

Stat 2911 Lecture Notes

Class II , 2017

Uri Keich

© Uri Keich, The University of  
Sydney

Method of moments estimation:  
definition and examples, MLE:  
definition (Rice 8.1, 8.2-8.4)

## Estimation

Given observations (sample)  $x_1, \dots, x_n$  from a dist.  $F(\theta)$ , where  $\theta$  is an unknown parameter, we are interested in estimating  $\theta$ .

Examples. Bernoulli( $p$ ), Poisson( $\lambda$ ), geometric( $p$ ), multinomial( $m; p_1, \dots, p_r$ ), Binomial( $m, p$ ).

Assumption: the sample  $x_1, \dots, x_n$  is a particular realization of the iid RVs  $X_1, \dots, X_n \sim F(\theta)$

## The Method of Moments

Example. Suppose we are given a sample  $x_1, \dots, x_n$  from a Bernoulli( $p$ ) dist.

How would you estimate  $p$ ?

The method of moments works by noting that

$$\mu_1 = E(\bar{X}_i) = p$$

so rather than estimating  $p$  we estimate  $\mu_1$ !



It's actually not as Homer-ish as it sounds:

$p$  is a parameter specific to the dist. whereas the mean  $\mu_1$  is a universal notion so we can instead concentrate on estimating universal quantities.

For example, how would you estimate  $\mu_1$  given a sample  $x_1, x_2, \dots, x_n$ ?

$\mu_1$  can be estimated by the sample mean  $\hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$

To generalize this we need to define the moments of the dist.

We can extend our definition of  $L^1$  and  $L^2$  to

$$L^k = \{ \bar{X} : \Omega \rightarrow \mathbb{R} \text{ is a RV with } \sum_i |x_i|^k p_{\bar{X}}(x_i) < \infty \} \quad (k \in \mathbb{N})$$

Note:  $\bar{X} \in L^k \Leftrightarrow E|\bar{X}|^k < \infty \Leftrightarrow |\bar{X}|^k \in L^1$

For  $\bar{X} \in L^k$  we define its  $k^{\text{th}}$  moment as

$$\mu_k = E(\bar{X}^k) = \sum_i x_i^k p_{\bar{X}}(x_i) \quad \text{e.g. } \mu_1 = E(\bar{X})$$

Clearly, the moments are property of the dist. of  $\bar{X}$  rather than of  $\bar{X}$  itself: the RHS is defined in terms of  $p_{\bar{X}}$ .

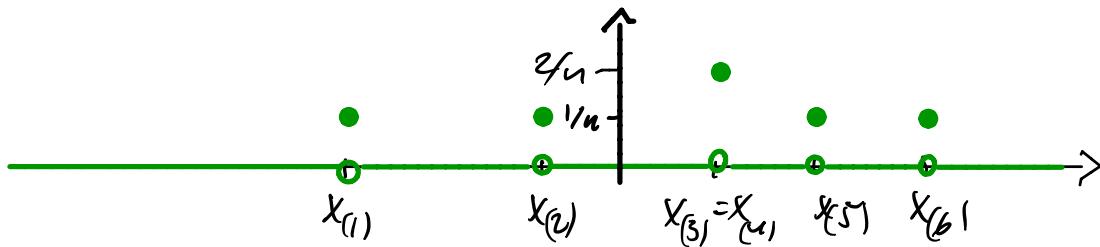
Hence we can just as well refer to the moments of the dist. or the pmf (for discrete RVs).

The moments can be estimated from the sample  $x_1, \dots, x_n$  by the sample moments:

$$\hat{\mu}_k := \frac{1}{n} \sum_i x_i^k$$

Note that the sample moments are the moments of the empirical dist. which is defined by assigning a probability of  $1/n$  to each sample point/observation  $x_i$  (cumulatively).

The point of the empirical dist. associated with the sample  $\{X_i\}_1^n$ :



$x_{(i)}$  are the ordered sample points (order statistics)  
 so  $\{x_{(i)}\}_1^n = \{x_i\}_1^n$  and  $x_{(1)} \leq x_{(2)} \leq x_{(3)} \dots \leq x_{(n)}$

In the method of moments we express the parameter  $\theta$  as a function of the moments, which yields an estimate of the parameter by ...

replacing the moments with their estimates: the sample moments.

### Examples

1)  $X_i \sim \text{Bernoulli}(p)$

$$p = \mu_1 \Rightarrow \hat{p} = \hat{\mu}_1 = \frac{1}{n} \sum_1^n x_i = \bar{x}$$

2)  $X_i \sim \text{Binom}(m, p)$   $m$  is known but  $p$  is not.

$$\mu_1 = E(\bar{X}_i) = mp \Rightarrow p = \mu_1/m$$

$$\Rightarrow \hat{p} = \bar{x} =$$

Ex. Given a sample  $y_1, \dots, y_m$  from a  $\text{Bernoulli}(p)$  dist.  
 verify that  $\hat{p}$  from (1) agrees with  $\hat{p}$  from (2) with  $x_i = \sum_j^m y_j$ .

$$3) X_i \sim \text{Pois}(\lambda) \quad \lambda > 0$$

$$\mu_1 = E(X_i) = \lambda \quad \Rightarrow \quad \hat{\lambda} = \hat{\mu}_1 = \bar{x}$$

Recall the Borfkiwicz horse kicks example.

Note: we could have used

$$\lambda = \sigma^2 = \mu_2 - \mu_1^2 \Rightarrow \hat{\lambda} = \hat{\mu}_2 - \hat{\mu}_1^2$$

Can you guess why that would be an inferior estimator?  
It involves the second as well as the first moment,  
hence its variability would be higher.

$$4) X_i \sim \text{geometric}(p) \quad p \in [0, 1]$$

$$\mu_1 = E(X_i) = 1/p$$

$$\Rightarrow p = 1/\mu_1$$

$$\Rightarrow \hat{p} = 1/\bar{x}$$

5)  $X_1, \dots, X_n \sim \text{Binom}(m, p)$  where both  $m$  and  $p$  are unknown.

Example: you try to estimate your competitor's production rate and/or failure rate by, say, going through their rejects bin.

Goal: express  $m$  and  $p$  in terms of moments of lowest possible order.

$$\mu_1 = E(\bar{X}_c) = mp \quad (\mu = \mu_1)$$

$$\mu_2 - \mu_1^2 = \sigma^2 = mp(1-p)$$

$$\Rightarrow \sigma^2 = \mu(1-p)$$

$$\Rightarrow p = 1 - \sigma^2/\mu = \frac{\mu - \sigma^2}{\mu}$$

$$\Rightarrow \hat{p} = 1 - \hat{\sigma}^2 / \hat{\mu}$$

$$m = \mu/p = \frac{\mu^2}{\mu - \sigma^2}$$

$$\hat{m} = \frac{\hat{\mu}^2}{\hat{\mu} - \hat{\sigma}^2}$$

where  $\hat{\mu} = \hat{\mu}_1 = \bar{X}$

$$\begin{aligned} \hat{\sigma}^2 &= \hat{\mu}_2 - \hat{\mu}_1^2 \\ &= \frac{1}{n} \sum_i x_i^2 - \left( \frac{1}{n} \sum_i x_i \right)^2 \end{aligned}$$

$$\begin{aligned} \text{ex} \quad &= \frac{1}{n} \sum_i^n (x_i - \bar{x})^2 \quad (\text{sample variance, or variance of the empirical dist.}) \end{aligned}$$

Note that the sample variance is often taken as

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i^n (x_i - \bar{x})^2$$

Can you spot a potential problem with these estimators?

More generally, estimating using the method of moments consists of expressing the moments of the distribution in terms of the unknown parameters  $\underline{\theta} = (\theta_1, \dots, \theta_r)$  and then inverting this map to obtain

$$\theta_1 = f_1(\mu_1, \dots, \mu_k)$$

$$\theta_2 = f_2(\mu_1, \dots, \mu_k)$$

 $\vdots$ 

$$\theta_r = f_r(\mu_1, \dots, \mu_k)$$

where  $k \geq r$  is minimal and  $f_i: \mathbb{R}^k \rightarrow \mathbb{R}$ .

$$\Rightarrow \hat{\theta}_i = f_i(\hat{\mu}_1, \dots, \hat{\mu}_n) \text{ where}$$

$$\hat{\mu}_k := \frac{1}{n} \sum_i x_i^k$$

Note: for a particular sample  $x_1, \dots, x_n$

$$\bar{x} = \sum_i x_i \cdot \frac{1}{n} \in \mathbb{R} \quad \text{but}$$

$$\bar{X} = \sum_i X_i \cdot \frac{1}{n} \text{ is a RV}$$

## Maximum Likelihood Estimation (MLE)

Given a sample  $x_1, \dots, x_n$  from a dist. with pmf  $f_\theta(x)$  we want to estimate  $\theta$ .

The **likelihood function** gives the prob. of the observed sample, assuming the sample was generated using  $f_\theta$ , where  $\theta$  is known:

$$L(\theta) = P_\theta(\bar{X}_1 = x_1, \dots, \bar{X}_n = x_n)$$

not the prob. of  $\theta$ !

If the sample is modeled as an iid then

$$L(\theta) = \prod_1^n f_\theta(x_i)$$

Example.  $\bar{X}_i \sim \text{Bernoulli}(i; \theta) \Rightarrow f_\theta(x) = \begin{cases} 1-\theta & x=0 \\ \theta & x=1 \\ 0 & \text{otherwise} \end{cases}$

In order to analyze  $L(\theta)$  it is beneficial to simplify:

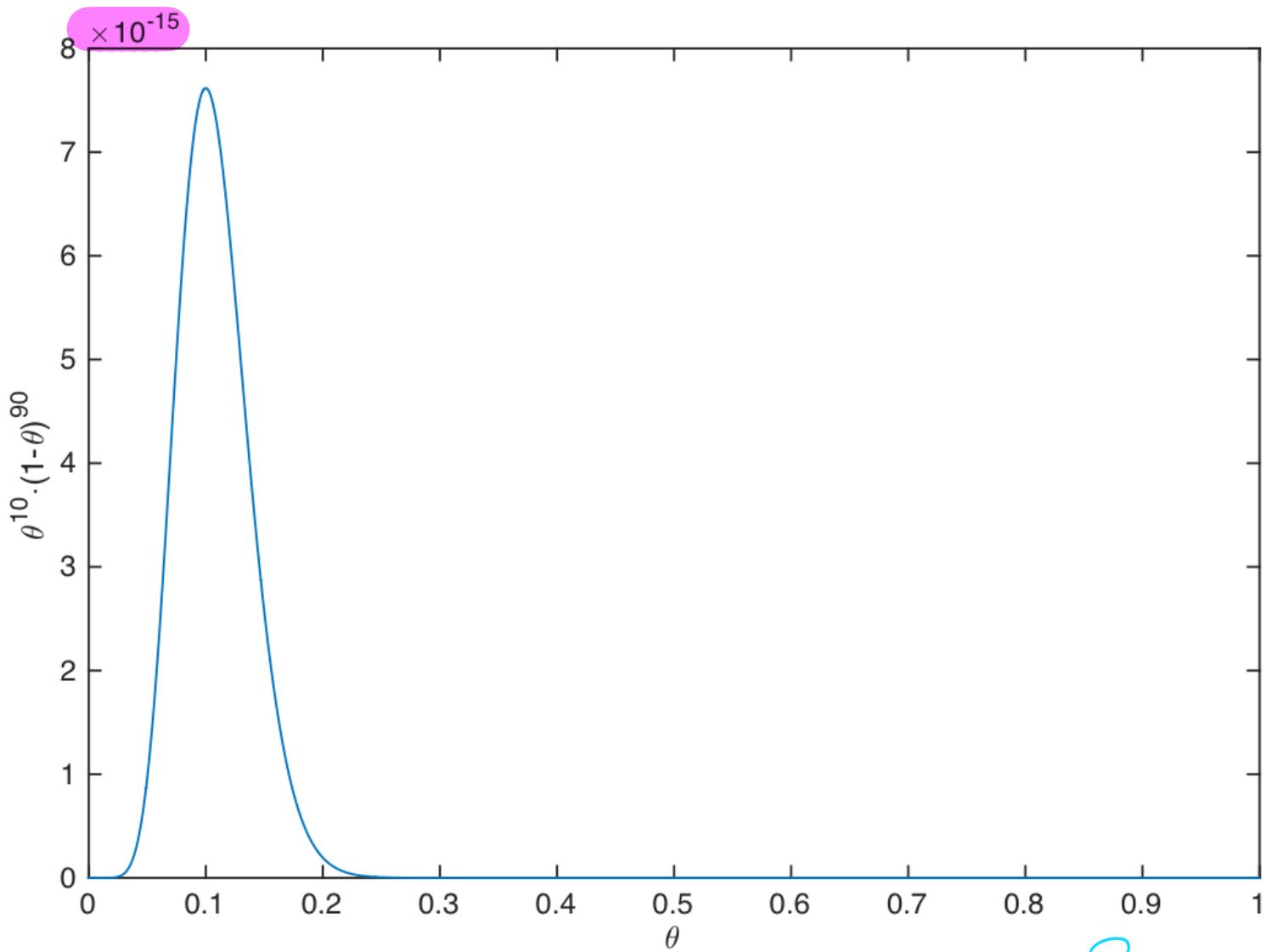
$$f_\theta(x) = \theta^x \cdot (1-\theta)^{1-x} \cdot \mathbb{1}_{x \in \{0,1\}}$$

Bernoulli sample  $\Rightarrow x_i \in \{0,1\}$

$$\begin{aligned} \Rightarrow L(\theta) &= \prod_1^n \theta^{x_i} (1-\theta)^{1-x_i} \\ &= \theta^{\sum x_i} (1-\theta)^{\sum (1-x_i)} \\ &= \theta^{\sum x_i} (1-\theta)^{n - \sum x_i} \end{aligned}$$

e.g.  $n=100, \sum x_i = 10 \Rightarrow L(\theta) = \theta^{10} (1-\theta)^{90}$

$$n=100, \sum_i x_i = 10 \Rightarrow L(\theta) = \theta^{10} (1-\theta)^{90}$$



How would you estimate  $\theta$ ?

The MLE is defined as

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

Since  $\log$  is a monotone increasing function

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

where  $l(\theta) = \log L(\theta) = \sum_{i=1}^n \log f_{\theta}(x_i)$  (more convenient for maximizing!)