

## 2. Introduction to M-estimation

**Video Note:** The video below is the last 10 minutes of the previous lecture when Prof Rigollet began the discussion of m-estimation.

### Introduction to M-estimation



□ (Caption will be displayed when you start playing the video.)

we said there's this function that my true parameter maximizes.

This function is the expectation of something.

So let me replace the expectation by an average

and optimize that function instead.

M-estimation is really this.

M means maximum or minimum, so it's something which is about maximizing a function.

And it's basically saying, well, this thing that you

did for the KL divergence, you could do in general, right?

The KL divergence was just a function

that was the expectation of something

so that you knew that the maximizer was theta star.

#### 视频

[下载视频文件](#)

#### 字幕

[下载 SubRip \(.srt\) file](#)

[下载 Text \(.txt\) file](#)

### M-estimation

Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be i.i.d. with some unknown distribution  $\mathbf{P}$  and an associated parameter  $\mu^*$  on a sample space  $\mathcal{E}$ . We make no modeling assumption that  $\mathbf{P}$  is from any particular family of distributions.

An **M-estimator**  $\hat{\mu}$  of the parameter  $\mu^*$  is the **argmin of an estimator of a function  $Q(\mu)$  of the parameter** which satisfies the following:

- $Q(\mu) = \mathbb{E}[\rho(\mathbf{X}, \mu)]$  for some function  $\rho : \mathcal{E} \times \mathcal{M} \rightarrow \mathbb{R}$ , where  $\mathcal{M}$  is the set of all possible values of the unknown true parameter  $\mu^*$ ;
- $Q(\mu)$  attains a **unique** minimum at  $\mu = \mu^*$ , in  $\mathcal{M}$ . That is,  $\text{argmin}_{\mu \in \mathcal{M}} Q(\mu) = \mu^*$ .

In general, the goal is to find the **loss function**  $\rho$  such  $Q(\mu) = \mathbb{E}[\rho(\mathbf{X}, \mu)]$  has the properties stated above.

Note that the function  $\rho(\mathbf{X}, \mu)$  is in particular a function of the random variable  $\mathbf{X}$ , and the expectation in  $\mathbb{E}[\rho(\mathbf{X}, \mu)]$  is to be taken against the **true distribution**  $\mathbf{P}$  of  $\mathbf{X}$ , with associated parameter value  $\mu^*$ .

Because  $Q(\mu)$  is an expectation, we can construct a (consistent) estimator of  $Q(\mu)$  by replacing the expectation in its definition by the sample mean.

#### Example: multivariate mean as minimizer

Let  $\mathbf{X} = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}$  be a continuous random vector with density  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ . Recall the mean of  $\mathbf{X}$  is

$$\mathbb{E}[\mathbf{X}] = \begin{pmatrix} \mathbb{E}[X^{(1)}] \\ \mathbb{E}[X^{(2)}] \end{pmatrix}$$

Recall the Euclidean square of the norm function on  $\mathbb{R}^2$ :

$$\|\cdot\|^2 : \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \mapsto (y_1)^2 + (y_2)^2.$$

We now show that the (multivariate) mean of  $\mathbf{X}$  satisfies:

$$\mathbb{E}[\mathbf{X}] = \operatorname{argmin}_{\vec{\mu} \in \mathbb{R}^2} \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right].$$

(We will use subscripts to label the components of vectors below.)

First, expand  $\mathcal{Q}(\vec{\mu}) = \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right]$  as an integral expression, and write down both partial derivatives  $\frac{\partial \mathcal{Q}}{\partial \mu_1}(\vec{\mu})$  and  $\frac{\partial \mathcal{Q}}{\partial \mu_2}(\vec{\mu})$ :

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left( (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 \right) f(x_1, x_2) \, dx_1 dx_2 \\ \implies \frac{\partial}{\partial \mu_1} \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right] &= -2 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_1 - \mu_1) f(x_1, x_2) \, dx_1 dx_2 \\ \frac{\partial}{\partial \mu_2} \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right] &= -2 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_2 - \mu_2) f(x_1, x_2) \, dx_1 dx_2. \end{aligned}$$

To find the argmin of  $\mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right]$ , we set both partial derivatives to  $\mathbf{0}$ , and obtain:

$$\operatorname{argmin}_{\vec{\mu} \in \mathbb{R}^2} \mathbb{E} \left[ \|\mathbf{X} - \vec{\mu}\|^2 \right] = \begin{pmatrix} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1, x_2) \, dx_1 dx_2 \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_2 f(x_1, x_2) \, dx_1 dx_2 \end{pmatrix} = \begin{pmatrix} \mathbb{E}[X^{(1)}] \\ \mathbb{E}[X^{(2)}] \end{pmatrix}.$$

### Concept check: M-estimators

1/1 point (graded)

Which of the following is true about M-estimation?

(Choose all that apply. Refer to the slides.)

- ☒ M-estimation involves estimating some parameter of interest related to the underlying, unknown distribution (e.g. its mean, variance, or quantiles) ☐
- ☒ Maximum likelihood estimation is a special case of M-estimation. ☐
- ☒ Unlike maximum likelihood estimation and the method of moments, no statistical model needs to be assumed to perform M-estimation. ☐
- ☐ M-estimation cannot be used for parametric statistical models.

☐

#### Solution:

We examine the choices in order.

- "M-estimation involves estimating some parameter of interested related to the underlying, unknown distribution (e.g. its mean, variance, or quantiles)" is correct. This is precisely the goal of M-estimation, as stated in the slides. It is a flexible approach that applies even outside of parametric statistical models.
- "Maximum likelihood estimation is a special case of M-estimation." is correct. If we set the loss function to be the negative log-likelihood, then the same optimization problem defining the MLE is the one considered for the M-estimator associated to this loss function.
- "Unlike maximum likelihood estimation and the method of moments, no statistical model needs to be assumed to perform M-estimation." is correct. As stated above, M-estimation is a flexible approach that can used to approximate relevant quantities of interest to a distribution, such as its moments.
- "M-estimation cannot be used for parametric statistical models." is incorrect. M-estimation can be used in both a parametric and non-parametric context, though in this lecture, we will only see it applied in parametric examples.

提交

你已经尝试了1次（总共可以尝试2次）

Relating M-estimation and Maximum Likelihood Estimation

1/1 point (graded)

Let  $(E, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$  denote a discrete statistical model and let  $X_1, \dots, X_n \stackrel{iid}{\sim} \mathbf{P}_{\theta^*}$  denote the associated statistical experiment, where  $\theta^*$  is the true, unknown parameter. Suppose that  $\mathbf{P}_\theta$  has a probability mass function given by  $p_\theta$ . Let  $\hat{\theta}_n^{\text{MLE}}$  denote the maximum likelihood estimator for  $\theta^*$ .

The maximum likelihood estimator can be expressed as an M-estimator– that is,

$$\hat{\theta}_n^{\text{MLE}} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \rho(X_i, \theta)$$

for some function  $\rho$ .

Which of the following represents the correct choice of the function  $\rho$  so that the equation above is satisfied?

☒  $-\ln p_\theta(X_i)$  ☐

☐  $\ln p_\theta(X_i)$

☐  $p_\theta(X_i)$

☐ None of the above.

Solution:

The correct response is " $-\ln p_\theta(X_i)$ ". Recall that the MLE is defined by

$$\hat{\theta}_n^{\text{MLE}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ln p_\theta(X_i).$$

By symmetry, we also have,

$$\hat{\theta}_n^{\text{MLE}} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n -\ln p_\theta(X_i).$$

Indeed, setting  $\rho(x, \theta) = -\ln p_\theta(x)$ , we recover the maximum likelihood estimator.

提交

你已经尝试了1次（总共可以尝试2次）

☐ Answers are displayed within the problem

Median as a Minimizer

2/3 points (graded)

Assume that  $X$  is a continuous random variable with density  $f : \mathbb{R} \rightarrow \mathbb{R}$ . Then a **median** of  $X$  is defined to be any point  $\mathbf{med}(X) \in \mathbb{R}$  such that

$$P(X > \mathbf{med}(X)) = P(X < \mathbf{med}(X)) = \frac{1}{2}.$$

(Recall that for a continuous distribution,  $P(X > \mathbf{med}(X)) = P(X \geq \mathbf{med}(X))$ .) Note: A median of a distribution is *not necessarily* unique.)

In this problem, you will show that any median satisfies

$$\mathbf{med}(X) = \operatorname{argmin}_{\mu \in \mathbb{R}} \mathbb{E}[|X - \mu|].$$

Which of the following correctly expresses  $\mathbb{E}[|X - \mu|]$  in terms of the density  $f(x)$ ?

☐  $\int_{-\infty}^{\infty} x f(x) dx - \mu$

☐  $\int_{-\infty}^{\infty} x f(x) dx - \mu \left( - \int_{\mu}^{\infty} f(x) dx + \int_{-\infty}^{\mu} f(x) dx \right)$

☐  $\int_{\mu}^{\infty} x f(x) dx - \int_{-\infty}^{\mu} x f(x) dx - \mu$

☒  $\int_{\mu}^{\infty} x f(x) dx - \int_{-\infty}^{\mu} x f(x) dx - \mu \left( \int_{\mu}^{\infty} f(x) dx - \int_{-\infty}^{\mu} f(x) dx \right) \square$

Let  $\mathcal{Q}(\mu) = \mathbb{E}[|X - \mu|]$  denote the expression obtained in the previous question. Then  $\mathcal{Q}(\mu)$  consists of a sum of terms, each of which can be differentiated with respect to  $\mu$ .

What is  $\mathcal{Q}'(\mu) = \frac{d}{d\mu} \mathcal{Q}(\mu)$ ?

*Hint:* Use the product rule and the fundamental theorem of calculus.

☐ 1

☐  $\int_{-\infty}^{\mu} f(x) dx - \int_{\mu}^{\infty} f(x) dx \square$

☒  $4\mu f(\mu) + \int_{-\infty}^{\mu} f(x) dx - \int_{\mu}^{\infty} f(x) dx \square$

☐  $4\mu f(\mu) + 1$

Using your response from the previous question and the definition of median, what is  $\mathcal{Q}'(\text{med}(X))$ ?

☒ 0  $\square$

☐ 1

☐  $4\text{med}(X) f(\text{med}(X)) + 1$

☐ Cannot be determined.

**Solution:**

For the first question, we have

$$\begin{aligned} \mathbb{E}[|X - \mu|] &= \int_{-\infty}^{\infty} |x - \mu| f(x) dx \\ &= \int_{\mu}^{\infty} (x - \mu) f(x) dx + \int_{-\infty}^{\mu} (-x + \mu) f(x) dx \\ &= \int_{\mu}^{\infty} x f(x) dx - \int_{-\infty}^{\mu} x f(x) dx - \mu \left( \int_{\mu}^{\infty} f(x) dx - \int_{-\infty}^{\mu} f(x) dx \right) \end{aligned}$$

Therefore, " $\int_{\mu}^{\infty} x f(x) dx - \int_{-\infty}^{\mu} x f(x) dx - \mu \left( \int_{\mu}^{\infty} f(x) dx - \int_{-\infty}^{\mu} f(x) dx \right)$ " is the correct answer to the first question.

For the second question, we differentiate the previous answer term by term. We have, by the fundamental theorem of calculus and the product rule that

$$\frac{d}{d\mu} \left( \int_{\mu}^{\infty} x f(x) dx \right) = -\mu f(\mu)$$

$$\frac{d}{d\mu}\left(-\int_{-\infty}^{\mu} x f(x) \, dx\right) = -\mu f(\mu)$$

$$\frac{d}{d\mu}\left(-\mu\left(\int_{\mu}^{\infty} f(x) \, dx - \int_{-\infty}^{\mu} f(x) \, dx\right)\right) = -\int_{\mu}^{\infty} f(x) \, dx + \int_{-\infty}^{\mu} f(x) \, dx + 2\mu f(\mu).$$

Adding these terms, we have cancellations, yielding

$$\frac{d}{d\mu}\mathcal{Q}(\mu) = -\int_{\mu}^{\infty} f(x) \, dx + \int_{-\infty}^{\mu} f(x) \, dx.$$

Therefore, the correct response to the second question is " $\int_{\mu}^{\infty} f(x) \, dx - \int_{-\infty}^{\mu} f(x) \, dx$ ".

For the third question, by definition, the median  $\mathbf{med}(X)$  of  $X$  is a real number that satisfies  $P(X > \mathbf{med}(X)) = P(X < \mathbf{med}(X))$ . Therefore,

$$\mathcal{Q}'(\mathbf{med}(X)) = \int_{\mathbf{med}(X)}^{\infty} f(x) \, dx - \int_{-\infty}^{\mathbf{med}(X)} f(x) \, dx = P(X > \mathbf{med}(X)) - P(X < \mathbf{med}(X)) = 0.$$

The correct response is "0".

提交

你已经尝试了2次（总共可以尝试2次）

☐ Answers are displayed within the problem

### Quantile as a Minimizer

6/7 points (graded)

Recall from the lecture that the **check function** is defined as

$$C_{\alpha}(x) \;=\; \begin{cases} -(1-\alpha)x & \text{if } x < 0 \\ \alpha x & \text{if } x \geq 0. \end{cases}$$

Assume that  $\boldsymbol{X}$  is a continuous random variable with density  $\boldsymbol{f} : \mathbb{R} \rightarrow \mathbb{R}$ . Define the  **$\alpha$ -quantile** of  $\boldsymbol{X}$  to be  $\boldsymbol{Q_X}(\alpha) \in \mathbb{R}$  such that

$$\mathbf{P}\left(\boldsymbol{X} \leq \boldsymbol{Q_X}(\alpha)\right) = \alpha.$$

(Here we have used a different convention of the definition of the quantile function from before, where for a standard normal distribution,  $\boldsymbol{q_{\alpha}}$  is such that  $\boldsymbol{P(X > q_{\alpha}) = \alpha.}$ )

Just like for the median, whether  $\boldsymbol{Q_{\alpha}}$  is unique depends on the distribution.

In this problem, you will convince yourself that any  $\alpha$ -quantile of  $\boldsymbol{X}$  satisfies

$$\boldsymbol{Q_{\alpha}}(\boldsymbol{X}) = \mathbf{argmin}_{\mu \in \mathbb{R}} \mathbb{E}\left[\boldsymbol{C_{\alpha}}(\boldsymbol{X} - \mu)\right].$$

First, compute  $\mathbb{E}\left[\boldsymbol{C_{\alpha}}(\boldsymbol{X} - \mu)\right]$ . Answer by entering the coefficients  $\boldsymbol{A}$ ,  $\boldsymbol{B}$ ,  $\boldsymbol{C}$ ,  $\boldsymbol{D}$  in terms of  $\alpha$  and  $\mu$  in the expression below:

$$\begin{aligned} \mathbb{E}\left[\boldsymbol{C_{\alpha}}(\boldsymbol{X} - \mu)\right] \;=\; &\boldsymbol{A} \int_{-\infty}^{\mu} x f(x) \, dx + \boldsymbol{B} \int_{\mu}^{\infty} x f(x) \, dx \\ &+ \boldsymbol{C} \int_{-\infty}^{\mu} f(x) \, dx + \boldsymbol{D} \int_{\mu}^{\infty} f(x) \, dx. \end{aligned}$$

A =

alpha-1

□ Answer: alpha-1

α−1

B =

alpha

□ Answer: alpha

α

C =

-(alpha-1)\*mu

□ Answer: -(alpha-1)\*mu

−(α−1)⋅μ

D =

-alpha\*mu

□ Answer: -alpha\*mu

−α⋅μ

Second, let  $F(\mu) = \mathbb{E}[C_\alpha(X - \mu)]$  denote the expression obtained in the question above. Find  $F'(\mu)$ . Answer by entering the coefficients  $E, G, H$ , in terms of  $\alpha$  and  $\mu$  below:

$$F'(\mu) = (\mathbb{E}[C_\alpha(X - \mu)])' = E + G(\mu f(\mu)) + H \int_{-\infty}^{\mu} f(x) \, dx.$$

$E =$

-alpha

Answer: -alpha

$-\alpha$

$G =$

0

Answer: 0

0

$H =$

-alpha+1

Answer: 1

$-\alpha + 1$

Finally, set  $F'(\mu) = 0$  to find  $\operatorname{argmin}_{\mu \in \mathbb{R}} F(\mu) = \operatorname{argmin}_{\mu \in \mathbb{R}} \mathbb{E}[C_\alpha(X - \mu)]$ . (There is no answer box for this question.)

STANDARD NOTATION

Solution:

Given the check function center about  $\mu$ :

$$C_\alpha(x - \mu) = \begin{cases} -(1 - \alpha)(x - \mu) & \text{if } x < \mu \\ \alpha(x - \mu) & \text{if } x \geq \mu, \end{cases}$$

compute  $F(\mu) = \mathbb{E}[C_\alpha(X - \mu)]$ :

$$\begin{aligned} F(\mu) = \mathbb{E}[C_\alpha(X - \mu)] &= - \int_{-\infty}^{\mu} (1 - \alpha)(x - \mu) f(x) \, dx + \int_{\mu}^{\infty} \alpha(x - \mu) f(x) \, dx \\ &= -(1 - \alpha) \int_{-\infty}^{\mu} x f(x) \, dx + \alpha \int_{\mu}^{\infty} x f(x) \, dx \\ &\quad + (1 - \alpha) \mu \int_{-\infty}^{\mu} f(x) \, dx - \alpha \mu \int_{\mu}^{\infty} f(x) \, dx. \end{aligned}$$

Then, the derivative of  $F$  with respect to  $\mu$  is:

$$\begin{aligned} F'(\mu) &= \frac{d}{d\mu} F(\mu) = -(1 - \alpha) \frac{d}{d\mu} \int_{-\infty}^{\mu} x f(x) \, dx + \alpha \frac{d}{d\mu} \int_{\mu}^{\infty} x f(x) \, dx \\ &\quad + (1 - \alpha) \frac{d}{d\mu} \left( \mu \int_{-\infty}^{\mu} f(x) \, dx \right) - \alpha \frac{d}{d\mu} \left( \mu \int_{\mu}^{\infty} f(x) \, dx \right) \\ &= -(1 - \alpha) (\mu f(\mu)) + \alpha (-\mu) f(\mu) \\ &\quad + (1 - \alpha) \left( \int_{-\infty}^{\mu} f(x) \, dx + \mu f(\mu) \right) - \alpha \left( \int_{\mu}^{\infty} f(x) \, dx - \mu f(\mu) \right) \quad (\text{by fundamental theorem of calculus 2}) \\ &= (1 - \alpha) \int_{-\infty}^{\mu} f(x) \, dx - \alpha \int_{\mu}^{\infty} f(x) \, dx \\ &= (1 - \alpha) \left( \int_{-\infty}^{\mu} f(x) \, dx \right) - \alpha \left( 1 - \int_{-\infty}^{\mu} f(x) \, dx \right) \\ &= \left( \int_{-\infty}^{\mu} f(x) \, dx \right) - \alpha. \end{aligned}$$

median split cdf to two half  
但是这里不需要center

Setting  $F'(\mu) = 0$  yields

$$\int_{-\infty}^{\mu} f(x) \, dx = \alpha.$$

Hence,  $\operatorname{argmin}_{\mu \in \mathbb{R}} F(\mu)$  is an  $\alpha$ -quantile of  $X$ .

提交

你已经尝试了3次（总共可以尝试3次）

(Optional) Convexity of the Expectation of the Loss Function

Strict convexity of  $\mathcal{Q}(\mu) = \mathbb{E}[\rho(X, \mu)]$  ensures that it has a unique minimum, and this is guaranteed by strict convexity of  $\rho(X, \mu)$  in  $\mu$ . We will explain the univariate case below.

Expectation of a convex function is convex

Let  $X$  be a random variable with some unknown distribution  $\mathbf{P}$  with some associated parameter  $\mu^*$  on some sample space  $E$ . Let  $\rho: E \times \mathcal{M} \rightarrow \mathbb{R}$ , where  $\mathcal{M}$  is the set of all possible values of the unknown true parameter  $\mu^*$  and let  $\mathcal{Q}(\mu) = \mathbb{E}[\rho(X, \mu)]$ .

$$\rho(X, \mu) \text{ strictly convex in } \mu \implies \mathbb{E}[\rho(X, \mu)] \text{ strictly convex in } \mu.$$

Proof:

Recall that  $\rho$  being strictly convex in  $\mu$  means that

$$t\rho(x, \mu_1) + (1 - t)\rho(x, \mu_2) - \rho(x, t\mu_1 + (1 - t)\mu_2) > 0 \text{ for all } x.$$

Taking the expectation of the above inequality gives:

$$\mathbb{E}[t\rho(X, \mu_1) + (1 - t)\rho(X, \mu_2) - \rho(X, t\mu_1 + (1 - t)\mu_2)] = t\mathbb{E}[\rho(X, \mu_1)] + (1 - t)\mathbb{E}[\rho(X, \mu_2)] - \mathbb{E}[\rho(X, t\mu_1 + (1 - t)\mu_2)] > 0$$

for any  $\mu_1 \neq \mu_2 \in \mathcal{M}$ , and  $t \in (0, 1)$ . This is because  $\mathbb{E}[f(X)] > 0$  if  $f(X) > 0$  (Think about why: in the discrete case, this can roughly be read as “the weighted average of a collection of positive numbers is positive”. ) The above inequality exactly implies strict convexity of  $\mathbb{E}[\rho(X, \mu)]$ .

Hide

讨论

显示讨论

主题： Unit 3 Methods of Estimation:Lecture 12: M-Estimation / 2. Introduction to M-estimation