


Probability and Statistics

Mahdi Roozbahani

Outline

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

Probability

$$P(\text{event}) = 1/6$$

we don't know

- 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 (arrow pointing to X)
specific outcome (arrow pointing to x)

~~$p(x) = p$~~

$$p(x) \geq 0$$

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

Any dist. function is a probability function if area under curve = 1

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

density or likelihood



$$p(x = x_t) = 0$$

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**.

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



$f(x) \rightarrow$ objective function needs to be optimized
 Parameters θ

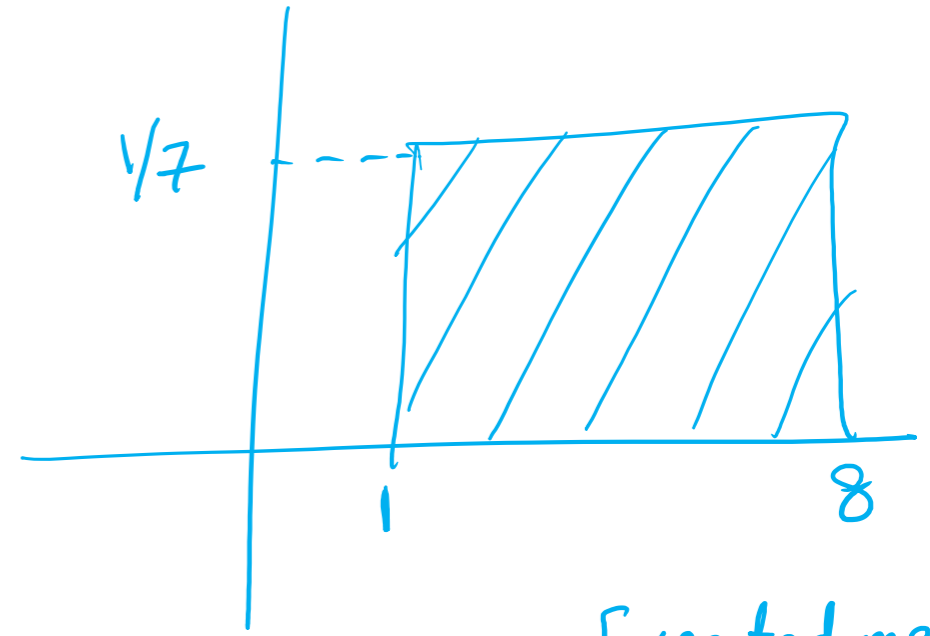
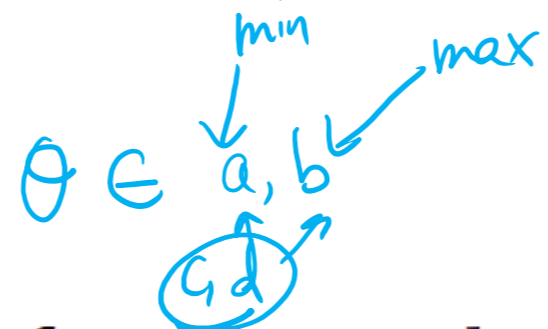
Continuous Probability Functions

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

• Examples:

• Uniform Density Function:

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



• Exponential Density Function:

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

$\theta \in a$
 for $x \geq 0$
 $f(x) = \frac{1}{a} e^{-x/a}$
 $a = \mu$

μ is average \rightarrow Expected mean or weighted average
 \rightarrow arithmetic mean

• Gaussian(Normal) Density Function

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

$\theta \in \mu, \sigma$

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

$y_{j=2} = tail$

$y_{j=1} = head$

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	N=35

JOINT PROBABILITY

$$P(x=1, y=t) = \frac{3}{35} = \frac{n_{ij}}{N}$$

MARGINAL PROBABILITY

$$P(x=4) = \frac{7}{35} = \frac{C_i}{N}$$

$$P(y=t) = \frac{20}{35} = \frac{C_j}{N}$$

SUM RULE: $P(x) = \sum_y P(x,y)$

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

CONDITIONAL PROBABILITY

$$P(x=3 | y=t) = \frac{2}{20} = \frac{n_{ij}}{C_j}$$

$$P(y=t | x=3) = \frac{2}{6} = \frac{n_{ij}}{C_i}$$

Sample shranked from 35 outcomes to 20

x is a more informative feature than y.

PRODUCT RULE:

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

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

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

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

$$P(X|Y) = P(X) \rightarrow X \text{ is C.I. of } Y.$$

Conditional Independence

- Examples:

Virus is independent conditionally of drinking beer

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

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:

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

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

Outline

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

$$\begin{array}{c}
 P(z|y) P(y) \\
 P(y|z) P(z) \\
 \hline
 P(x|y, z) \underbrace{P(y, z)} \\
 \uparrow \\
 \frac{P(x, y, z)}{P(x, z)} \\
 \hline
 P(x|z) P(z)
 \end{array}$$

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

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

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

$$\underbrace{P(y|x)} = \frac{P(x, y)}{P(x)} = \frac{P(x|y) P(y)}{P(x)} = \sum_y P(x, Y=y) = \sum_{Y=\text{cat, dog}} P(x|Y=y) P(Y=y)$$

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:

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

Outline

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

$E[\cdot] \rightsquigarrow$ expected value
weighted avg

$E(\cdot) \rightsquigarrow$ function

Mean and Variance

$g(\cdot)$
 $h(\cdot)$

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

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

$$\sum g(x)p(x)$$

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

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

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

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

$$E[a+b+c] = E[a] + E[b+c]$$

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

- Variance(Second central moment): $Var(x) = E[(x - E[x])^2]$

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

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

$$= \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

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

$$Var(x) = E[x^2] - E[x]^2$$

Expected value and average

$$f(x) = x$$

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

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

$$\Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N (x^{(i)}) =$$

$$E[f(x)] = \sum_{i=1}^3 P(x) f(x) = P(x=1)f(x=1) + \dots + P(x=3)f(x=3)$$

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

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

$$E[f(x)] \neq \mu$$

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

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

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:

$$X = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}_{n \times d} \quad h = \text{height}$$

$$\Rightarrow \mu_{x_h} = 2$$

$$\sigma_h^2 = \frac{1}{N} \sum_{i=1}^N (x_i^{\{h\}} - \mu_{x_h})^2$$

$$\sigma_h^2 = \left(\frac{1}{N} \sum_{i=1}^N (\cdot) \right) = E[\cdot]$$

$$x^{\{1\}} = 1$$

$$x_1 = [1, 2, 3]$$

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

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

$\mu_{\bar{h}} = 0$

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

Covariance:

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

$$X = \begin{matrix} \{1\} \\ \{1, 4\} \end{matrix} \quad X_1 = [1, 2, 3]$$

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

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

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

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

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

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

$$\sigma_{hw} = 0$$

Correlation:

$$\text{Cor} = \frac{1}{N} \bar{X}^{*T} \bar{X}^*$$

$$\bar{X} = \begin{bmatrix} \bar{x} & \bar{w} \\ 1 - \mu_h & 1 - \mu_w \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix} \xrightarrow{\text{Standardization}} \bar{X}^* = \begin{bmatrix} \frac{1 - \mu_h}{\sigma_h} \\ \vdots \\ \frac{1 - \mu_w}{\sigma_w} \\ \vdots \end{bmatrix}$$

$$\text{Cor} = \frac{1}{N} \begin{bmatrix} \frac{1 - \mu_h}{\sigma_h} & \dots \\ \frac{1 - \mu_w}{\sigma_w} \\ \vdots \\ \vdots \end{bmatrix}$$

$$\text{Cor} = \begin{matrix} h & w \\ \left[\begin{matrix} \frac{\sigma_w^2}{\sigma_h^2} = 1 \\ \frac{\sigma_h^2}{\sigma_w^2} = 1 \end{matrix} \right] & -1 \leq r_{hw} \leq 1 \end{matrix}$$

EDA

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

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.

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

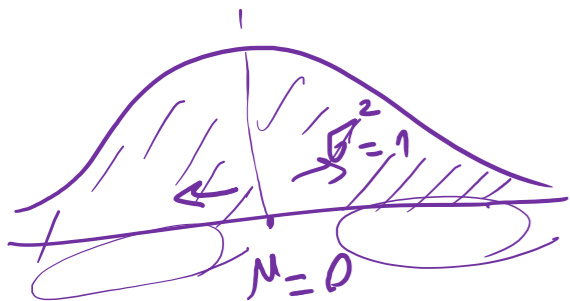
For Joint Distributions

$$\text{Cov}(x, y) = \cancel{E[z^3]} - \cancel{E[z]} \cdot E[z^2] = 0$$

Expectation and Covariance:

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

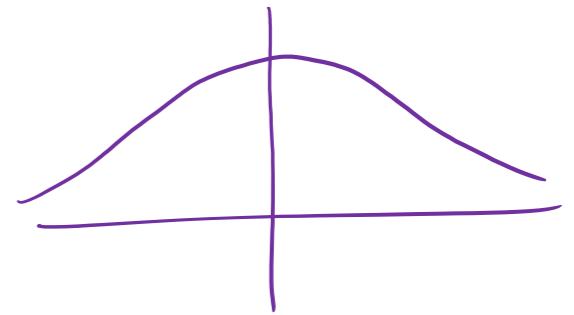
$X = z \rightsquigarrow$ Standard gaussian
 $\mu = 0 \quad \sigma^2 = 1$



$Y = z^2 \rightsquigarrow$ chi-squared



$h = z^3$



$\text{Cov}(X, Y) =$

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

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

$$E[z^2] = 1$$

Outline

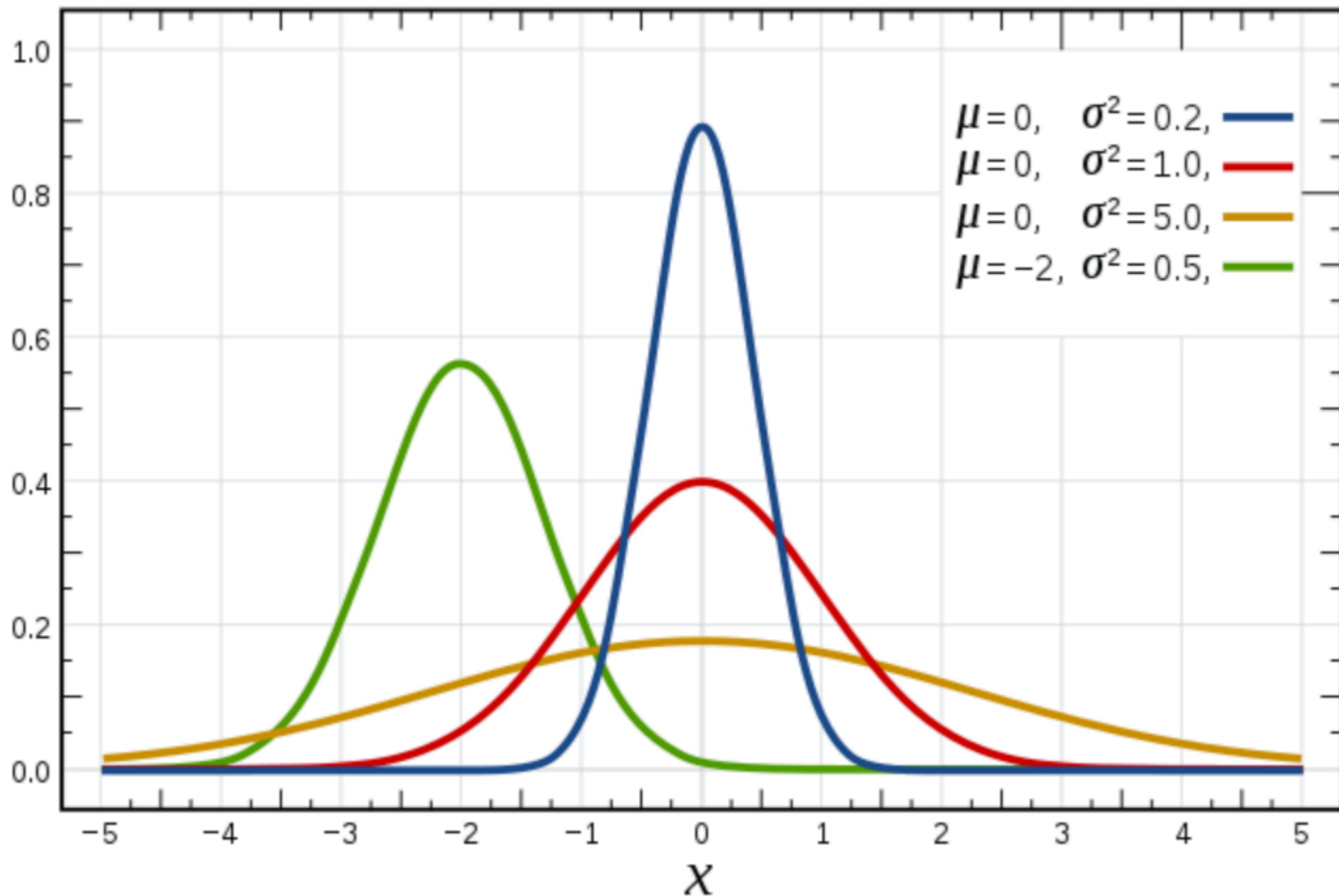
- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution ←
- Maximum Likelihood Estimation

Gaussian Distribution

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

- 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



=h

Probability versus likelihood

Multivariate Gaussian Distribution

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

Handwritten notes: A purple circle is drawn around Σ^{-1} in the exponent, with an arrow pointing to the word "Cov" written above it. Another purple circle is drawn around $|\Sigma|$ in the denominator, with an arrow pointing to the word "Cov" written below it.

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

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

- 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

$$(AX) + b = Y$$

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

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

↓ gaussian

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

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

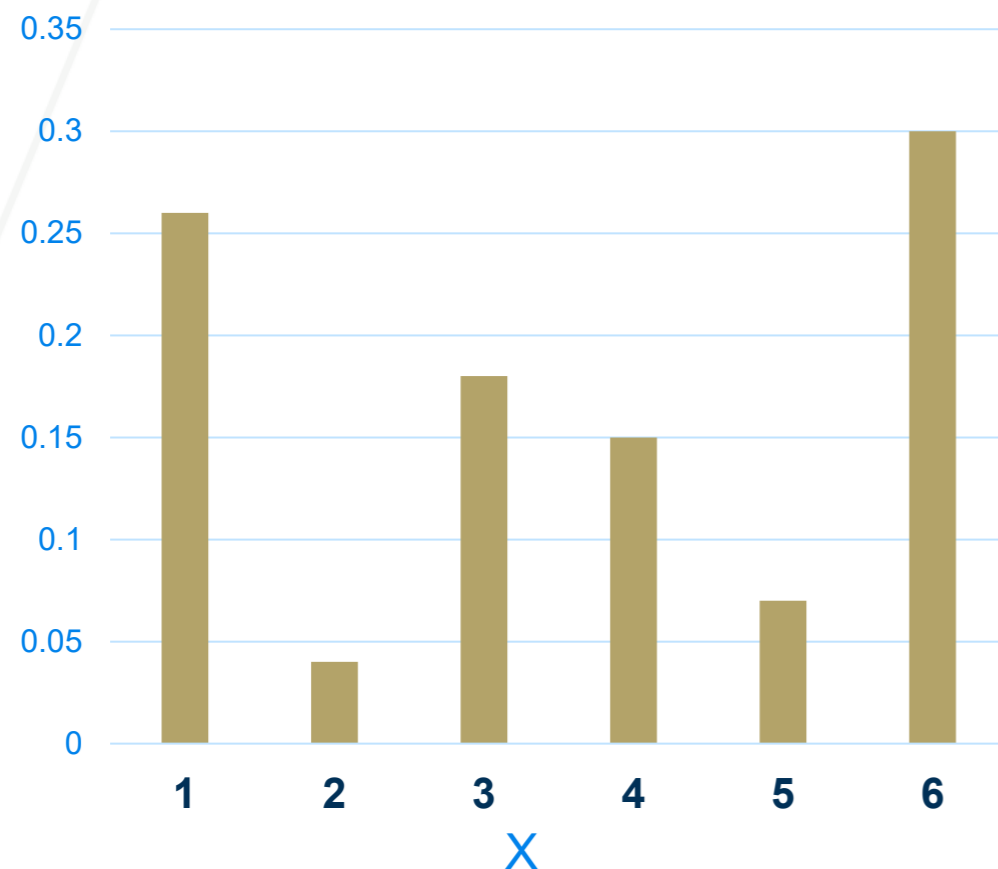
- The **multiplication** of two Gaussian functions is another Gaussian function (although no longer normalized)

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

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

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 $n = 4$

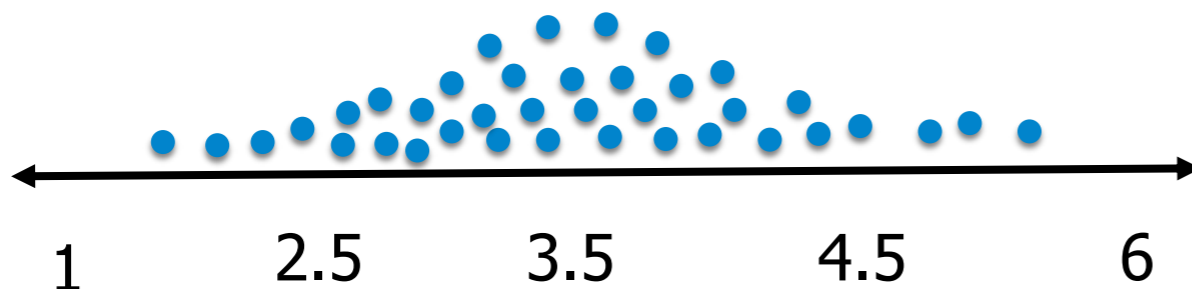
$$S_1 = \{1,1,1,6\} \Rightarrow E(S_1) = 2.25$$

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

⋮

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

- 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.

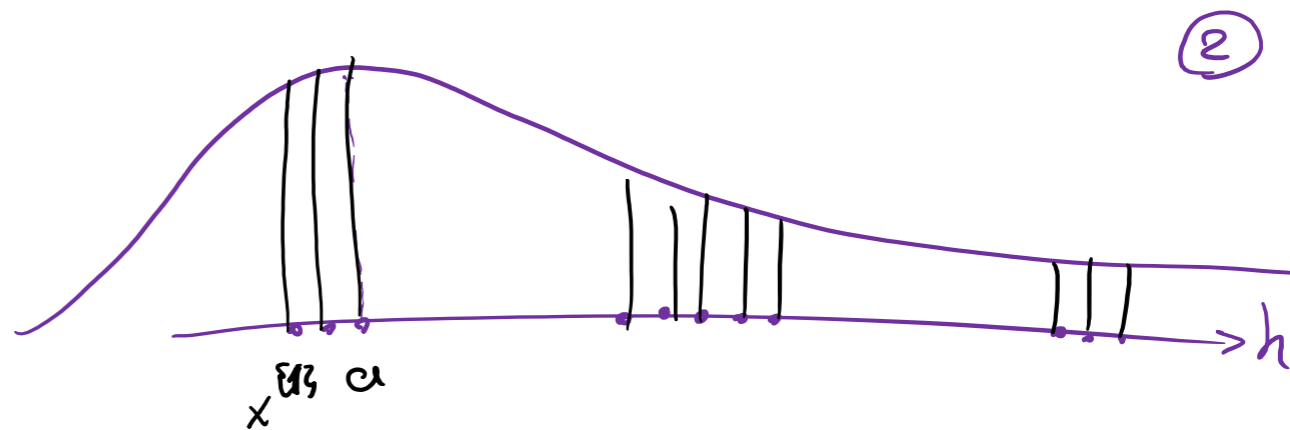
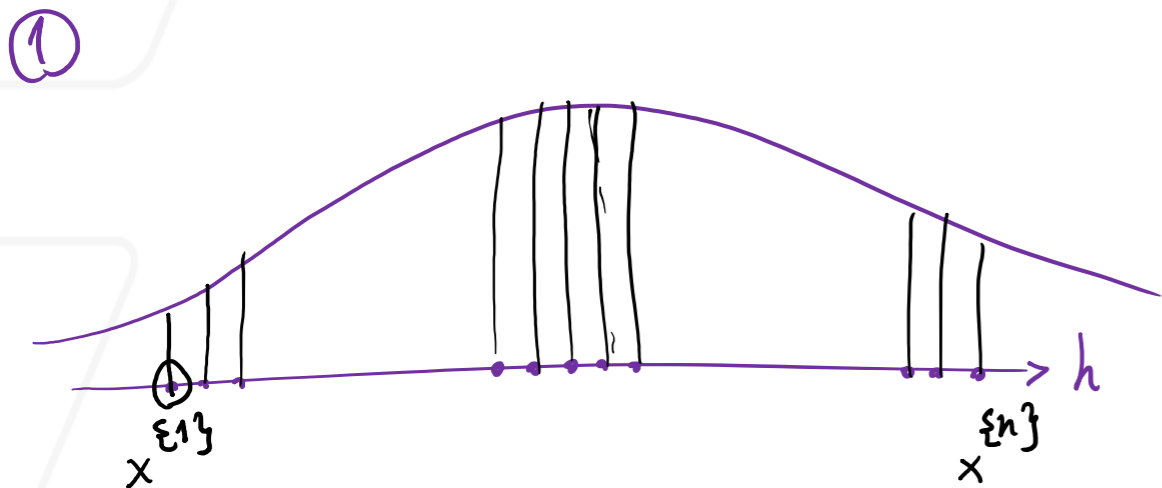
•Likelihood

- 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

$$P(X=T) = \frac{1}{2}, \quad T, T, T, H \Rightarrow P(X=T) = \frac{3}{4}$$

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



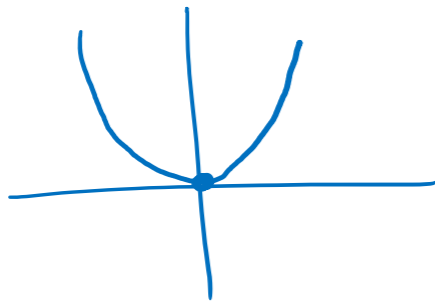
$$L(\theta | x^{i3}, \dots, x^{n3}) = L(a, b | x) \rightarrow \underline{f(x^{i3}, \dots, x^{n3} | \theta) = f(x | \theta) = f(x | a, b)}$$

$$\left[f(a, b | c) = f(a | c) f(b | c) \right]$$

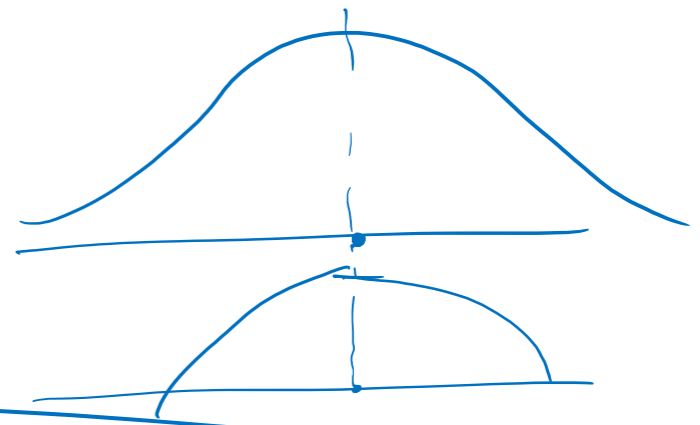
$$\text{Max } L(a, b | x) = \text{Max} f(x^{i3} | \theta) f(x^{23} | \theta) \dots f(x^{n3} | \theta) = \text{Max} \prod_{i=1}^N f(x^{i3} | \theta) \uparrow \{a, b\}$$

Prob vs Likelihood

$$f(x) = x^2$$



$$\frac{\partial f(x)}{\partial x} = 2x = 0 \Rightarrow x = 0$$



Max
 $L(\theta|x) = \prod_{i=1}^N f(x^{i3} | a, b)$ Maximizing Likelihood

Maximizing log likelihood

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

$$\max l(\theta|x) = \max \log L(\theta|x)$$

Objective function

$$= \max \sum_{i=1}^N \log f(x^{i3} | a, b)$$

always scalar

\Rightarrow MLE

$$\frac{\partial l(a, b|x)}{\partial a} = 0 \Rightarrow a = \dots$$

$$a = \frac{1}{N} \sum_{i=1}^N x^{i3} = \mu$$

$$\frac{\partial l(a, b|x)}{\partial b} = 0 \Rightarrow b = \dots$$

$$b = \frac{1}{N} \sum_{i=1}^N (x^{i3} - \mu)^2 = \sigma^2$$

~~$X = \begin{bmatrix} h \\ 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$~~

cut
data point

$f(x_{\text{cut}} | \mu, \sigma)$

~~$X_{\text{dogs}} = \begin{bmatrix} h \\ 4 \\ 5 \\ 6 \\ 7 \end{bmatrix}$~~

$f(x_{\text{dog}} | \mu, \sigma)$

- ML \rightarrow
- ① what is the problem you are trying to solve
 - ② we need to come up with an Objective function
 - ③ training \rightsquigarrow Optimization
 - ④ Testing

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

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

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

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

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

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

What's the trick?
 $\log a^b = b \log a$

Take the log

$$\frac{\partial \log(\theta)}{\partial \theta} = \frac{1}{\theta}$$

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

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

How to optimize θ ?

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

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

1, 1, 1, 1, 1, 1 0, 0, 0

$$\theta = \frac{1+1+1 \dots + 0+0}{10} = 0.7$$