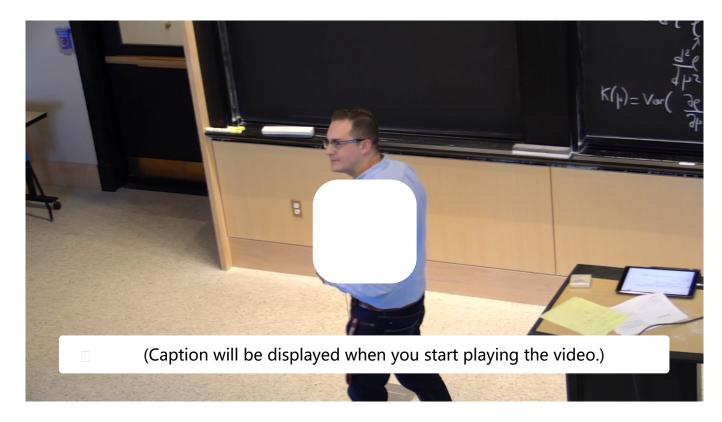


8. Applying Huber's loss to the

课程 □ Unit 3 Methods of Estimation □ Lecture 12: M-Estimation □ Laplace distribution

# 8. Applying Huber's loss to the Laplace distribution Applying Huber's loss to the Laplace distribution



Start of transcript. Skip to the end.

If I was really in a great mood, and I felt like integrating 1 plus 1 over x squared

and things like this, and do some nasty integration

by parts, I would actually show you the example of what happens when you apply this m estimator to the [? cowshed ?] distribution.

In the meantime, let's just do it for a simpler

0:00 / 0:00

□ 1.0x

视频

下载视频文件

字幕

下载 SubRip (.srt) file

下载 Text (.txt) file



### The Laplace distribution

3/3 points (graded)

The **Laplace distribution** (also known as the **double-exponential distribution** ) is a continuous distribution with **location parameter**  $m \in \mathbb{R}$  and density given by

$$f_{m}\left( x
ight) =rac{1}{2}e^{-\leftert x-m
ightert }.$$

Let X denote a Laplace random variable with location parameter set to be m=0.

What is  $\mathbb{E}\left[ oldsymbol{X}
ight]$ ?

0

☐ **Answer**: 0

Does the variance  $\sigma^2 = \mathbb{E}\left[\left(X - \mathbb{E}\left[X
ight]\right)^2
ight]$  exist?

Yes

No

Hint: The function  $x^k e^{-x}$  is integrable, i.e.  $\int_{-\infty}^{\infty} x^k e^{-x} dx$  is finite for all k.

The distribution of X is symmetric in the sense that X and X have the same distribution. X

- lacksquare The function  $\ln f_m\left(x
  ight)$  has a continuous first derivative.
- $^{ullet}$  For any integer k>0, the k-th moment  $\mathbb{E}\left[X^k
  ight]$  exists.  $\square$

#### **Solution:**

For the first question, observe that the function  $xe^{-|x|}$  is odd and also integrable. Therefore,

$$\mathbb{E}\left[X
ight] = \int_{-\infty}^{\infty} rac{1}{2} x e^{-|x|} \, dx = 0.$$

For the second question, the function  $x^2e^{-|x|}$  is integrable. Hence,

$$\mathbb{E}\left[\left(X-\mathbb{E}\left[X
ight]
ight)^{2}
ight]=\mathbb{E}\left[X^{2}
ight]=\int_{-\infty}^{\infty}rac{1}{2}x^{2}e^{-\left|x
ight|}\,dx$$

For the third question, we examine the choices in order.

- "The distribution of X is symmetric in the sense that X and -X have the same distribution." is correct. This is because the density  $\frac{1}{2}e^{-|x|}$  is an even function.
- "The function  $\ln f_m(x)$  has a continuous first derivative." is incorrect. This is because  $\ln f_m(x) = -|x-m|$ , which is not differentiable at x=m.
- "For any integer k>0, the k-th moment  $\mathbb{E}\left[X^k\right]$  exists." is correct. In general, the function  $x^ke^{-|x|}$  is integrable on  $\mathbb{R}$ , so the k-th moment  $\mathbb{E}\left[X^k\right]$  exists for all k>0.

提交 你已经尝试了3次(总共可以尝试3次)

☐ Answers are displayed within the problem

### The Sample Median

4/4 points (graded)

Let  $S = x_1 < x_2 < \ldots < x_n$  denote a sorted list of numbers. We define the **elementary median \operatorname{med}\_e(S)** to be

$$\operatorname{med}_e\left(S
ight) := \left\{egin{array}{ll} x_{\lceil n/2 
ceil} & ext{if} \ n ext{ is odd} \ rac{1}{2}(x_{n/2} + x_{n/2+1}) & ext{if} \ n ext{ is even} \end{array}
ight.$$

In other words, when n is odd, the median is the middle number when the set S is sorted from smallest to largest. If n is even, we can just define the median to be the average of both middle numbers. This definition is likely familiar from prior math classes.

A more advanced definition, useful for statistical purposes, is to define the **sample median**  $\operatorname{med}_s(S)$  of a sample  $S:=X_1,X_2,\ldots,X_n$  to be

$$\operatorname{med}\left(S
ight) := \operatorname{argmin}_{m} \sum_{i=1}^{n} \left|X_{i} - m
ight|.$$

While the elementary median is unique, this is not always the case for the **sample median**, as you will see in the next few questions.

#### **Solution:**

None of the above.

For the first question, the elementary median of  $m{S}$  is  $m{2}$ . For the second question, we want to find  $m{m}$  such that

$$F(m) = |1-m| + |2-m| + |3-m|$$

is as small as possible. Note that F(m) is a piecewise linear function with discontinuities in its first derivative precisely at the points (1,3),(2,2), and (3,3) (one corresponding to each summand of the form |x-m|). Drawing the line segments connecting these three points, we see that the minimizer of F(m) is at m=2. Therefore, for the second question, the sample median is 2.

For the third question, the elementary median is 2.5. For the fourth question, we want to find m such that

$$G\left(m
ight) = \left|1-m
ight| + \left|2-m
ight| + \left|3-m
ight| + \left|4-m
ight|$$

is as small as possible. As before, G(m) is a piecewise linear function. Its discontinuities are at the points (1,6), (2,4), (3,4), and (4,6). Hence, G(m) is a horizontal line segment on the interval [2,3], and G(m) has positive slope otherwise. Hence, the correct response to the fourth question is "Any number in the closed interval [2,3].".

**Remark:** In general, one can show that for an ordered sample  $S = x_1 < \ldots < x_n$  that

$$\operatorname{med}\left(S
ight) = egin{cases} x_{\lceil n/2 
ceil} & ext{if $n$ is odd} \ \operatorname{Any number in the interval } \left[x_{n/2}, x_{n/2+1}
ight] & ext{if $n$ is even} \end{cases}$$

• Any number in the closed interval [2,3].  $\Box$ 

□ Answers are displayed within the problem

## Maximum Likelihood Estimator for the Laplace Distribution

2/2 points (graded)

Consider a Laplace statistical model  $(\mathbb{R},\{P_m\}_{m\in\mathbb{R}})$  where  $P_m$  denotes the Laplace distribution with location parameter m. Let  $X_1,\ldots,X_n\stackrel{iid}{\sim}P_{m^*}$  denote a sample from a Laplace distribution with unknown parameter  $m^*$ . Recall that the density of  $P_m$  is given by

$$f_{m}\left( x
ight) =rac{1}{2}e^{-\leftert x-m
ightert }.$$

What is the log-likelihood  $\ell_n\left(X_1,\ldots,X_n;m
ight)$  for this statistical model?

$$left -n\ln{(2)} - \sum_{i=1}^n |X_i - m| \; \Box$$

$$0 - \sum_{i=1}^n |X_i - m|$$

$$\bigcirc -n\ln{(2)} - \sum_{i=1}^n{(X_i-m)^2}$$

$$-n \ln \left(2\right) + \sum_{i=1}^{n} |X_i - m|$$

Recall that the maximum likelihood estimator  $\widehat{m}_n^{ ext{MLE}}$  is given by

$$\widehat{m}_{n}^{ ext{MLE}} = \mathop{
m argmin}_{m \in \mathbb{R}} - \ell_{n}\left(X_{1}, \ldots, X_{n}; m
ight).$$

Suppose you observe the sample

$$S = 0.5, 1.2, 0.6, -0.7, -0.2.$$

What is the value of the MLE for  $m^*$  for this data set? *Hint:* Use the previous question, in particular the remark at the end of the solution.

#### **Solution:**

For the first question, the likelihood for  $m{n}$  observations is given by

$$\prod_{i=1}^n f_m\left(X_i
ight) = rac{1}{2^n} \prod_{i=1}^n e^{-|x-m|}.$$

Therefore,

$$\ell_n\left(X_1,\ldots,X_n;m
ight) = -n\ln\left(2
ight) - \sum_{i=1}^n |X_i-m|.$$

For the second question, we need to minimize the quantity

$$5 \ln{(2)} + |m - 0.5| + |m - 1.2| + |m - 0.6| + |m + 0.2| + |m + 0.7|.$$

with respect to m. As stated in the previous quantity, any m that minimizes the above is a sample median of the data set S. Since S is odd, we have that  $\widehat{m}_n^{\text{MLE}} = 0.5$ .

提交 你已经尝试了1次(总共可以尝试3次)

Concept Question: Maximum Likelihood Estimator for the Laplace distribution 1/1 point (graded)
As in the previous problem, let  $\widehat{m}_n^{\text{MLE}}$  denote the MLE for an unknown parameter  $m^*$  of a Laplace distribution.

Can we apply the theorem for the asymptotic normality of the MLE to  $\widehat{m}_n^{\text{MLE}}$ ? (You must choose the correct answer that also has the correct explanation.)

- No, because the log-likelihood is not concave.
- $^{ullet}$  No, because the log-likelihood is not twice-differentiable, so the Fisher information does not exist.  $\Box$
- Yes, because the log-likelihood is concave.
- Yes, because the other technical conditions required to apply the theorem are satisfied.

#### **Solution:**

We examine the choices in order.

• "No, because the log-likelihood is not concave." is incorrect. A sum of concave functions is concave, and for any constant c, the function  $x \to -|x-c|$  is concave. Therefore, the log-likelihood

$$\ell_{n}\left(X_{1},\ldots,X_{n};m
ight)=-n\ln\left(2
ight)-\sum_{i=1}^{n}\left|X_{i}-m
ight|$$

is also concave. Hence, the reasoning for this response is incorrect.

- "No, because the log-likelihood is not twice-differentiable, so the Fisher information does not exist." is correct. This is because  $\ell_n\left(X_1,\ldots,X_n;m\right)$  has discontinuities in its first derivative with respect to m at  $m=X_i$  for  $i=1,\ldots,n$ .
- "Yes, because the log-likelihood is concave." is incorrect. Although the log-likelihood is concave, the Fisher information does not exist, as discussed in the analysis of the previous two responses.
- "Yes, because the other technical conditions required to apply the theorem are satisfied." is incorrect. The remaining technical conditions are not enough to guarantee asymptotic normality since the Fisher information does not exist, as discussed above.

提交

你已经尝试了1次(总共可以尝试2次)

☐ Answers are displayed within the problem

### Applying Huber's loss to a Laplace distribution I

2/2 points (graded)

As above, let  $m^*$  denote an unknown parameter for a Laplace distribution. In this problem, we will use the principles of M-estimation and the smoothness of Huber's loss function to construct an asymptotically normal estimator for  $m^*$ . Let  $P_m$  denote the Laplace distribution with location parameter m.

Recall Huber's loss is defined as

$$h_{\delta}\left(x
ight)=\left\{egin{array}{ll} rac{x^{2}}{2} & ext{if} \;\; |x|<\delta \ \delta\left(|x|-\delta/2
ight) & ext{if} \;\; |x|>\delta \end{array}
ight..$$

As computed in lecture, the derivative of Huber's loss is the **clip function**:

$$\operatorname{clip}_{\delta}\left(x
ight) := rac{d}{dx} h_{\delta}\left(x
ight) = \left\{egin{array}{ll} \delta & ext{if } x > \delta \ x & ext{if } -\delta \leq x \leq \delta \ -\delta & ext{if } x < -\delta \end{array}
ight.$$

Find the value of

$$\left.rac{\partial}{\partial m}\mathbb{E}_{X\sim P_{m^*}}\left[h_\delta\left(X-m
ight)
ight]
ight|_{m=m^*}.$$

Hint: You are allowed to switch the derivative and expectation.

0 □ **Answer:** 0

In the framework of M-estimation, our loss function is not Huber's loss itself, but rather

$$\rho\left(x,m\right):=h_{\delta}\left(x-m\right)$$

Recall the functions

$$egin{aligned} J\left(m
ight) &= \mathbb{E}\left[rac{\partial^2 
ho}{\partial m^2}(X_1,m)
ight] \ K\left(m
ight) &= \mathrm{Var}\left[rac{\partial 
ho}{\partial m}(X_1,m)
ight] \end{aligned}$$

Do the functions  $oldsymbol{K}$  and  $oldsymbol{J}$  exist for a Laplace statistical model?

- No, because the log-likelihood is not twice-differentiable.
- lacksquare No, because J(m) exists but K(m) does not.
- Yes, because the Fisher information is well-defined for a Laplace statistical model.
- $^{ullet}$  Yes, because the function  $ho\left(x,m
  ight)$  as defined above is twice-differentiable.  $\Box$

#### **Solution:**

The answer to the first question is  ${f 0}$ . To see this, observe that

$$egin{aligned} rac{\partial}{\partial m} \mathbb{E}_{X \sim P_m^*} \left[ h_\delta \left( X 
ight) 
ight] &= \mathbb{E}_{X \sim P_m^*} \left[ rac{\partial}{\partial m} h_\delta \left( X 
ight) 
ight] \ &= rac{1}{2} \int_{-\infty}^{\infty} \operatorname{clip}_\delta \left( x - m 
ight) e^{-|x - m^*|} \, dx \ &= rac{1}{2} igg( \delta \int_{m + \delta}^{\infty} e^{-|x - m^*|} \, dx - \delta \int_{-\infty}^{-\delta + m} e^{-|x - m^*|} \, dx + \int_{-\delta + m}^{\delta + m} x e^{-|x - m^*|} \, dx igg) \, . \end{aligned}$$

Applying the change of variables y = x - m, we have

$$=rac{1}{2}igg(\delta\int_{\delta}^{\infty}e^{-|y+m-m^{st}|}\,dy-\delta\int_{-\infty}^{-\delta}e^{-|y+m-m^{st}|}\,dy+\int_{-\delta}^{\delta}ye^{-|y+m-m^{st}|}\,dyigg)\,.$$

Setting  $m=m^*$  , we have

$$\left. rac{\partial}{\partial m} \mathbb{E}_{X \sim P_m^*} \left[ h_\delta \left( X 
ight) 
ight] 
ight|_{m=m^*} = rac{1}{2} \Biggl( \delta \int_\delta^\infty e^{-|y|} \ dy - \delta \int_{-\infty}^{-\delta} e^{-|y|} \ dy + \int_{-\delta}^\delta y e^{-|y|} \ dy \Biggr) = 0.$$

**Remark:** The function  $m \mapsto \mathbb{E}_{X \sim P_m^*}[h_\delta(X)]$  is strictly convex, so this means the loss function has a unique critical point, and this is where the minimum is attained. The above calculation guarantees that the minimum is at  $m = m^*$ , the value of the true parameter.

"No, because the log-likelihood is not twice-differentiable.", "No, because J(m) exists but K(m) does not.", "Yes, because the Fisher information is well-defined for a Laplace statistical model.", "Yes, because the function  $\rho(x,m)$  as defined above is twice-differentiable."

For the second question, we consider the responses in order.

- "No, because the log-likelihood is not twice-differentiable." is incorrect. In the problem "Huber's loss" on the page "Robust Statistics and Huber's Loss", we showed that  $\rho(x,m) = h_\delta(x-m)$  is twice-differentiable with respect to m.
- "No, because J(m) exists but K(m) does not." is also incorrect. Both K(m) and J(m) exist because ho is twice-differentiable, and its derivatives are integrable.
- "Yes, because the Fisher information is well-defined for a Laplace statistical model." is incorrect. The Fisher information does not exist for a Laplace statistical model, as was shown in the problem "Concept Question: Maximum Likelihood Estimator for the Laplace Distribution" on the page "Applying Huber's loss to the Laplace distribution."
- "Yes, because the function ho(x,m) as defined above is twice-differentiable." is the correct answer. In a previous problem, we showed that  $ho(x,m)=h_\delta(x-m)$  is twice-differentiable with respect to m.

提交

你已经尝试了2次(总共可以尝试3次)

□ Answers are displayed within the problem

### Applying Huber's loss to a Laplace distribution II

1/1 point (graded)

We use the same set-up as in the previous problem. Recall that  $m^*$  is an unknown location parameter for a Laplace distribution.

The M-estimator  $\widehat{m}$  for  $m^*$  associated to the loss function  $ho\left(x,m
ight)=h_{\delta}\left(x-m
ight)$  is given by

$$\widehat{m} = \mathop{
m argmin}_{m \in \mathbb{R}} rac{1}{n} \sum_{i=1}^n h_\delta \left( X_i - m 
ight).$$

Consider the slide Asymptotic normality in Unit 3. Suppose that the technical conditions in 3. of slide are satisfied. Also, note that for any fixed x, the function  $m\mapsto h_\delta\left(x-m\right)$  is strictly convex. (You can see this by observing  $h_\delta\left(x-m\right)=h_\delta\left(m-x\right)$ , so the graph of  $h_\delta\left(x-m\right)$  as a function of x for a fixed x is the same as the graph of x for a fixed x for a fixed x is the calculation of x for a fixed x for a fixed x is the calculation of x for a fixed x form

Can we apply the theorem on the slide Asymptotic normality to conclude that  $\widehat{m{m}}$  is asymptotically normal?

(Choose the correct answer, 'Yes' or 'No', that also has a correct explanation.)

- ullet No, because  $m^*$  is not the unique minimizer of the function  $m\mapsto \mathbb{E}_{X\sim P_{m^*}}\left[
  ho\left(X,m
  ight)
  ight].$
- ullet No, because J(m) is not invertible.
- extstyle ext

$^ullet$ Yes, because $m^*$ is the unique minimizer of the function $m\mapsto \mathbb{E}_{X\sim P_{m^*}}\left[ ho\left(X,m ight) ight]$ and $J\left(m ight)$ is invertible. $\Box$	

#### Solution:

"Yes, because  $m^*$  is the unique minimizer of the function  $m \mapsto \mathbb{E}_{X \sim P_{m^*}} \left[ \rho \left( X, m \right) \right]$  and  $J \left( m \right)$  is invertible." is the correct answer. As shown in the Remark in the solution to the previous problem,  $m^*$  is the unique minimizer of the loss function  $m \mapsto \mathbb{E}_{X \sim P_{m^*}} \left[ \rho \left( X, m \right) \right]$ . Moreover,

$$J\left( m
ight) =1-e^{-\delta}$$

as was shown in the lecture. Hence J(m) is invertible. Assuming that the required technical conditions are satisfied, we conclude that the estimator  $\widehat{m}$  is asymptotically normal.

提交

你已经尝试了1次(总共可以尝试2次)

☐ Answers are displayed within the problem

# Applying Huber's loss to the Laplace distribution: Computation

Statistical analysis

► Let  $J(\mu) = +\frac{\partial^2 Q}{\partial \mu \partial \mu^\top}(\mu)$  (= +\mathbb{E}\left[\frac{\partial^2 \rho}{\partial \mu \partial \mu^\T}(X\_1, \mu)\right] under

some regularity conditions).

Let  $K(\mu) = \operatorname{Cov} \left[ \frac{\partial \rho}{\partial \mu} (X_1, \dots) \right]$ 

**Remark:** In the log-likelihood case (write  $\mu = \theta$ ),

 $J(\theta) = K(\theta) = \mathcal{I}(\theta)$  (Fisher information)

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

So it's just delta squared, f0 of x dx.

[INAUDIBLE] with me?

All right, so now I have to do limit of competition.

So let's just write what this thing actually is.

So it's 2 times the integral between 0

and delta of x squared times e to the minus-

sorry, that's delta-- e to the minus x dx divided by 2.

So those two cancel.

That's just my Laplace.

I just plugged in f0 for what it is.

And then I have plus-- well, those two

are going to cancel again.

So I have delta squared integral between delta and infinity of e

to the minus lambda dx.

Agreed?

OK, let's split work.

□ 0:00 / 0:00

下载视频文件

**宁**昔

下载 SubRip (.srt) file

□ 1.0x

<u>下载 Text (.txt) file</u>

xuetangX.com 学堂在线

讨论

主题: Unit 3 Methods of Estimation:Lecture 12: M-Estimation / 8. Applying Huber's loss to the Laplace distribution

隐藏讨论

Add a Post

☐ All Posts

Applying Huber's loss to a Laplace distribution I; what does it mean to be differentiable?

question posted 3 days ago by **SergK** (Community TA)

We have seen in the forums, thanks mrBB, example of a function that is differentiable but not continuously differentiable; but this is tricky.



I thought that being differentiable means to have derivative at each point, so derivative of a differentiable function can't have simple jumps, because having different left and right derivatives means derivative does not exist, and so the function is not differentiable.

Now after reading the solution I feel like an idiot. What does it mean to be differentiable, in context of the solution ???

此帖对所有人可见。

添加回复	1 response
sudarsanvsr_mit (Staff) 3 days ago	
This is a great question. I feel like we've been a bit casual when dealing with the second derivative in the	e Huber loss case.
The second derivative $h_\delta''$ does not exist at $m-\delta$ and $m+\delta$ but that does not matter because we are $h_\delta''$ with respect to a continuous density.	interested in the expected value of
So is it OK to understand that "twice differentiable" really means twice differentiable <i>almost surely?</i>	
kentaro nakaishi 在3 days ago前发表	
"Assume that $m{\ell}$ is a.s. twice differentiable" Slide 35/54, lecture 3.	
Mark B2 在3 days ago前发表	
@kentaro_nakaishi: That is correct.	
<u>sudarsanvsr mit</u> (Staff) 在3 days ago前发表	
Thank both of you, @Mark_B2 and @sudarsanvsr_mit , for clarification!	
kentaro nakaishi 在3 days ago前发表	
That still makes me wonder how that differs from the absolute value function, that also is twice differentiable almost eve a single point. Is there an essential difference in having also a single discontinuity in the first derivative? I'd say, but didn't over, that that also doesn't stop us from calculating an expected value (or variance) w.r.t. a continuous density function.	•
Edit: gave it a bit more thought (but perhaps not enough yet). Is the problem with using the Maximum Likelihood Estimated distribution not so much the (twice-)differentiability of $ ho$ , but rather that we find a $J\left(m\right)$ that is zero and therefore not in	
mrBB (Community TA) 在about 11 hours ago前发表	
This topic belongs to measure theory. If $F(x)$ is differentiable just almost everywhere, we cannot expect the fundament calculus $\int_a^b F'(x)dx = F(b) - F(a)$ in general. It actually holds if and only If $F(x)$ is absolutely continuous (a stro	
continuity). So if $F(x)$ has a jump like the derivative of $ x $ , we cannot use various techniques of calculus.	inger notion of
Anyway this will explain why the staff won't go into the details here:)	
kentaro nakaishi 在about 5 hours ago前发表	
Thanks for that note, kentaro_nakaishi. I will chew on it a bit. It doesn't immediately answer the question to me why we be discontinuity in $\rho'$ but not in $\rho''$ when calculating $E\left[\left(\rho'\right)^2\right]$ and $E\left[\rho''\right]$ respectively. In both instances piecewise integration a long way. But guess you're right that the subject is probably more technical than we can expect the course material to $\rho'$	ion seems to get us
mrBB (Community TA) 在about 2 hours ago前发表	
Discontinuity in $\rho'$ prevents us from finding its extremum at the points of discontinuity whereas discontinuity in $\rho''$ does finding its expectation (or variance of $\rho'$ ).	n't prevent us from
Mark B2 在about an hour ago前发表	

Sure, but these questions, and my pondering, are related to calculating the Fisher information and not so much to finding the extremum of the log likelihood. To me it isn't immediately obvious that we <i>can't</i> calculate the Fisher information in the absence of a continuous derivative, but that the absence of a continuous second derivative is no impediment for doing that. Somehow the concept questions take this as something obvious, but it still isn't obvious to me.	
mrBB (Community TA) 在18 minutes ago前发表	
Does Fisher information has the meaning if extrema of likelihood couldn't be found?	
Mark B2 在about a minute ago前发表	
添加评论	
显示所有的回复	
添加一条回复:	
Preview	//
提交	

© 保留所有权利