

# MAXIMUM LIKELIHOOD ESTIMATION

## Motivation example:

Given a biased coin whose probability of heads is unknown

$$P(X = H) = \theta$$

$$P(X = T) = 1 - \theta$$

And given these independence observations : H H H H T

Problem: Estimate  $\theta$

## Solution:

• We find the probability of observing the dataset given some value of parameter  $\theta$ . For example:

$$\bullet P(HHHHT \mid \theta = 0.9) = \theta^4 (1 - \theta)^1$$

$$\bullet P(HHHHT \mid \theta = 0.8) = \dots$$

$$\bullet P(HHHHT \mid \theta = 0.7) = \dots$$

↓  
Continuing this pattern, we some function over  $\theta$

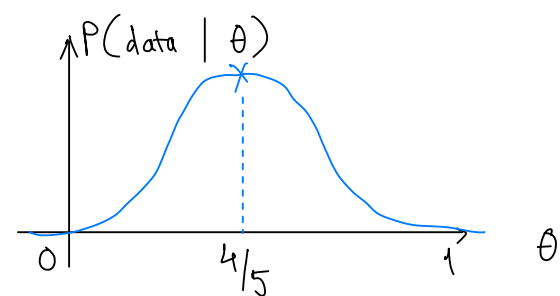
$$L(\theta) = P(HHHHT \mid \theta) \rightarrow \text{Likelihood Function} \\ = \theta^4 (1 - \theta)^1$$

To best estimate  $\theta$ , we maximize this function:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$$

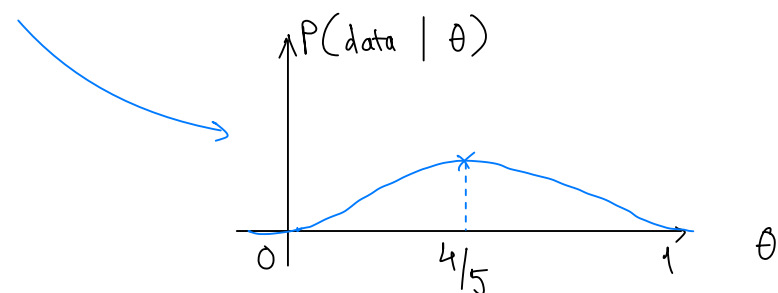
$$= \underset{\theta}{\operatorname{argmax}} \theta^4 (1 - \theta)^1$$

→ Maximum Likelihood Estimation  $L(\theta)$



In practice, we often maximize the log function:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log(\theta^4 (1 - \theta)^1) \rightarrow \text{Log-likelihood Function } l(\theta)$$



Solve for best value of  $\theta$ :

$$\nabla_{\theta} [\log(\theta^4 (1 - \theta))] = 0$$

$$\Leftrightarrow \nabla_{\theta} [4 \log \theta + \log(1 - \theta)] = 0$$

$$\Leftrightarrow \frac{1}{\theta} 4 + \frac{-1}{1 - \theta} = 0$$

$$\Rightarrow \hat{\theta} = \frac{4}{5}$$

## Generalized approach of MLE:

- Given a set of observation / dataset  $S = \{x_i\}_{i=1}^n$
- Independent and identically distributed (i.i.d.) following an unknown distribution, from a parametric family of distributions:  
 $\{p(\cdot | \theta) : \theta \in \Theta\}$

Problem: Estimate  $\theta$

### Solution:

- Write down likelihood function:

$$\begin{aligned} L(\theta) &= P(S = \{x_1, x_2, \dots, x_n\} | \theta) &<\text{definition}> \\ &= P(x_1 | \theta) \cdot P(x_2 | \theta) \dots P(x_n | \theta) &<\text{i.i.d.}> \\ &= \prod_{i=1}^n P(x_i | \theta) \end{aligned}$$

- Then write down the log-likelihood function:

$$\begin{aligned} l(\theta) &= \log L(\theta) \\ &= \log \left( \prod_{i=1}^n P(x_i | \theta) \right) \\ &= \sum_{i=1}^n \log P(x_i | \theta) \end{aligned}$$

- Define maximum likelihood estimation (MLE):

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmax}} l(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P(x_i | \theta) \end{aligned}$$

- Solve the optimization problem, there are 2 ways:

$\left\{ \begin{array}{ll} \text{Closed form} & \text{if we have "nice form" function} \\ \text{Numerical algorithm} & \text{otherwise} \end{array} \right.$

↓ a common algorithm  
is Gradient Ascent

### Example: Gaussian Distribution

- Given  $\{x_i\}_{i=1}^n$  i.i.d. drawn from Gaussian Distribution  $N(\mu, \sigma^2)$ :
$$p(x | \theta) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^2\right), \quad \theta = \{\mu, \sigma\}$$

- Problem: Estimate  $\theta$

### Solution:

- Write down log-likelihood function:

$$\begin{aligned} l(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^n \log p(x_i | \theta) \\ &= \sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2\sigma^2} (x_i - \mu)^2\right) \right] \\ &= \sum_{i=1}^n \left[ \log\left(\frac{1}{\sqrt{2\pi} \sigma}\right) - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right] \\ &= n \cdot \log\left(\frac{1}{\sqrt{2\pi} \sigma}\right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\ &= -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

- Define maximum likelihood estimation (MLE):

- Since  $\theta = \{\mu, \sigma\}$ :

- First, find  $\hat{\mu}$ :

$$\hat{\mu} = \operatorname{argmax}_{\mu} l(\theta)$$

Solve using closed form:

$$\begin{aligned} \nabla_{\mu} l(\theta) &= 0 \\ \Leftrightarrow \nabla_{\mu} \left[ -\sum_{i=1}^n (x_i - \mu)^2 \right] &= 0 \end{aligned}$$

$$\Leftrightarrow 2 \left[ \sum_{i=1}^n (x_i - \hat{\mu}) \right] = 0$$

$$\Leftrightarrow \sum_{i=1}^n (x_i - \hat{\mu}) = 0$$

$$\Leftrightarrow \sum_{i=1}^n x_i - \sum_{i=1}^n \hat{\mu} = 0$$

$$\Leftrightarrow \sum_{i=1}^n x_i - n \cdot \hat{\mu} = 0$$

$$\Leftrightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{mean value of dataset } S)$$

- Second, find  $\hat{\sigma}$ :

$$\hat{\sigma} = \operatorname{argmax}_{\sigma} l(\theta)$$

Also solve using closed form:

$$\nabla_{\sigma} l(\theta) = 0$$

$$\Leftrightarrow \nabla_{\sigma} \left[ -n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] = 0$$

$$\Leftrightarrow -\frac{n}{\hat{\sigma}} + \frac{1}{\hat{\sigma}^3} \cdot \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Leftrightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad (\text{mean variance of dataset } S)$$

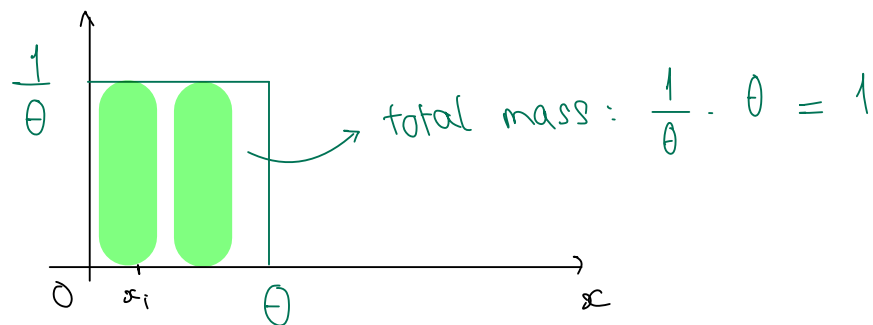
• Conclusion :

$$\text{MLE gives } \begin{cases} \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \end{cases}$$

### Example: Uniform Distribution

- Given  $\{x_i\}_{i=1}^n$  i.i.d. drawn from a uniform distribution  $\text{Uniform}([0, \theta])$   
$$p(x | \theta) = \begin{cases} \frac{1}{\theta} & \text{if } x \in [0, \theta] \\ 0 & \text{otherwise} \end{cases}$$

Visualize:



- Problem: Estimate  $\theta$

### Solution:

- Since some probability distribution is 0, using likelihood function is easier in this case (because likelihood function use products):

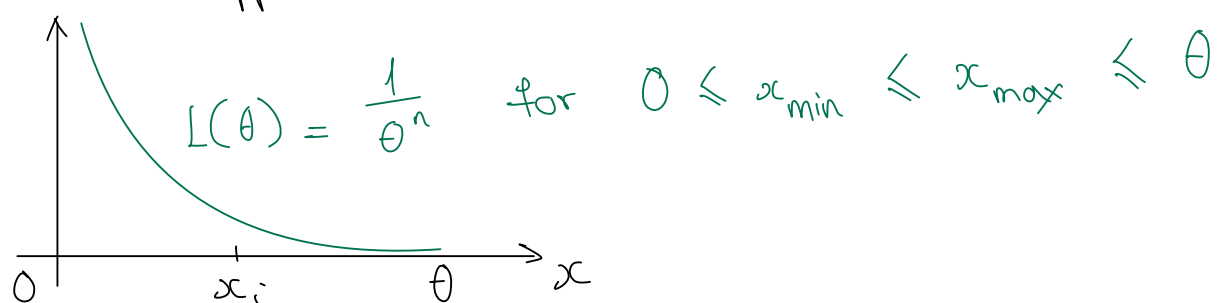
$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(x_i | \theta) \\ &= \prod_{i=1}^n \frac{1}{\theta} \mathbb{1}(x_i \in [0, \theta]) \\ &= \left(\frac{1}{\theta}\right)^n \prod_{i=1}^n \mathbb{1}(x_i \in [0, \theta]) \\ &= \begin{cases} \left(\frac{1}{\theta}\right)^n & \text{if } x_i \in [0, \theta] \quad \forall i \\ 0 & \text{otherwise} \end{cases} \\ \text{or} &= \frac{1}{\theta^n} \quad \text{for } 0 \leq x_{\min} \leq x_{\max} \leq \theta \end{aligned}$$

- Try to solve using closed form:

$$\begin{aligned} \nabla_{\theta} L(\theta) &= 0 \\ \Leftrightarrow \nabla_{\theta} \left(\frac{1}{\theta}\right)^n &= 0 \\ \Leftrightarrow -n \left(\frac{1}{\theta}\right)^{n+1} &= 0 \end{aligned}$$

↓ always  $\neq 0$

- Try another approach, let's visualize  $L(\theta)$ :



$L(\theta)$  is a decreasing function of  $\theta$ , so it doesn't have a max value.

But we can still find the "allowed" maximum, the "allowed" maximum occurs at the minimum  $\theta$  (as showed in graph):

$$\max L(\theta) = \min \theta$$

And since the constraint on  $\theta$  is:

$$\theta \geq x_{\max}$$

Therefore:

$$\max L(\theta) = \min \theta = x_{\max}$$

meaning:

$$\hat{\theta} = x_{\max}$$

or:

$$\hat{\theta} = \{x_1, \dots, x_n\}$$

• Conclusion:

The best estimate of  $\theta$  is one that:

- covers all data point :  $\theta \geq x_{\max}$
- as tight as possible :  $\theta = x_{\max}$

## Example: MLE for (Linear) Regression

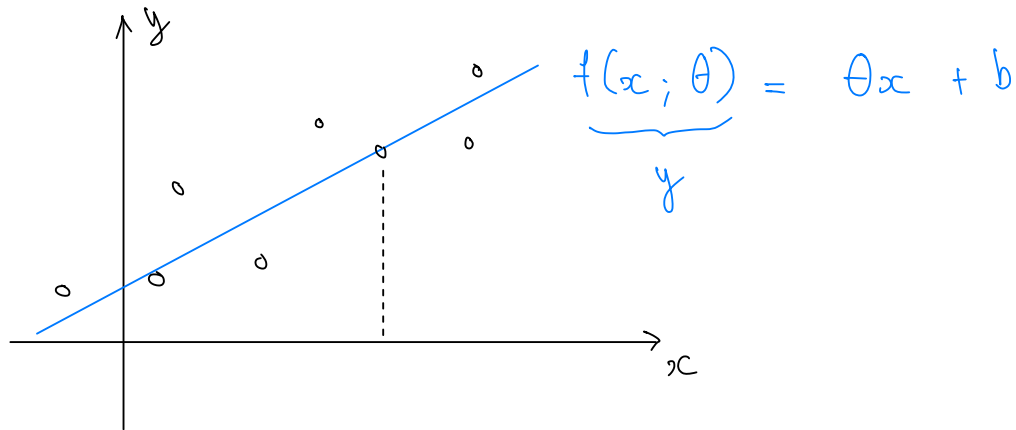
- Given  $\{x_i, y_i\}_{i=1}^n$ , where  $x_i$ , such that:

$$p(y|x, \theta) = N(y | f(x; \theta), \sigma^2) \rightarrow \text{why?}$$

some line  
to fit  
the dataset

We usually assume  
 $y_i$  follow Gaussian  
Distribution

Visualize:



- Problem: Estimate  $\theta$  and  $\sigma^2$

Solution:

- Write down log-likelihood function:

$$l(\theta, \sigma) = \sum_{i=1}^n \log P(y_i | x_i, \theta)$$

$$= -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n [y_i - f(x_i, \theta)]^2$$

- Define maximum likelihood estimation (MLE):

- For  $\hat{\theta}$ :

$$\hat{\theta} = \arg\max_{\theta} l(\theta, \sigma)$$

$$= \arg\max_{\theta} \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n [y_i - f(x_i, \theta)]^2 \right)$$

$$= \arg\min_{\theta} \left( \sum_{i=1}^n [y_i - f(x_i, \theta)]^2 \right)$$

$$= \arg\min_{\theta} \|y - f(x, \theta)\|^2$$

Linear Least Square

- For  $\hat{\sigma}$ :

$$\hat{\sigma} = \arg\max_{\sigma} l(\theta)$$

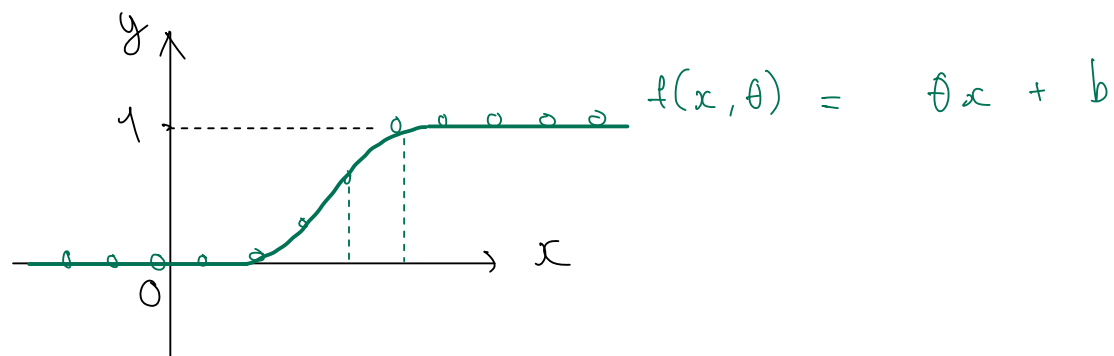
$$\text{Solve for } \hat{\sigma}: \quad \nabla_{\sigma} l(\theta, \sigma) = 0$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n [y_i - f(x_i, \hat{\theta})]^2 \rightarrow \text{mean variance of } y$$

## Example: MLE for Logistic Regression

- Given  $\{x_i, y_i\}_{i=1}^n$ , where  $y_i \in [0, 1]$  such that:
$$p(y | x, \theta) = \frac{\exp(y \cdot f(x, \theta))}{1 + \exp(f(x, \theta))}$$

Visualize:



- Problem: Estimate  $\theta$

Solution:

- Write down log-likelihood function:

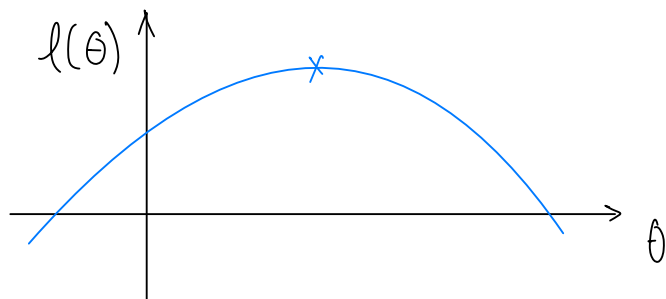
$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \log P(y_i | x_i, \theta) \\ &= \sum_{i=1}^n \log \left[ \frac{\exp(y_i f(x_i, \theta))}{1 + \exp(f(x_i, \theta))} \right] \\ &= \sum_{i=1}^n \left[ y_i f(x_i, \theta) - \log(1 + \exp(f(x_i, \theta))) \right] \end{aligned}$$

- Define maximum likelihood estimation:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} l(\theta)$$

Since  $l(\theta)$  has a term with  $\log$ , this function curves exponentially, so you won't find a nice solution using closed form method.

Instead, we converge slowly to the optimal solution using gradient "ascent", this is called numerical method.



Gradient update rule:

$$\theta_{t+1} = \theta_t + \eta \nabla_{\theta} l(\theta)$$