

4. Consistency of Maximum Likelihood Estimator

Review: Definition of MLE

0/1 point (graded)

Let $\{E, (\mathbf{P}_\theta)_{\theta \in \Theta}\}$ be a statistical model associated with a sample of i.i.d. random variables X_1, X_2, \dots, X_n . Assume that there exists $\theta^* \in \Theta$ such that $X_i \sim \mathbf{P}_{\theta^*}$.

Recall the **Kullback-Leibler (KL) divergence** between two distributions \mathbf{P}_{θ^*} and \mathbf{P}_θ , with pdfs p_{θ^*} and p_θ respectively, is defined as

$$\text{KL}(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta) = \mathbb{E}_{\theta^*} \left[\ln \left(\frac{p_{\theta^*}(X)}{p_\theta(X)} \right) \right],$$

and a consistent estimator of $\text{KL}(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta)$ is

$$\widehat{\text{KL}}_n(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta) = \text{a constant} - \frac{1}{n} \sum_{i=1}^n \ln p_\theta(X_i).$$

Which of the following represents the maximum likelihood estimator of θ^* ? (Choose all that apply).

☒ $\text{argmin}_{\theta \in \Theta} \widehat{\text{KL}}_n(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta)$ □

☒ $\text{argmax}_{\theta \in \Theta} \sum_{i=1}^n \ln p_\theta(X_i)$ □

☒ $\text{argmax}_{\theta \in \Theta} \ln \left(\prod_{i=1}^n p_\theta(X_i) \right)$ □

☐ $\text{argmax}_{\theta \in \Theta} \ln(L_n(X_1, X_2, \dots, X_n; \theta))$ □

□

Solution:

Recall the **maximum likelihood estimator** can be defined as the

$$\hat{\theta}_n^{MLE} = \text{argmin}_{\theta \in \Theta} \widehat{\text{KL}}_n(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta).$$

In other words, the maximum likelihood estimator is the (unique) θ that minimizes $\widehat{\text{KL}}(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta)$ over the parameter space $\theta \in \Theta$. (The minimizer of the KL divergence is unique due to it being strictly convex in the space of distributions once \mathbf{P}_{θ^*} is fixed.)

All choices are equivalent to this definition:

$$\begin{aligned} \hat{\theta}_n^{MLE} &= \text{argmin}_{\theta \in \Theta} \widehat{\text{KL}}_n(\mathbf{P}_{\theta^*}, \mathbf{P}_\theta) = \text{argmin}_{\theta \in \Theta} \left(\text{Constant} - \frac{1}{n} \sum_{i=1}^n \ln p_\theta(X_i) \right) \\ &= \text{argmax}_{\theta \in \Theta} \left(\frac{1}{n} \sum_{i=1}^n \ln p_\theta(X_i) \right) \quad (\text{drop additive constant and negative sign}) \\ &= \text{argmax}_{\theta \in \Theta} \left(\sum_{i=1}^n \ln p_\theta(X_i) \right) \quad (\text{drop positive scaling factor}) \end{aligned}$$

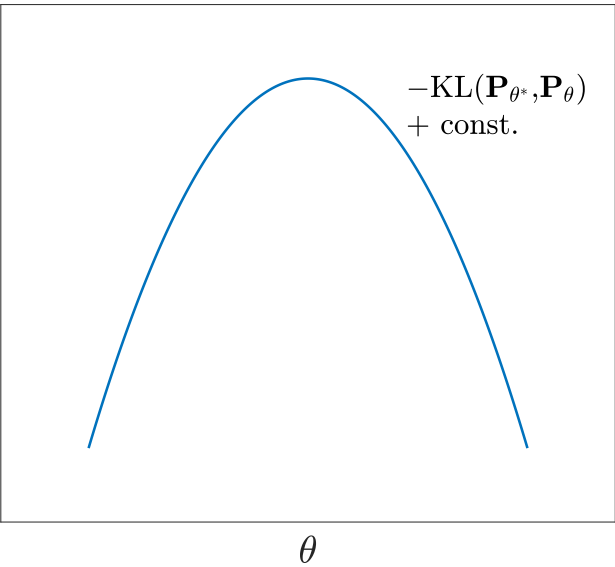
$$= \operatorname{argmax}_{\theta \in \Theta} \left(\ln \left(\prod_{i=1}^n p_{\theta}(X_i) \right) \right) \quad (\text{log property})$$
$$= \operatorname{argmax}_{\theta \in \Theta} \ln(L_n(X_1, X_2, \dots, X_n; \theta)) \quad (\text{definition of likelihood}).$$

提交

你已经尝试了3次（总共可以尝试3次）

☐ Answers are displayed within the problem

Note: In the following video, at around the 3:20 mark, the plot of $-\text{KL}(\mathbf{P}_{\theta^*}, \mathbf{P}_{\theta})$, with θ^* fixed and as a function of θ , is presented incorrectly as a convex curve while it should be concave. This error propagates until the end of the video and we request you to keep the following picture in mind instead:



Consistency of the Maximum Likelihood Estimator

[Start of transcript. Skip to the end.](#)



OK, so now I have this MLEs.
And I have two ways of computing MLEs, either setting derivative equal to 0 or just looking at the plot.
And once I do this, which is really just one way, which
is taking the maximum of the likelihood, once I have this,
I would like you--
you might question whether this estimator is

视频
[下载视频文件](#)

字幕
[下载 SubRip \(.srt\) file](#)
[下载 Text \(.txt\) file](#)

Consistency of MLE

Given i.i.d samples $X_1, \dots, X_n \sim \mathbf{P}_{\theta^*}$ and an associated statistical model $(E, \{\mathbf{P}_{\theta}\}_{\theta \in \Theta})$, the maximum likelihood estimator $\hat{\theta}_n^{\text{MLE}}$ of θ^* is a **consistent** estimator under mild regularity conditions (e.g. continuity in θ of the pdf p_{θ} almost everywhere), i.e.

$$\hat{\theta}_n^{\text{MLE}} \xrightarrow[p]{n \rightarrow \infty} \theta^*.$$

Note that this is true even if the parameter θ is a vector in a higher dimensional parameter space Θ , and $\hat{\theta}_n^{\text{MLE}}$ is a multivariate random variable, e.g. if $\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \in \mathbb{R}^2$ for a Gaussian statistical model.

Multivariate Random Variables

A **multivariate random variable**, or a **random vector**, is a vector-valued function whose components are (scalar) random variables on the same underlying probability space. More specifically, a random vector $\mathbf{X} = (X^{(1)}, \dots, X^{(d)})^T$ of dimension $d \times 1$ is a vector-valued function from a probability space Ω to \mathbb{R}^d :

$$\begin{aligned} \mathbf{X} : \Omega &\rightarrow \mathbb{R}^d \\ \omega &\mapsto \begin{pmatrix} X^{(1)}(\omega) \\ X^{(2)}(\omega) \\ \vdots \\ X^{(d)}(\omega) \end{pmatrix} \end{aligned}$$

where each $X^{(k)}$ is a (scalar) random variable on Ω . We will often (but not always) use the bracketed superscript (k) to denote the k -th component of a random vector, especially when the subscript is already used to index the samples.

The **probability distribution** of a random vector \mathbf{X} is the **joint distribution** of its components $X^{(1)}, \dots, X^{(d)}$.

The **cumulative distribution function (cdf)** of a random vector \mathbf{X} is defined as

$$\begin{aligned} F : \mathbb{R}^d &\rightarrow [0, 1] \\ \mathbf{x} &\mapsto \mathbf{P}(X^{(1)} \leq x^{(1)}, \dots, X^{(d)} \leq x^{(d)}). \end{aligned}$$

Convergence in Probability in Higher Dimension

To make sense of the consistency statement $\hat{\theta}_n^{\text{MLE}} \xrightarrow[p]{n \rightarrow \infty} \theta^*$ where the MLE $\hat{\theta}_n^{\text{MLE}}$ is a random vector, we need to know what convergence in probability means in higher dimensions. But this is no more than convergence in probability for **each component**.

Let $\mathbf{X}_1, \mathbf{X}_2, \dots$ be a sequence of random vectors of size $d \times 1$, i.e. $\mathbf{X}_i = \begin{pmatrix} X_i^{(1)} \\ \vdots \\ X_i^{(d)} \end{pmatrix}$.

Let $\mathbf{X} = \begin{pmatrix} X^{(1)} \\ \vdots \\ X^{(d)} \end{pmatrix}$ be another vector of size $d \times 1$.

Then

$$\mathbf{X}_n \xrightarrow[p]{n \rightarrow \infty} \mathbf{X} \iff X_n^{(k)} \xrightarrow[p]{n \rightarrow \infty} X^{(k)} \text{ for all } 1 \leq k \leq d.$$

In other words, the sequence $\mathbf{X}_1, \mathbf{X}_2, \dots$ **converges in probability** to \mathbf{X} if and only if each component sequence $X_1^{(k)}, X_2^{(k)}, \dots$ converges in probability to $X^{(k)}$.

Hence, for example, in the Gaussian model $((-\infty, \infty), \{\mathcal{N}(\mu, \sigma^2)\}_{(\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_{>0}})$, consistency of the MLE $\hat{\theta}_n^{\text{MLE}} = \begin{pmatrix} \hat{\mu} \\ \widehat{\sigma^2} \end{pmatrix}$ means that $\hat{\mu}$ and $\widehat{\sigma^2}$ are consistent estimators of μ^* and $(\sigma^2)^*$, respectively.

Remark: You can check that this condition is equivalent to the following definition of convergence in probability, which is a straightforward generalization of the 1-dimensional case:

$$P(\{\omega \in \Omega : |\mathbf{X}_n(\omega) - \mathbf{X}(\omega)| < \epsilon\}) \xrightarrow{n \rightarrow \infty} 1 \quad \text{for any } \epsilon > 0.$$

Consistency of the MLE of a Uniform Model

1/1 point (graded)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}[0, \theta^*]$ where θ^* is an unknown parameter. We construct the associated statistical model $(\mathbb{R}_{\geq 0}, \{\text{Unif}[0, \theta]\}_{\theta > 0})$

Consider the maximum likelihood estimator $\hat{\theta}_n^{\text{MLE}} = \max_{i=1, \dots, n} X_i$.

Which of the following are true about $\hat{\theta}_n^{\text{MLE}}$. (Choose all that apply.)

- ☒ $\max_{i=1, \dots, n} X_i$ is a consistent estimator ☐
- ☒ For any $0 < \epsilon \leq \theta^*$, $\mathbf{P}\left(\left|\max_{i=1, \dots, n} X_i - \theta^*\right| \geq \epsilon\right) \rightarrow 0$ as $n \rightarrow \infty$ ☐
- ☐ For any $0 < \epsilon \leq \theta^*$, $\mathbf{P}\left(\left|\max_{i=1, \dots, n} X_i - \theta^*\right| \geq \epsilon\right) \rightarrow c$ as $n \rightarrow \infty$, where $c > 0$ is a constant
- ☒ For any $0 < \epsilon \leq \theta^*$, $\mathbf{P}\left(\left|\max_{i=1, \dots, n} X_i - \theta^*\right| \geq \epsilon\right) = \left(\frac{\theta^* - \epsilon}{\theta^*}\right)^n$ ☐

☐

Solution:

Choices 1, 2, and 4 are true because of the following proof for consistency of this ML estimator. Let $0 < \epsilon \leq \theta^*$:

$$\begin{aligned} \mathbf{P}\left(\left|\max_{i=1, \dots, n} X_i - \theta^*\right| \geq \epsilon\right) &= \mathbf{P}\left(\theta^* - \max_{i=1, \dots, n} X_i \geq \epsilon\right) \\ &= \mathbf{P}\left(\max_{i=1, \dots, n} X_i \leq \theta^* - \epsilon\right) \\ &= \left(\frac{\theta^* - \epsilon}{\theta^*}\right)^n \longrightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

Choice 3 is not true because if a sequence (the relevant sequence here is $\left(\frac{\theta^* - \epsilon}{\theta^*}\right)^n$) converges to a limit, then the limit is unique.

提交

你已经尝试了1次（总共可以尝试2次）

☐ Answers are displayed within the problem

讨论

显示讨论

主题: Unit 3 Methods of Estimation:Lecture 10: Consistency of MLE, Covariance Matrices, and Multivariate Statistics / 4. Consistency of Maximum Likelihood Estimator