dx$$

$$= \sum_{x_i} g(x_i) p(x_i)$$

Mean and Variance

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

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

In ML,
Expectation =
Arithmetic Mean

- N-th moment: $g(x) = x^n$
- N-th central moment: $g(x) = (x - \mu)^n$
- Mean: $E_X[X] = \int_{-\infty}^{\infty} x p_X(x) dx$

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

$$\text{Var} = \frac{1}{n} \sum_i (x_i - \mu)^2$$

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

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

$$E[X - E[X]]^2$$

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

Expected value and average

Why arithmetic mean is same as expected value:

3 sided biased die: $X = [1, 2, 3]$

$p(x) = [1/6, 1/3, 1/2]$
valid pmf??

$$g(x) = x$$

$$E[g(x)] = \sum_i g(x_i) \cdot p(x_i)$$

$$= 1 \times \frac{1}{6} + 2 \times \frac{1}{3} + 3 \times \frac{1}{2} = \frac{1+4+9}{6} = \frac{14}{6}$$

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

Roll die six times $\rightarrow [1, 2, 2, 3, 3, 3]$

$$E[g(x)] = 14/6$$
$$\mu(x) = \frac{1+2+2+3+3+3}{6} = 14/6$$

Expectation determined by probability distribution. Arithmetic average determined by observed outcomes of trials.

Both are the same in our class because we assume we have sufficient data that models the distribution. So, the arithmetic average is a good estimate of the true expectation (Law of Large Numbers).

Variance and average:

$$X = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 3 \end{bmatrix} \begin{matrix} n \\ \times \\ d \end{matrix}$$

$$\text{Var} = \frac{1}{n} X^T X$$

$$\bar{x} = \begin{bmatrix} 1 - \mu_h \\ 2 - \mu_h \\ 3 - \mu_h \end{bmatrix} 3 \times 1$$

$$\frac{1}{n} X^T X = \frac{1}{2} \begin{bmatrix} 1 - \mu_h & 2 - \mu_h & 3 - \mu_h \end{bmatrix} \begin{bmatrix} 1 - \mu_h \\ 2 - \mu_h \\ 3 - \mu_h \end{bmatrix}$$
$$\text{Var} = \frac{1}{2} \left((1 - \mu_h)^2 + (2 - \mu_h)^2 + (3 - \mu_h)^2 \right)$$

2' 1×1

$$\text{Var} = \frac{1}{n} \sum_i (x_i - \mu)^2$$
$$E \left[x - E[x] \right]^2$$

Covariance:

Cov_{dxd} for X with d features

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

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

$$\text{Cov} = \frac{1}{N} X^T X$$

$$\frac{1}{N} \begin{bmatrix} 1-\mu_h & 2-\mu_h & 3-\mu_h \\ 4-\mu_w & 5-\mu_w & 6-\mu_w \end{bmatrix}$$

$$\begin{bmatrix} 1-\mu_h & 4-\mu_w \\ 2-\mu_h & 5-\mu_w \\ 3-\mu_h & 6-\mu_w \end{bmatrix}$$

$$\text{Cov}(1,1) = \frac{1}{N} \left[(1-\mu_h)^2 + (2-\mu_h)^2 + (3-\mu_h)^2 \right] = \sigma_h^2$$

$$\text{Cov}(2,2) = \frac{1}{N} \left[(4-\mu_w)^2 + (5-\mu_w)^2 + (6-\mu_w)^2 \right] = \sigma_w^2$$

$$\text{Cov}(1,2) = \frac{1}{N} \left[(1-\mu_h)(4-\mu_w) + (2-\mu_h)(5-\mu_w) + (3-\mu_h)(6-\mu_w) \right]$$

$$\text{Cov}(2,1)$$

$$\sigma_{hw} = \sigma_{wh}$$

$$\text{if } \sigma_h^2 = \sigma_{hw} \Rightarrow h = w$$

Covariance:

$$\text{Cov} = \begin{bmatrix} \sigma_h^2 & \sigma_{hw} \\ \sigma_{wh} & \sigma_w^2 \end{bmatrix}$$

- Off-diagonal terms of covariance matrix same? *Yes*
- Diagonal terms higher in value or off diagonal term? *✓ Equal*
- Will eigenvectors of covariance matrix be orthonormal? *✓*
- Do we know its eigenvectors? *✓ from SVD of X*
- Any upper and lower bound in values of covariance matrix?
No bound

Covariance measures how two features vary together. Problem: covariance depends on scale

Correlation:

$$\sum = \frac{1}{N} X^{*T} X^*$$

$$X = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

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

$$\text{Cor}(i, i) = 1$$

Exploratory
Data
Analysis

$$\sum dxd = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

$X^* \rightarrow$ Standardized

Normalization $\rightarrow [0, 1]$
Standardization $\rightarrow [-1, 1]$

$$\frac{X - \mu}{\sigma}$$

$$\sum = \frac{1}{N} X^{*T} X^*$$

$$\begin{aligned} \sum (1,1) &= \frac{1}{N} \frac{(1 - \mu_h)^2 + (2 - \mu_h)^2 + (3 - \mu_h)^2}{\sigma_h^2} \\ &= \frac{\sigma_h^2}{\sigma_h^2} = 1 \end{aligned}$$

Useful in EDA. If features are correlated, then data may have redundancy.

For Joint Distributions

$$\text{cov}(X, Y) = E[XY] - E[X]E[Y]$$

- Expectation and Covariance:

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

- $\text{cov}(X, Y) = E[(X - E_X[X])(Y - E_Y(Y))] = E[XY] - E[X]E[Y]$

- $\text{Var}(X + Y) = \text{Var}(X) + 2\text{cov}(X, Y) + \text{Var}(Y)$

Uncorrelated vs Independent RV

Uncorrelated (Definition: $\text{Cov}(X, Y) = 0$)

- Means no **linear relationship** between X and Y.
- Does **not** rule out non-linear dependence.
- Example: $X \sim U(-1, 1)$, $Y = x^2 \rightarrow$ Uncorrelated but dependent.

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

Independent (Definition: $P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$)

- Stronger condition: knowing one variable gives **no information** about the other.
- Independence \Rightarrow Uncorrelated (if variances exist).

Key Differences

- Independence is a **stronger** property.
- Uncorrelated only removes **linear relationships**.
- Uncorrelated $\not\Rightarrow$ Independent.

ML Relevance

- Most ML models care about uncorrelatedness because they only model linear relationships.
- True independence is rarer and much harder to check, but it's what we assume in stronger probabilistic models

Outline

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

Gaussian Distribution

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

$$a = \mu$$
$$b = \sigma^2$$

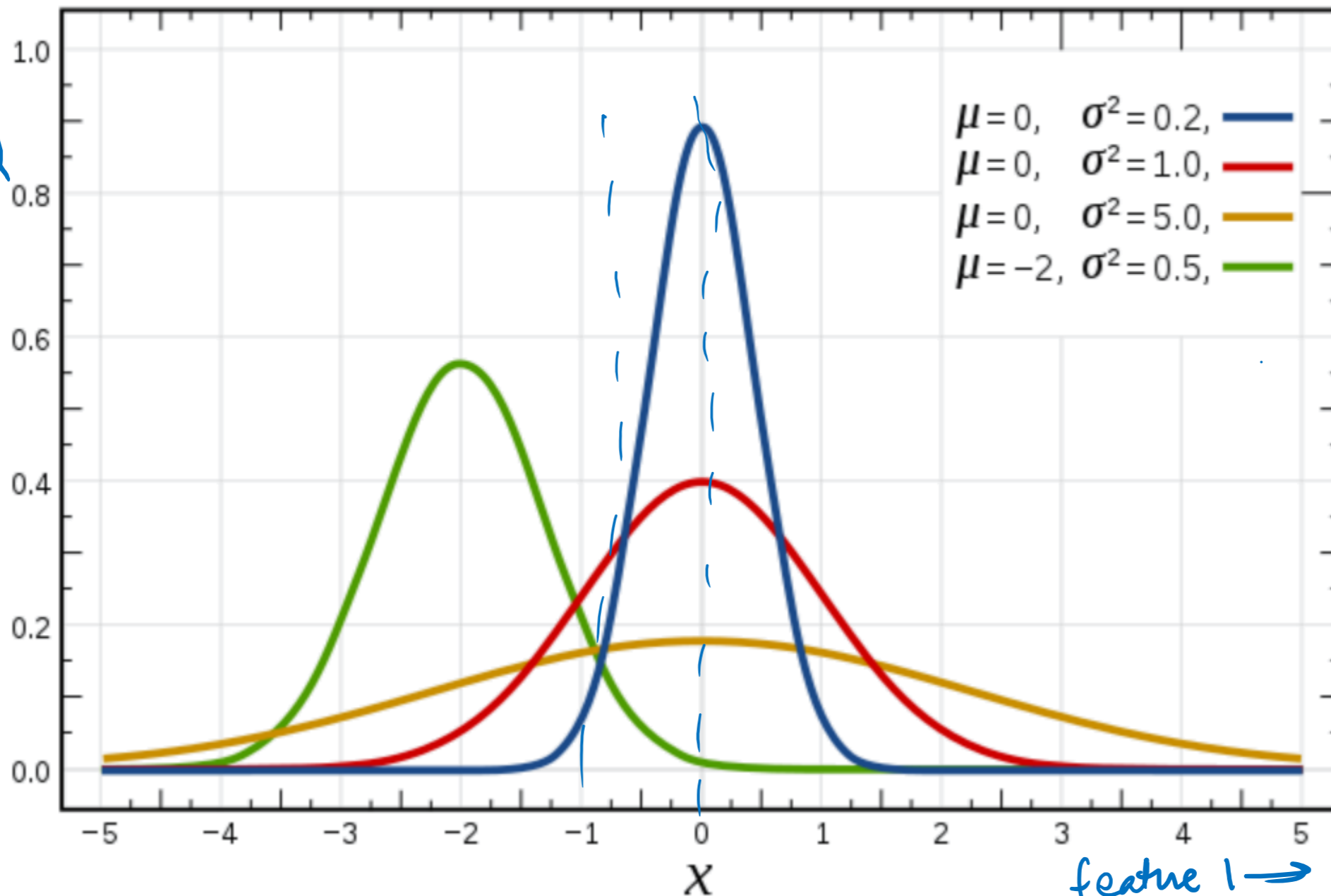
- Gaussian Distribution:

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

Probability density function

within 1 standard deviation

likelihood
or
density



$$x \pm \sigma \rightarrow 68\%$$

$$x \pm 2\sigma \rightarrow 95\%$$

$$x \pm 3\sigma \rightarrow 99.7\%$$

feature 1 \rightarrow

Multivariate Gaussian Distribution

$$X = \begin{bmatrix} \\ \end{bmatrix}_{n \times d}$$
$$\mu = \begin{bmatrix} \\ \end{bmatrix}_{1 \times d}$$
$$\sigma^2 = \text{Cov} = \begin{bmatrix} \\ \end{bmatrix}_{d \times d}$$

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

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

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

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

- 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

$$E(\overset{\text{r.v.} \sim N}{AX} + b) = AE(X) + b$$

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

this means that for Gaussian distributed quantities:

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

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

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

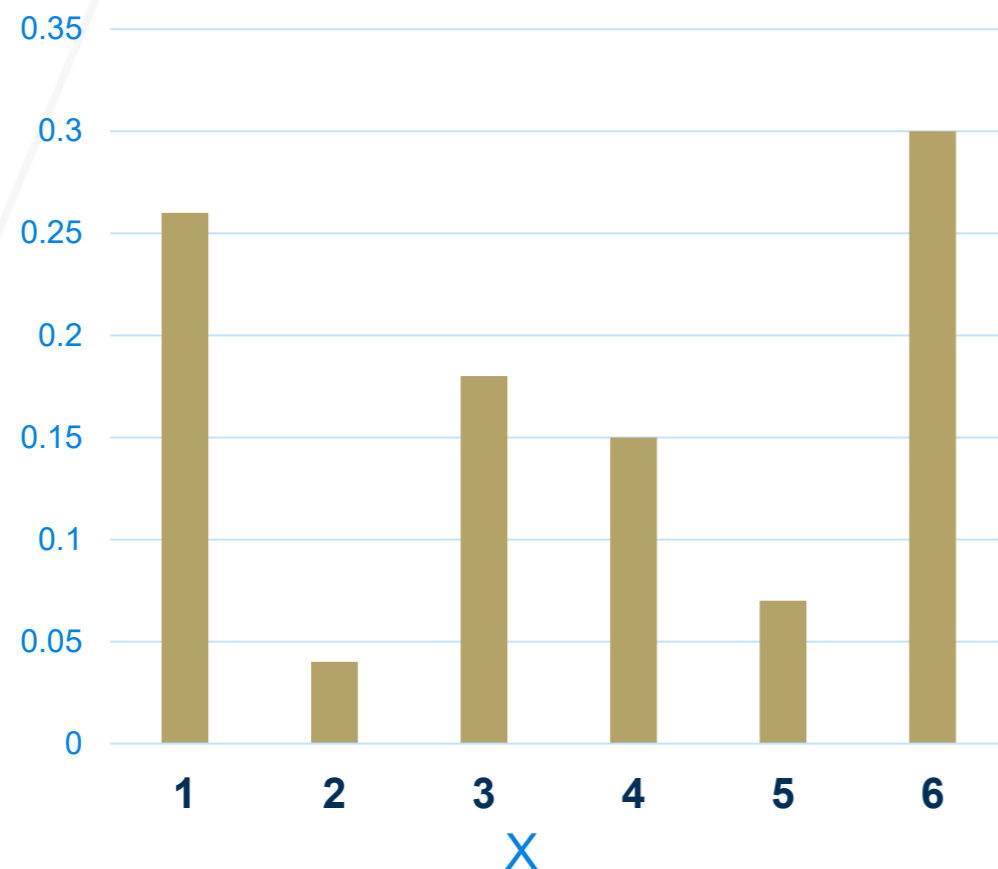
- The **multiplication** of two Gaussian functions is another Gaussian function (although no longer normalized)

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

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

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 $n = 4$

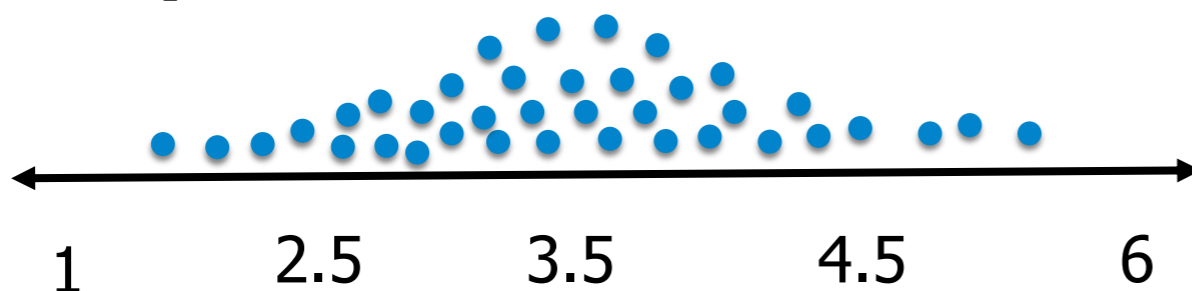
$$S_1 = \{1,1,1,6\} \Rightarrow E(S_1) = 2.25$$

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

⋮

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

- According to CLT, if you sample enough from any distribution with finite variance you will get an approximate Gaussian distribution.
- No matter what the population looks like, the average of many samples looks Normal.
- Explains **why the Normal distribution is everywhere** in statistics and ML.



Prob vs Likelihood

*Probability predicts data from parameters.
Likelihood evaluates parameters from data.*

•Probability

- Forward direction: given parameters, what's the chance of data?
- $P(data | \theta)$
- Example: *If the coin has bias $\theta = 0.7$, what's the probability of 8 heads in 10 tosses?*
- Varies with **different outcomes** of data.

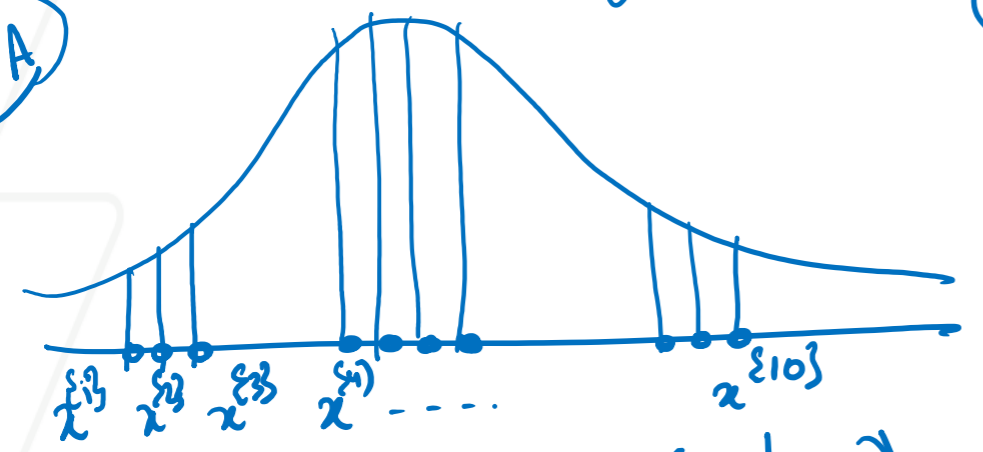
•Likelihood

- Reverse view: given data, how plausible are parameters?
- $L(\theta | data) \propto P(data | \theta)$
- Example: *I observed 8 heads in 10 tosses. How likely is it that the coin's bias is $\theta = 0.7$?*
- Varies with **different parameter values**.

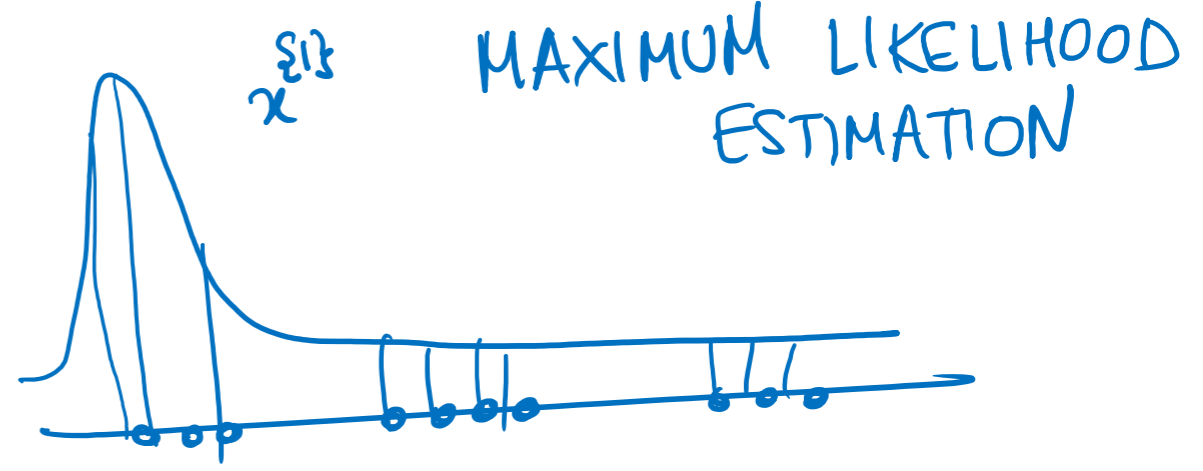
Prob vs Likelihood

$$\theta = \mu, \sigma^2$$

Case A



Case B



Likelihood $\rightarrow L(\theta | x)$

$$= L(a, b | x)$$

$$= P(x^{(1)}, x^{(2)}, \dots, x^{(n)} | \theta)$$

maximize

$$L(\theta | x) = \max P(x^{(1)}, x^{(2)}, \dots, x^{(n)} | \theta)$$

Do NOT LIKE LIST

- ① for loop
- ② Inverse
- ③ Matrix Mult.
- ④ Joint Probs

$$= P(x^{(1)} | \theta) \cdot P(x^{(2)} | \theta) \dots P(x^{(n)} | \theta)$$

$$= \prod_i P(x^{(i)} | \theta)$$

iid assumption
 \rightarrow all datapoints are independent & identically distributed

Prob vs Likelihood

$$L(\theta | x)$$

$\log abc$
 $\log a + \log b + \log c$

$$\prod_i P(x^{\epsilon_{i3}} | \theta)$$

Maximize this probability
↓
most optimal values of a, b

$$\log L(\theta | x) = \log \prod_i P(x^{\epsilon_{i3}} | \theta)$$

$$\log L(\theta | x) = \sum_i \log P(x^{\epsilon_{i3}} | \theta)$$

$$\frac{\partial \log L}{\partial a} = 0$$

$$\frac{\partial \log L}{\partial b} = 0$$

a, b

Maximum Likelihood Estimation

Outline

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

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

$$\theta = 0.5$$

Maximum Likelihood Estimation

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

Objective function: $P(x^{i} | \theta) = \theta^{x^{i}} (1 - \theta)^{1-x^{i}}$ $x^{i} \in \{0,1\}$ 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

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

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

$$\begin{aligned} & a^b \cdot a^c \cdot a^d \\ & = a^{b+c+d} \end{aligned}$$

$$\begin{aligned} L(\theta | X) &= \theta^{x^{1}} (1 - \theta)^{1-x^{1}} \times \theta^{x^{2}} (1 - \theta)^{1-x^{2}} \dots \times \theta^{x^{n}} (1 - \theta)^{1-x^{n}} = \\ &= \theta^{\sum x^{i}} (1 - \theta)^{\sum (1-x^{i})} \end{aligned}$$

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

What's the trick?

Take the log

$$\log(a^b) = b \log a$$

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

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

How to optimize θ ?

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

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

- If your data has **70% ones**, the probability that best explains it is **$\theta = 0.7$**
- Any other θ would make the observed pattern *less plausible*