

Parameter Estimation

Sunita Sarawagi

CS 215. Fall 2024

So far..

- Computing probabilities of outcomes given a fixed distribution.
- Distributions were given to us as a function..
- Functions had parameters with fixed values

What are Parameters?

Consider some probability distributions:

- $\text{Ber}(p)$
- $\text{Poi}(\lambda)$
- $\text{Uni}(\alpha, \beta)$
- $\text{Normal}(\mu, \sigma^2)$
- $Y = mX + b$
- etc...

$X \sim N(0, 1)$

$$\theta = p$$

$$\theta = \lambda$$

$$\theta = (\alpha, \beta)$$

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

$$\theta = (m, b)$$

Call these “parametric models”

Non parametric model example - histogram

Given model, **parameters** yield actual distribution

- Usually refer to parameters of distribution as θ
- Note that θ that can be a vector of parameters

Today's class

How to determine the values of the parameters.



Parameters differ based on the task and application. These are not fixed like the speed of light.

The setup for parameter estimation in real-life

- Step 1: A real-life problem:

- 1. Estimating the probability that at least two out of four servers will be alive next day ✓
- 2. The probability that stock price will rise by 10% in the next week
- 3. The expected number of clicks on an advertisement in the next 3 hours

- Step 2: Model the problem: Choose a functional form of the uncertainty.

1. Binomial?

Assume that servers fail independently
 $X = \#$ of failures in a day $X \sim \text{Bin}(C)$

2. Gaussian?

$X =$ change from ^{one} day to the next

3. Poisson?

$X = \#$ of clicks on the ad per hour

The setup for parameter estimation in real-life

Step 3: Collect a training sample by observing over several days.

1. Sample server failure data observed over 3 days

$\left\{ \begin{array}{l} \text{day 1 ser 1} \\ \text{day 1 ser 2} \end{array} \right. \begin{array}{l} x_1 \\ x_2 \end{array} \quad \begin{array}{l} \text{day 1 ser 2} \\ \text{day 3 ser 3} \end{array} \begin{array}{l} x_3 \\ x_{12} \end{array}$
 $\begin{array}{cccccccccccc} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & \dots & 0 \end{array}$

2. Stock price change over a 10 days

change $1 \rightarrow 2$ $2 \rightarrow 3$ - - - - - $9 \rightarrow 10$
 1% -2%

3. Number of clicks on the ad over the last 20 hour

Hour $\begin{array}{ccc} 1 & 2 & 3 \end{array}$ - - - - - 20
 $\begin{array}{ccc} 10 & 15 & 5 \end{array}$ 7
 x_1 x_2 x_3 x_{20}

• Step 4: Estimate the unknown parameters using the training sample

The overall setup in parameter estimation

- Given: a density or distribution function with parameters $f(x, \theta)$ density of pmf
- Given: sample: $D = \{x_1, x_2, \dots, x_N\}$
 - The i -th sample is a random variable X_i assumed to be independently identically distributed as per the unknown $f(x, \theta)$
- Find θ .

- Since D is a finite sample, we cannot really know the actual θ . Best we can do is obtain an estimate of θ .
- We will denote the estimate as $\hat{\theta}$
- Goodness of estimate will be discussed later.

Types of estimators

- Maximum likelihood: sample D is all you got.

$\hat{\theta}$ point estimator

- Bayesian estimation: in addition to sample, we got prior beliefs.

Maximum Likelihood Estimation

- If θ were known we could have calculated the probability of getting the N outcomes in $D = \{x_1, x_2, \dots, x_N\}$ from the distribution as

- $P(D|\theta) = P(x_1, \dots, x_N|\theta) = \prod_i P(x_i|\theta) = \prod_i f(x_i; \theta)$ ** for both continuous & discrete*

- Likelihood refers to the above function. Often denoted as $L(\theta)$

- Maximum likelihood estimator:

- Choose the parameter θ for which the above likelihood is maximized

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N f(x_i, \theta) \rightarrow$$

Finding θ that maximizes likelihood

- Use log-likelihood instead of likelihood to convert products into sums

- $LL(\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log f(x_i, \theta)$
sum over observations

$\max_{\theta} LL(\theta)$

- Maximum likelihood estimator

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log f(x_i | \theta)$$

Solved using numerical optimization methods applying calculus.

MLE for Bernoulli

$$X \sim \text{Bern}(\underline{p}) \quad x \in \{0, 1\}$$

$$\underline{f(x; p)} = p^x (1 - p)^{1-x}$$

Data sample: D $N=10$

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
0	1	1	0	0	0	0	0	1	1

$$\begin{aligned} \max_P \underline{LL}(\theta \equiv [p]) &= LL_D(p) = \max_P \sum_{i=1}^N \log P^{x_i} (1-P)^{1-x_i} \\ &= \max_P \sum_{i=1}^N x_i \log P + \left(N - \sum_{i=1}^N x_i\right) \log(1-P) \end{aligned}$$

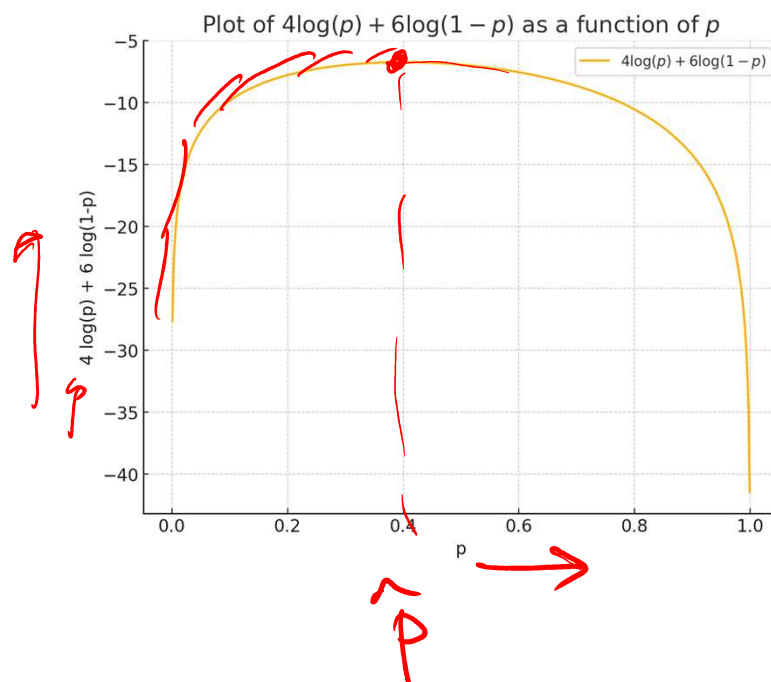
$$\text{Let } - \sum_{i=1}^N x_i = N_1$$

$$\max_p \underbrace{N_1 \log p + (N - N_1) \log(1 - p)}_{LL(p)}$$

$$\frac{\partial LL}{\partial p} = \frac{N_1}{\hat{p}} - \frac{N - N_1}{1 - \hat{p}} = 0$$

$$\Rightarrow \hat{p} = \frac{N_1}{N}$$

concave in p
 \Rightarrow unique maximum at the p where $\frac{\partial LL}{\partial p} = 0$



Examples: MLE for Poisson

$$X \sim \text{exp}(\lambda)$$
$$f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$D = \{x_1, x_2, \dots, x_N\}$$

$$LL(\lambda) = \sum_{i=1}^N \log \left(\frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \right) = \left(\sum_{i=1}^N x_i \right) \log \lambda - \lambda N - \sum_{i=1}^N \log x_i!$$

$$\hat{\lambda} = \arg \max_{\lambda} \left(\sum_{i=1}^N x_i \right) \log \lambda - \lambda N$$

$$\frac{\partial LL}{\partial \lambda} = \sum_{i=1}^N \frac{x_i}{\lambda} - N \quad \therefore \hat{\lambda} =$$

$$\frac{\sum x_i}{N} \text{ sample mean}$$

$$f(k, \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$$

MLE for Gaussian

Home work

$$x \sim N(\mu, \sigma^2)$$
$$f(x, \theta = [\mu, \sigma^2]) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\frac{\partial LL}{\partial \mu} = 0$$

← at μ calculated above

$$\frac{\partial LL}{\partial \sigma} = 0$$

$$\begin{aligned}
 f(x_1, \dots, x_n | \mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[\frac{-(x_i - \mu)^2}{2\sigma^2}\right] \\
 &= \left(\frac{1}{2\pi}\right)^{n/2} \frac{1}{\sigma^n} \exp\left[\frac{-\sum_1^n (x_i - \mu)^2}{2\sigma^2}\right]
 \end{aligned}$$

The logarithm of the likelihood is thus given by

$$\log f(x_1, \dots, x_n | \mu, \sigma) = -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{\sum_1^n (x_i - \mu)^2}{2\sigma^2}$$

In order to find the value of μ and σ maximizing the foregoing, we compute

$$\begin{aligned}
 \frac{\partial}{\partial \mu} \log f(x_1, \dots, x_n | \mu, \sigma) &= \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma^2} \\
 \frac{\partial}{\partial \sigma} \log f(x_1, \dots, x_n | \mu, \sigma) &= -\frac{n}{\sigma} + \frac{\sum_1^n (x_i - \mu)^2}{\sigma^3}
 \end{aligned}$$

Equating these equations to zero yields that

$$\hat{\mu} = \sum_{i=1}^n x_i / n$$

and

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

Example 7.2.d. The number of traffic accidents in Berkeley, California, in 10 randomly chosen nonrainy days in 1998 is as follows:

4, 0, 6, 5, 2, 1, 2, 0, 4, 3

Use these data to estimate the proportion of nonrainy days that had 2 or fewer accidents that year.

• Most difficult question: what distribution to use to model accidents in a city?

- Binomial? Will need to know total number of drivers
- Gaussian?
- Poisson?

Solution in textbook

Homework