

EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the [Privacy Policy](#).



[Lecture 8: Distance measures](#)

13. Parameter Estimation via KL

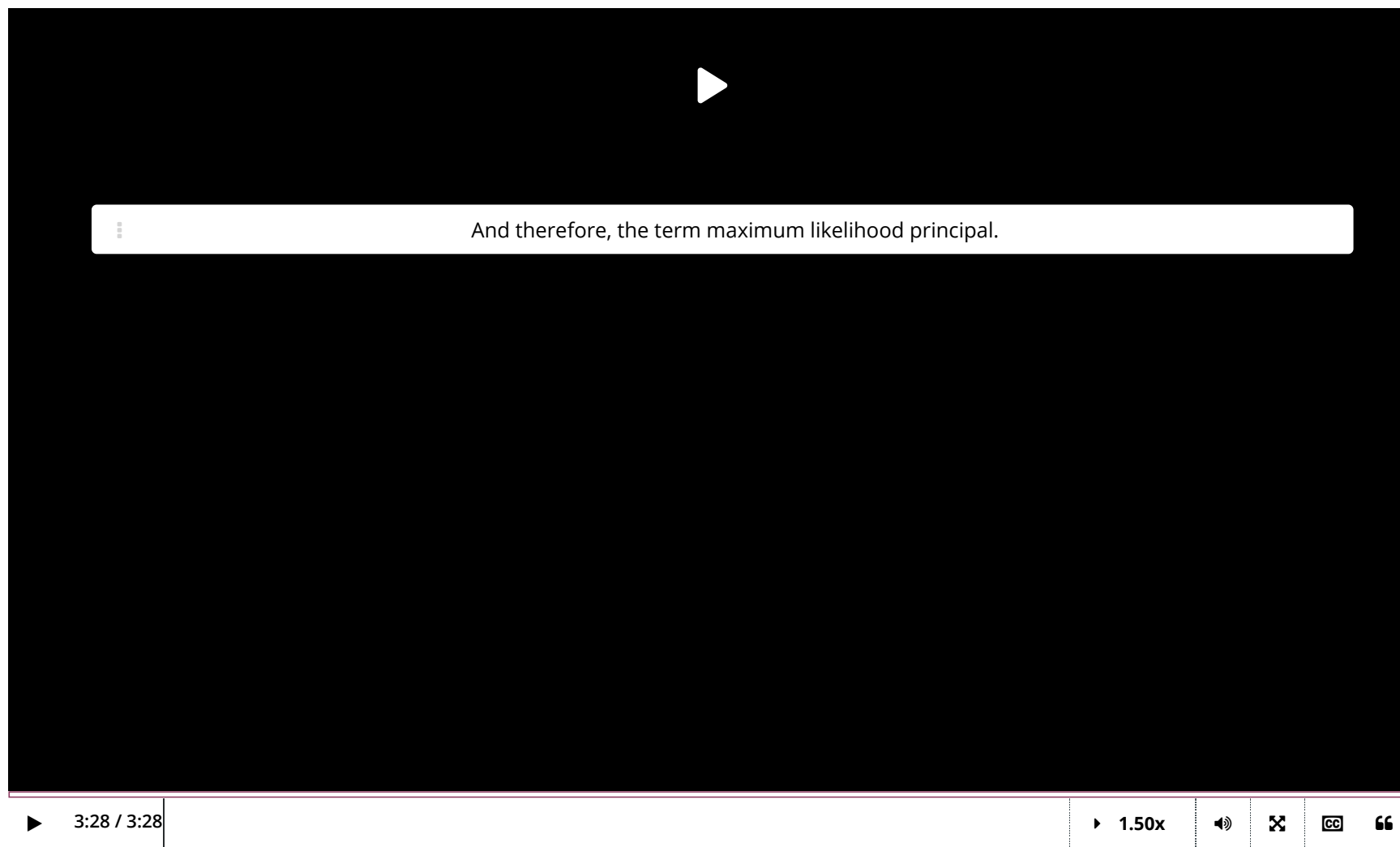
[Course](#) > [Unit 3 Methods of Estimation](#) > [between distributions](#)

> Divergence

Currently enrolled in **Audit Track** (expires December 25, 2019) [Upgrade \(\\$300\)](#)

13. Parameter Estimation via KL Divergence

Deriving the Maximum Likelihood Estimator



Video

[Download video file](#)

Transcripts

[Download SubRip \(.srt\) file](#)

[Download Text \(.txt\) file](#)

The next four problems concern the following statistical set-up.

You observe discrete random variables

$$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta^*}$$

where θ^* is the true parameter. You construct an associated statistical model $(E, \{P_\theta\}_{\theta \in \mathbb{R}})$. The sample space E is discrete.

Intuitively, your goal is to find an estimator $\hat{\theta}_n \in \mathbb{R}$ so that the distributions $P_{\hat{\theta}_n}$ and P_{θ^*} are close. Precisely, you want to find an estimator $\hat{\theta}_n \in \mathbb{R}$ so that the quantity

$$\text{KL}(P_{\theta^*}, P_{\hat{\theta}_n})$$

is as small as possible.

This approach will naturally lead to the construction of the **maximum likelihood estimator**.

Finding a Minimizer of KL Divergence

1/1 point (graded)

Consider the optimization problem in which we minimize the KL divergence between P_{θ^*} , the true distribution, and P_θ . Formally, we want to solve

$$\min_{\theta \in \mathbb{R}} \text{KL}(P_{\theta^*}, P_\theta).$$

We are not so much interested in the minimum value attained by the objective function $\text{KL}(P_{\theta^*}, P_\theta)$, but rather the value of θ where the minimum is attained. We refer to such a θ as a **minimizer**.

Let's suppose that there is a unique minimizer for the above optimization problem—i.e., if m is the minimum value of $\text{KL}(P_{\theta^*}, P_\theta)$, there is only one point θ_{\min} such that

$$m = \text{KL}(P_{\theta^*}, P_{\theta_{\min}}).$$

For which θ is the minimum value of $\text{KL}(P_{\theta^*}, P_{\theta})$ attained? (Equivalently, what is θ_{\min} ?)

☒ θ^* ☐ θ ☐ 0☐ None of the above.**Solution:**

The KL divergence is nonnegative, so $\text{KL}(P_{\theta^*}, P_{\theta}) \geq 0$. The right-hand side is achieved if we set $\theta = \theta^*$: $\text{KL}(P_{\theta^*}, P_{\theta^*}) = 0$. Since the minimizer is unique by assumption, we conclude that the minimum value is attained at $\theta = \theta^*$.

Remark: The assumption that there is a unique minimizer holds if we are given that the parameter θ is identified. Here is why: since KL divergence is definite, $\text{KL}(P_{\theta^*}, P_{\theta}) = 0$ if and only if P_{θ^*} and P_{θ} are the same distribution. And if θ is identified, this implies that $\theta = \theta^*$.

You have used 1 of 2 attempts

 Answers are displayed within the problem

Can we Minimize KL Divergence Directly?

2/2 points (graded)

Let's use the same statistical set-up as above. Recall that you have access to the iid samples X_1, \dots, X_n . You use these samples to build an estimator $\hat{\theta}_n$. Can you compute

$$\text{KL}(P_{\hat{\theta}_n}, P_{1/2})$$

without knowing θ^* , the true parameter?

☒ Yes

☐ No



Can you compute

$$\text{KL}(P_{\theta^*}, P_{1/2})$$

without knowing θ^* ?

☐ Yes

☒ No



Solution:

In general, we can compute $\text{KL}(\mathbf{P}, \mathbf{Q})$ if and only if we know both distributions \mathbf{P} and \mathbf{Q} . Moreover, by our statistical model, we can compute P_θ if and only if we know the real number θ . Putting these last two facts together, we can compute

$$\text{KL}(P_{\hat{\theta}_n}, P_{1/2})$$

because $\hat{\theta}_n$ is known—it is an estimator so its expression does not depend on θ^* , the true parameter. However, regardless of how many samples we take, we cannot compute $\text{KL}(P_{\theta^*}, P_{1/2})$ exactly because the distribution P_{θ^*} is unknown.

Remark: Since we cannot even compute the function $\text{KL}(P_{\theta^*}, P_{\theta})$ for general θ , this implies that the optimization problem

$$\min_{\theta \in \mathbb{R}} \text{KL}(P_{\theta^*}, P_{\theta})$$

cannot be solved exactly, regardless of the number of samples we have. So to estimate the minimizer of this optimization problem (which is the true parameter θ^*) we will have to consider an approximation for $\text{KL}(P_{\theta^*}, P_{\theta})$.

Submit

You have used 1 of 2 attempts

i Answers are displayed within the problem

Finding the Minimizer for an Approximation of KL Divergence

1/1 point (graded)

We use the same statistical set-up as above. Recall that $X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta^*}$. Let p_{θ} be the pmf of P_{θ} .

Which of the following is a (weakly) **consistent** estimator for

$$\mathbb{E}_{\theta^*} [\ln p_{\theta}(X)] = \sum_{x \in E} p_{\theta^*} \ln p_{\theta}(x) ?$$

☐ $\frac{1}{n} \sum_{i=1}^n X_i$

☒ $\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i))$

☐ $\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta^*}(X_i)) - \frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i))$

☐ $\theta^* - \mathbb{E}_{\theta^*}[\ln p_{\theta^*}]$



Solution:

By the law of large numbers, $\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i)) \rightarrow \mathbb{E}_{\theta^*}[\ln p_{\theta}]$ in probability. Hence, the second choice is correct.

Remark 1: The KL divergence between P_{θ^*} and P_{θ} can be written

$$\text{KL}(P_{\theta^*}, P_{\theta}) = \sum_{x \in E} p_{\theta^*} \ln p_{\theta^*}(x) - \sum_{x \in E} p_{\theta^*} \ln p_{\theta}(x) = \mathbb{E}_{\theta^*}[\ln p_{\theta^*}(X)] - \mathbb{E}_{\theta^*}[\ln p_{\theta}(X)]$$

where $X \sim P_{\theta^*}$.

Remark 2: While we can't find θ that minimizes $\text{KL}(P_{\theta^*}, P_{\theta})$, we can find θ that minimizes

$$\hat{\text{KL}}(P_{\theta^*}, P_{\theta}) := \mathbb{E}_{\theta^*}[\ln p_{\theta^*}] - \frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i)).$$

Here's why: the first term on the RHS, $\mathbb{E}_{\theta^*}[\ln p_{\theta^*}]$, does not depend on θ . Hence, the θ that minimizes $\hat{\text{KL}}(P_{\theta^*}, P_{\theta})$ is the same as the θ that minimizes $-\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i))$.

Submit

You have used 1 of 2 attempts

i Answers are displayed within the problem

Deriving the Maximum Likelihood Estimator

1/1 point (graded)

We use the same statistical set-up as above. Recall that p_{θ} is the pmf of P_{θ} and $X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta^*}$.

Suppose that θ_{\min} is a minimizer for the function

$$f(\theta) := -\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i))$$

Which of the following functions is also minimized at θ_{\min} ?

☐ $g_1(\theta) = -\prod_{i=1}^n p_{\theta}(X_i)$

☐ $g_2(\theta) = 25 - \prod_{i=1}^n p_{\theta}(X_i)$

☐ $g_3(\theta) = h(\theta^*) - \prod_{i=1}^n p_\theta(X_i)$ where h is a function of θ^* that does **not** depend on θ .

☐ $g_4(\theta) = \theta^* - \prod_{i=1}^n p_\theta(X_i)$

☒ All of the above



Solution:

Observe that rescaling by n does not change where the minimum of a function is attained. Hence, $f(\theta)$ and $nf(\theta)$ have the same minimizer. Next, by the addition property of logarithms,

$$nf(\theta) = \sum_{i=1}^n \ln(p_\theta(X_i)) = \ln\left(\prod_{i=1}^n p_\theta(X_i)\right).$$

Since \ln is an increasing function, the function

$$\theta \mapsto \prod_{i=1}^n p_\theta(X_i)$$

has the same minimizer as $\ln\left(\prod_{i=1}^n p_\theta(X_i)\right)$. Thus the first choice is correct.

Moreover, the second and third choices are also correct. Whenever we have an optimization problem

$$\min_{\theta \in \mathbb{R}} C + g(\theta)$$

where C does not depend on θ , then the above will have the same minimizer as the optimization problem

$$\min_{\theta \in \mathbb{R}} g(\theta).$$

In the second choice, $C = 25$ (which is independent of θ), and in the third choice, $C = h(\theta^*)$ (which by assumption is independent of θ).

Remark 1: The quantity

$$\hat{\theta}_n := \text{maximizer of } \prod_{i=1}^n p_{\theta}(X_i)$$

is referred to as the **maximum likelihood estimator**. Note that this is the same as the estimator

$$\hat{\theta}_n := \text{minimizer of } -\frac{1}{n} \sum_{i=1}^n \ln(p_{\theta}(X_i))$$

considered in Remark 2 in the solution of the previous problem.

Remark 2: Under certain technical conditions, the maximum likelihood estimator is guaranteed to (weakly) converge to the true parameter θ^* .

Submit

You have used 1 of 3 attempts

i Answers are displayed within the problem

Discussion

Topic: Unit 3 Methods of Estimation:Lecture 8: Distance measures between distributions / 13. Parameter Estimation via KL Divergence

[Hide Discussion](#)[Add a Post](#)

Show all posts ▼

by recent activity ▼

There are no posts in this topic yet.

✕

Audit Access Expires Dec 24, 2019

You lose all access to this course, including your progress, on Dec 24, 2019.

Upgrade by Nov 4, 2019 to get unlimited access to the course as long as it exists on the site. [Upgrade now](#)

[Learn About Verified Certificates](#)

© All Rights Reserved