

Lecture 8

①

• An alternative to MM, and one that is in some way, more principled, is the notion of maximum likelihood estimation (MLE).

• The idea here is to choose parameter $\hat{\theta} \in \Theta \subset \mathbb{R}^D$ that "fit" the data best, in the sense of being the "likeliest" choice of parameter.

• To make this precise, we need to focus on the case that $\mathcal{F} = \{f(x; \theta) \mid \theta \in \Theta\}$ is a family of probability density functions. Let $\{X_i\}_{i=1}^n$ i.i.d. \mathcal{F} be an iid sample from X , unknown. We do not assume the density of X is in \mathcal{F} , though it's possible.

Defn. Let $\{f(x; \theta) \mid \theta \in \Theta\}$ be a parametric family of densities. Let data $\{X_i\}_{i=1}^n \subset \mathbb{R}$ be observed. The likelihood function associated to \mathcal{F} and the data $\{X_i\}_{i=1}^n$ is $L_n: \Theta \rightarrow \mathbb{R}_{\geq 0}$

$$\Theta \mapsto \prod_{i=1}^n f(X_i; \theta).$$

The idea is that even if $f(X_i; \theta) \neq P(X_i \text{ occurs under parameter } \theta)$, it is true that $f(X_i; \theta)$ large suggests X_i is likely under parameter θ .

So, $L_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$ is large if the observed data $\{x_i\}_{i=1}^n$ is likely under θ .

Notice we need a parametric family and data to define a likelihood function.

Remark: While for all θ , $\int_{\mathbb{R}} f(x; \theta) dx = 1$, it need not be the case (and generically is not) that $\int_{\Theta} f(x; \theta) d\theta = 1$ for any x .

The data variable (x) and the parameter variable (θ) behave very differently.

Given the likelihood function L_n , the idea of MLE is simple: make

L_n maximal:

$$\hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmax}} L_n(\theta), \text{ i.e. } L_n(\hat{\theta}_n) = \max_{\theta \in \Theta} L_n(\theta).$$

~~the maximum~~

So, given data and a model, all we have to do is do calculus on L_n ! Easier said than done.

Because taking $\frac{\partial}{\partial \theta}$ over a product is tedious, it is preferable often to consider not $L_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$, but $\log(L_n(\theta)) = \sum_{i=1}^n \log(f(x_i; \theta))$.

This log-likelihood is often easier to handle. Since $\log(y)$ is increasing, the optima don't change.

ex: Let $\mathcal{F} = \{f_\lambda \mid \lambda > 0\}$ be the collection of exponential densities,

$$f_\lambda(x) = \begin{cases} 0, & x < 0 \\ \lambda e^{-\lambda x}, & x \geq 0 \end{cases}$$

→ assumed $x_i > 0 \forall i$

Then given data $\{x_i\}_{i=1}^n \sim \underline{X}$, we can compute the likelihood and

log-likelihood: $L_n(\lambda) = \prod_{i=1}^n f(x_i; \lambda)$

$$= \prod_{i=1}^n \lambda e^{-\lambda x_i}$$

$$= \lambda^n \exp\left(-\lambda \sum_{i=1}^n x_i\right)$$

$$\log(L_n(\lambda)) = \log\left(\lambda^n \exp\left(-\lambda \sum_{i=1}^n x_i\right)\right)$$

$$= n \log(\lambda) - \lambda \sum_{i=1}^n X_i$$

it's typically to see some kind of "normalizing constant" appearing

To find the MLE estimator $\hat{\lambda}$, we can solve

$$\frac{d}{d\lambda} \left[\log(L_n(\lambda)) \right] = 0, \text{ i.e. } \frac{d}{d\lambda} \left[n \log(\lambda) - \lambda \sum_{i=1}^n X_i \right] = 0$$

$$\Leftrightarrow \frac{n}{\lambda} - \sum_{i=1}^n X_i = 0$$

$$\Leftrightarrow \hat{\lambda} = \frac{n}{\sum_{i=1}^n X_i}, \text{ same as MM.}$$

In the previous example, the log-likelihood was differentiable, so we were able to use calculus. When it is not differentiable, we need to think harder!

ex: Let $\mathcal{F} = \{f_\theta \mid \theta \in \mathbb{R}_+\}$, where $f_\theta = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta] \\ 0, & \text{else} \end{cases}$

be the family of uniform densities on $[0, \theta]$. Then given data $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} X$,

assumed to be positive, what is $\hat{\theta}_{MLE}$? Writing out the likelihood gives

$$\begin{aligned}
 \mathcal{L}_n(\theta) &= \prod_{i=1}^n f(x_i; \theta) \\
 &= \prod_{i=1}^n \frac{1}{\theta} \cdot \begin{cases} 1, & x_i \leq \theta \\ 0, & x_i > \theta \end{cases} \\
 &= \frac{1}{\theta^n} \cdot \prod_{i=1}^n \begin{cases} 1, & x_i \leq \theta \\ 0, & x_i > \theta \end{cases}
 \end{aligned}$$

Notice that $\prod_{i=1}^n \begin{cases} 1, & x_i \leq \theta \\ 0, & x_i > \theta \end{cases} = \begin{cases} 1, & x_i \leq \theta \text{ for all } i \\ 0, & \text{else.} \end{cases}$

$$\text{So, } \mathcal{L}_n(\theta) = \underbrace{\frac{1}{\theta^n}}_{\textcircled{I}} \cdot \underbrace{\begin{cases} 1, & \max_i x_i \leq \theta \\ 0, & \text{else} \end{cases}}_{\textcircled{II}}$$

Notice the tension between \textcircled{I} and \textcircled{II} : to keep \textcircled{I} large, we want θ small. But \textcircled{II} switches from 0 to 1 only once $\theta \geq \max_i x_i$. So,

$$\hat{\theta}_{MLE} = \max_i x_i$$