

Stat 2911 Lecture Notes

Class 12 , 2017

Uri Keich

© Uri Keich, The University of  
Sydney

MLE: Bernoulli, Poisson,  
Binomial( $p$ ), Binomial( $m, p$ ),  
Multinomial (Rice 8.5)

## 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  $\theta$  is known:

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

If the sample is modeled as an iid then

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

The **MLE** is defined as

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

Since  $\log$  is a monotone increasing function

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

where  $\ell(\theta) = \log L(\theta) \stackrel{\text{in iid case}}{=} \sum_1^n \log f_\theta(x_i)$

Example.  $\bar{X}_i \sim \text{Bernoulli}(i; \theta)$

$$L(\theta) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

Examples 1)  $\bar{X}_i \sim \text{Bernoulli}(\theta) \Rightarrow L(\theta) = \theta^{\sum X_i} (1-\theta)^{n-\sum X_i}$

$$\Rightarrow l(\theta) = s \cdot \log \theta + (n-s) \log (1-\theta), \text{ where } s = \sum X_i$$

How shall we maximize  $l(\theta)$ ?

$$\begin{aligned} l'(\theta) &= \frac{s}{\theta} - \frac{n-s}{1-\theta} & \theta \in (0,1) \\ &= \frac{s-n\theta}{\theta(1-\theta)} \\ &= \frac{n}{\theta(1-\theta)} (s/n - \theta) \end{aligned}$$

$$\text{Sgn}(l'): \quad \begin{array}{c} + \quad \overset{\circ}{\underset{\nearrow \bar{x}=s/n}{\underset{\searrow}{\text{---}}}} \quad - \\ \theta \qquad \qquad \qquad \theta \end{array}$$

$\Rightarrow$  Assuming  $\bar{x} \in (0,1)$ ,  $\theta = s/n = \bar{x}$  is a unique critical point which is a global max for  $\theta \in [0,1]$  (why?), so the MLE is  $\hat{\theta} = \bar{x}$ . (moment estimator?) What if  $\bar{x} \in \{0,1\}$ ?

If  $\bar{x} = 0$ ,  $s = 0$  so  $l(\theta) = n \log(1-\theta)$  which is clearly maximized for  $\theta = 1$ .

$\Rightarrow \hat{\theta} = \bar{x}$  in this case as well.

Finally if  $\bar{x} = 1$ ,  $s = n$  so  $l(\theta) = n \log(\theta)$  which is maximized for  $\hat{\theta} = 1$ .

$\Rightarrow \hat{\theta} = s/n = \bar{x}$  is the Bernoulli( $\theta$ ) MLE.

Note that the above analysis shows that for  $x_i > 0$

and  $s = \sum_i^n x_i$ ,  $\underset{\theta \in [0,1]}{\operatorname{argmax}} \theta^s (1-\theta)^{N-s} = \frac{1}{N}s$ .

2)  $\bar{X}_i \sim \text{Pois}(\lambda)$   $\lambda > 0$

$$\begin{aligned} l(\lambda) &= \sum_i^n \log (e^{-\lambda} \lambda^{x_i} / x_i!) \\ &= -n\lambda + (\sum_i^n x_i) \log \lambda - \sum_i^n \log(x_i!) \end{aligned}$$

$$\begin{aligned} \Rightarrow l'(\lambda) &= -n + \frac{1}{\lambda} \sum_i^n x_i \quad \lambda > 0 \\ &= \frac{n}{\lambda} (\bar{x} - \lambda) \end{aligned}$$

$$\operatorname{sgn}(l'): \quad \begin{array}{c} + \quad 0 \quad - \\ \nearrow \quad \searrow \\ 0 \quad \bar{x} \end{array} \quad \rightarrow \lambda$$

If  $\bar{x} > 0$  then  $\lambda = \bar{x}$  is a unique critical point of  $l$  which is a global maximum (why?)

$\Rightarrow \hat{\lambda} = \bar{x}$  is the MLE (moment estimator?)

If  $\bar{x} = 0$ , then  $\sum_i x_i = 0$  so  $l(\lambda) = -n\lambda$  which is clearly maximized for  $\lambda = 0$ , so again  $\hat{\lambda} = \bar{x}$ .

$\Rightarrow \hat{\lambda} = \bar{x}$  is always the MLE

3) Binomial  $(m, p)$  where  $m$  is known and  $p$  is not.

$$\begin{aligned} L(p) &= \prod_{i=1}^n \binom{m}{x_i} p^{x_i} (1-p)^{m-x_i} \\ &= p^{\sum_{i=1}^n x_i} (1-p)^{m - \sum_{i=1}^n x_i} \prod_{i=1}^n \binom{m}{x_i} \end{aligned}$$

If we denote by  $N = m \cdot n$  (total number of iid Bernoulli trials), then

$$L(p) = C(m; x_1, \dots, x_n) p^{\sum_{i=1}^n x_i} (1-p)^{N - \sum_{i=1}^n x_i}$$

so ignoring  $C()$  this is the same function we maximized in the Bernoulli case (check that all we used there is that  $x_i \geq 0$ ), and therefore the MLE is:

$$\hat{p} = \frac{1}{N} \sum_{i=1}^n x_i = \frac{1}{m \cdot n} \sum_{i=1}^n x_i = \bar{x}$$

(moment estimator?)

4) Binomial  $\binom{m}{k} p^k (1-p)^{m-k}$  both  $m$  and  $p$  are unknown

Moment estimators:

$$\tilde{m} = \bar{x}^2 / (\bar{x} - \hat{\sigma}^2), \quad \tilde{p} = \bar{x} / \tilde{m} = (\bar{x} - \hat{\sigma}^2) / \bar{x} \quad (\tilde{m}\tilde{p} = \bar{x})$$

If  $m$  is known then we saw the MLE is  $\hat{p} = \bar{x}/m$ .

In other words, if  $\ell(p, m) = \sum_{k=1}^n \log \left[ \binom{m}{x_k} p^{x_k} (1-p)^{m-x_k} \right]$ ,

then  $\operatorname{argmax}_p \ell(p, m) = \bar{x}/m =: \hat{p}(m)$ . (fixed sample)

$$\begin{aligned} \Rightarrow \max_{p, m} \ell(p, m) &= \max_m \max_p \ell(p, m) \\ &= \max_m \ell(\hat{p}(m), m) \\ &= \max_m \sum_{k=1}^n \log \left[ \binom{m}{x_k} \left( \frac{\bar{x}}{m} \right)^{x_k} \left( 1 - \frac{\bar{x}}{m} \right)^{m-x_k} \right] \end{aligned}$$

This maximization can be partially achieved numerically:

For each value  $m = \dots, \dots, \dots$  compute  $\ell(\bar{x}/m, m)$  and choose  $m$  that gives the value.

How far should the "..." extend to?

Unfortunately, we cannot prescribe an upper bound on  $m$  (example in next tutorial). In practice we bound  $m$ .

Note that the MLE and moment estimators differ!  
can be negative!

5) A sample  $(x_1, \dots, x_r)$  is drawn from a multinomial  $(n; p_1, \dots, p_r)$  distribution where  $n$  is known but  $\underline{\theta} := (p_1, \dots, p_r)$  is not.

E.g. in  $n=1000$  rolls of a die we observed

$$\underline{x} = (105, 207, 253, 152, 182, 101)$$

How would you estimate  $p_i$ ?

$$\begin{aligned} L(\underline{\theta}) &= P_{\underline{\theta}}(\bar{X}_1 = x_1, \dots, \bar{X}_r = x_r) \\ &= \binom{n}{x_1 \dots x_r} \prod_{i=1}^r \theta_i^{x_i} \end{aligned}$$

$$\Rightarrow \ell(\underline{\theta}) = \underbrace{\log \binom{n}{x_1 \dots x_r}}_{\text{const.}} + \underbrace{\sum_{i=1}^r x_i \log \theta_i}_{g(\underline{\theta})} \quad (x_i \text{ are fixed})$$

$$\hat{\underline{\theta}} = \underset{\underline{\theta} \in \mathbb{R}^r : \theta_i \geq 0, \sum \theta_i = 1}{\operatorname{argmax}} g(\underline{\theta})$$

General method for constrained optimization: Lagrange Multipliers. Alternatively, use ad-hoc optimization.

$$\begin{aligned} \text{Consider } \underline{x} = (5, 0, 0) \Rightarrow g(\underline{\theta}) &= 5 \log \theta_1 \Rightarrow \hat{\underline{\theta}} = (1, 0, 0) \\ \underline{x} = (4, 6, 0) \Rightarrow g(\underline{\theta}) &= 4 \log \theta_1 + 6 \log \theta_2 \Rightarrow \hat{\underline{\theta}} = (0.5, 0.5) \end{aligned}$$

WLOG (without loss of generality) all  $x_i > 0$ : otherwise, if  $x_i = 0$  the maximum is attained at  $\theta_i = 0$  and we can ignore the  $i^{\text{th}}$  coord.

Note that we optimize over  $\underline{\theta} \in \mathbb{R}^r$  with  $\sum_i \theta_i = 1$  (and  $\theta_i \geq 0$ ), therefore for those  $\underline{\theta}$ ,

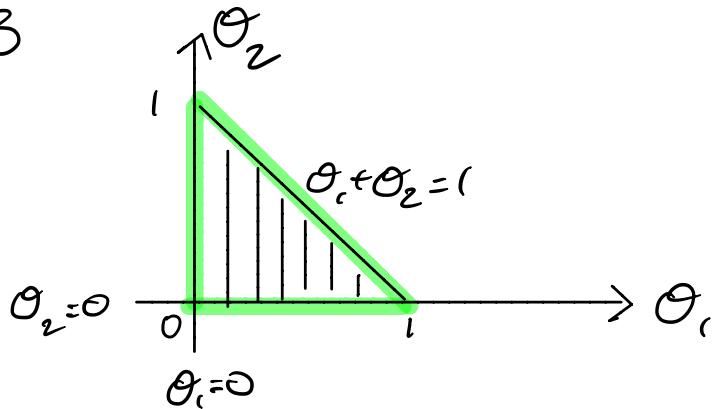
$$\begin{aligned} g(\theta_0, \dots, \theta_r) &= g(\theta_0, \dots, \theta_{r-1}, 1 - \sum_{i=0}^{r-1} \theta_i) \\ &= \sum_{i=0}^{r-1} x_i \log \theta_i + x_r \log (1 - \sum_{i=0}^{r-1} \theta_i) \end{aligned}$$

Therefore, with  $h(\theta_0, \dots, \theta_{r-1}) := g(\theta_0, \dots, \theta_{r-1}, 1 - \sum_{i=0}^{r-1} \theta_i)$

$$\hat{\theta} = \underset{\substack{\underline{\theta} \in \mathbb{R}^{r-1}; \theta_i \geq 0, \sum_i \theta_i \leq 1}}{\operatorname{argmax}} h(\theta_0, \dots, \theta_{r-1})$$

The latter is essentially an unconstrained optimization.

Example.  $r=3$



In particular, the maximum of  $h$  on the simplex  $\{\underline{\theta} \in \mathbb{R}^{r-1} : \theta_i \geq 0, \sum_i \theta_i \leq 1\}$  is attained either at the **boundary** ( $\theta_i = 0$  for  $i = 1, 2, \dots, r-1, \underline{r}$ ), or at a critical point ( $\sum_i \theta_i = 1 \Rightarrow \theta_r = 1$ )

$$\frac{\partial h}{\partial \theta_i} = 0 \quad i = 1, 2, \dots, r-1$$

$$\text{Recall } h(\theta_1, \dots, \theta_{r-1}) = \sum_{i=1}^{r-1} x_i \log \theta_i + x_r \log \left(1 - \sum_{i=1}^{r-1} \theta_i\right)$$

Since  $x_i > 0 \ \forall i$  the maximum cannot be attained at the boundary. (why?)

At a critical point  $\frac{\partial h}{\partial \theta_k} = 0 \text{ for } k=1, 2, \dots, r-1.$

$$\Leftrightarrow \frac{x_k}{\theta_k} - \frac{x_r}{1 - \sum_{i=1}^{r-1} \theta_i} = 0 \text{ for } k=1, 2, \dots, r-1.$$

$$\Leftrightarrow \frac{x_k}{\theta_k} = \frac{x_r}{\theta_r} = C \quad k=1, \dots, r-1.$$

But

$$= \sum_{k=1}^r x_k = \sum_{k=1}^r C \theta_k =$$

$$\Rightarrow \text{At a critical point } \theta_k = \frac{x_k}{C} = \frac{x_k}{n}.$$

$\Rightarrow$  The critical point is unique and therefore has to be a global maximum. (why?)

$\Rightarrow \hat{\theta}_k = \frac{x_k}{n} \quad k=1, \dots, r$  is the MLE

Same for the moment estimate (tutorial).