

# Unit 3 - part 3

## Fisher Information, Asymptotic Normality of MLE; Method of Moments

### 2. Review: Covariance Matrices and the Log-Likelihood Function

Exercises due Jul 14, 2020 20:59 JST

[Bookmark this page](#)

Let  $\mathbf{X}$  be a random vector of dimension  $d \times 1$  with expectation  $\mu_{\mathbf{X}}$ . Recall from [Lecture 10](#) that the covariance matrix  $\Sigma$  is defined as the following matrix outer product:

$$\Sigma = \mathbb{E}[(\mathbf{X} - \mu_{\mathbf{X}})(\mathbf{X} - \mu_{\mathbf{X}})^T].$$

It can be shown (similar to the covariance of random variables  $X, Y$  in [Lecture 10](#)) that

$$\begin{aligned}\Sigma &= \mathbb{E}[\mathbf{X}\mathbf{X}^T] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^T \\ &= \mathbb{E}[\mathbf{X}\mathbf{X}^T] - \mu_{\mathbf{X}}\mu_{\mathbf{X}}^T.\end{aligned}$$

---

#### Review of Covariance Matrices

1/3 points (graded)

Consider the following random vector of dimension  $d \times 1$ :  $\mathbf{X} = [X^{(1)}, X^{(2)}, \dots, X^{(d)}]^T$  is equally likely to be one of  $[1, 0, \dots, 0]^T, [0, 1, \dots, 0]^T, \dots, [0, 0, \dots, 1]^T$ . That is,  $\mathbf{X}$  is equal to any of the unit vectors along the coordinate axes with probability  $\frac{1}{d}$ .

Let us compute the entries of the covariance matrix  $\Sigma_{ij} = \text{Cov}(X^{(i)}, X^{(j)})$ .

$$\text{Cov}(X^{(i)}, X^{(i)}) = \boxed{\frac{1}{d}} \quad \times \text{Answer: } \frac{1}{d} - \frac{1}{d^2}$$

$$\text{With } i \neq j, \text{Cov}(X^{(i)}, X^{(j)}) = \boxed{\frac{2}{d}} \quad \times \text{Answer: } -\frac{1}{d^2}$$

Is  $\Sigma$  a singular (i.e. not invertible) covariance matrix? **Note:** A matrix  $\Sigma$  is singular if  $\det(\Sigma) = 0$ .

Yes

No



**Solution:**

For any  $i \in \{1, 2, \dots, d\}$ ,

$$\begin{aligned}\text{Cov}(X^{(i)}, X^{(i)}) &= \text{Var}(X^{(i)}) \\ &= \frac{1}{d} - \frac{1}{d^2},\end{aligned}$$

as each  $X^{(i)}$  is equal to 1 with probability  $\frac{1}{d}$  and equal to 0 with probability  $1 - \frac{1}{d}$ .

For any  $i \neq j$ ,  $\mathbb{E}[X^{(i)}X^{(j)}] = 0$  as  $X^{(i)}$  and  $X^{(j)}$  are never both equal to 1 at the same time. Therefore,

$$\begin{aligned}\text{Cov}(X^{(i)}, X^{(j)}) &= \mathbb{E}[X^{(i)}X^{(j)}] - \mathbb{E}[X^{(i)}]\mathbb{E}[X^{(j)}] \\ &= -\frac{1}{d^2}.\end{aligned}$$

The covariance matrix looks as follows:

$$\Sigma = \begin{bmatrix} \frac{1}{d} - \frac{1}{d^2} & -\frac{1}{d^2} & \cdots & -\frac{1}{d^2} \\ -\frac{1}{d^2} & \frac{1}{d} - \frac{1}{d^2} & \cdots & -\frac{1}{d^2} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{d^2} & -\frac{1}{d^2} & \cdots & \frac{1}{d} - \frac{1}{d^2} \end{bmatrix}.$$

Adding all the rows and replacing row 1 with the result yields

$$\widehat{\Sigma} = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ -\frac{1}{d^2} & \frac{1}{d} - \frac{1}{d^2} & \cdots & -\frac{1}{d^2} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{d^2} & -\frac{1}{d^2} & \cdots & \frac{1}{d} - \frac{1}{d^2} \end{bmatrix}.$$

From the above, we can see that the determinant of  $\widehat{\Sigma}$  is equal to 0. This means that  $\Sigma$ , which is row-equivalent to  $\widehat{\Sigma}$ , is a singular covariance matrix.

## Dimensions of Gradient of Log-Likelihood Function

1/2 points (graded)

Let  $(E, (\mathbf{P}_\theta)_{\theta \in \Theta})$  be a statistical model associated with a random vector  $\mathbf{X}$  of dimension  $k \times 1$ . Let  $f_\theta(\mathbf{x})$  be the joint pdf of  $\mathbf{X}$  and let  $\theta \in \mathbb{R}^d$ .

Let the log-likelihood function associated with one observation of  $\mathbf{X}$  be denoted  $\ell_1(\mathbf{x}, \theta)$ . For simplicity, let  $\ell_1(\mathbf{x}, \theta)$  be denoted  $\ell(\theta)$ , where it is assumed that  $\mathbf{x}$  is fixed.

Assuming that  $\ell(\theta)$  is differentiable with respect to  $\theta$  for almost all  $\mathbf{x}$ , what are the dimensions of the gradient  $\nabla \ell(\theta)$ ?

Number of rows in  $\nabla \ell(\theta)$ :  ✖ Answer: d + 0\*k

Number of columns in  $\nabla \ell(\theta)$ :  ✓ Answer: 1 + 0\*d + 0\*k

**Solution:**

$\ell(\theta)$ , at any given  $\mathbf{x}$ , is a real-valued function of  $d$  variables in the parameter  $\theta \in \mathbb{R}^d$ .

Therefore, the gradient vector  $\nabla \ell(\theta)$  is of size  $d \times 1$ .

## Log-Likelihood Function of a Bernoulli-like Random Variable

0/1 point (graded)

Consider the following experiment: You take a coin that lands a head (H) with probability  $0 < p < 1$  and you toss it twice. Define  $X$  as the following random variable:

$$X = \begin{cases} 1 & \text{if outcome is HH} \\ 0 & \text{otherwise} \end{cases}$$

Let  $\ell(p)$  be the log-likelihood function of  $X$  when written as a random function, i.e. all of the  $x$  in the function written as  $X$ . What is  $\ell(p)$ ?

*Hint:* Write the pmf of  $X$  as a one-line formula.

(Enter  $\mathbf{X}$  for  $X$ , and  $\ln(y)$  for  $\ln(y)$ . Do not enter "log".)

$\ell(p) =$   ✖ Answer: 2\*X\*ln(p) + (1-X)\*ln(1-p^2)

**Solution:**

First, note that  $X$  takes on the value 1 with probability  $p^2$  and the value 0 with probability  $1 - p^2$ .

Finding the log-likelihood function involves writing down the pmf of  $X$  as a one-line equation:

$$f(x; p) = (p^2)^x \cdot (1 - p^2)^{1-x}, \quad \text{where } x \in \{0, 1\}.$$

Taking logarithm and replacing all  $x$  with  $X$  yields the desired log-likelihood function written as a random function.

**Note:** because  $X$  takes only two values  $\{0, 1\}$ , there is more than one way to write down the pmf, and consequently the log-likelihood function.

just showing the definition here

$$\ell(\theta) = \log L_1(X; \theta)$$

$$L_1(x; \theta) = f_\theta(x)$$

$$\nabla \ell(\theta) \in \mathbb{R}^d \text{ if } \theta \in \mathbb{R}^d$$

$\nearrow$   
pdf & pmf

$$I(\theta) = \text{Cov}(\nabla \ell(\theta)) = E[(\nabla \ell(\theta) (\nabla \ell(\theta))^T) - E[\nabla \ell(\theta)] E[\nabla \ell(\theta)]^T]$$

$I(\theta)$  is a  $d \times d$  matrix called Fisher information

Theorem :  $I(\theta) = -E[H \ell(\theta)]$

Let  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  denote a continuous statistical model. Let  $f_\theta(x)$  denote the pdf (probability density function) of the continuous distribution  $\mathbf{P}_\theta$ . Assume that  $f_\theta(x)$  is twice-differentiable as a function of the parameter  $\theta$ .

In the next few problems, you will derive the formula

$$I(\theta) = \int_{-\infty}^{\infty} \frac{\left( \frac{\partial f_\theta(x)}{\partial \theta} \right)^2}{f_\theta(x)} dx$$

using the definition  $I(\theta) = \text{Var}(\ell'(\theta))$  and the basic formula  $\text{Var}(X) = E[X^2] - E[X]^2$  for any random variable  $X$ .

For computations, it is sometimes convenient to use the above formula for the Fisher information.

**Note:** The derivation in the next set of problems is presented as a proof in the video that follows, but we encourage you to attempt these problems before watching the video.

## Deriving a Useful Formula for the Fisher Information I

2/2 points (graded)

Let  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  denote a statistical model for a continuous distribution  $\mathbf{P}_\theta$ . Let  $f_\theta$  denote the pdf (probability density function) of the continuous distribution  $\mathbf{P}_\theta$ . Recall that

$$\int_{-\infty}^{\infty} f_\theta(x) dx = 1$$

for all  $\theta \in \mathbb{R}$ .

For the next two questions, assume that you are allowed to interchange derivatives and integrals.

What is

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_\theta(x) dx ?$$

✓ Answer: 0.0

What is

$$\int_{-\infty}^{\infty} \frac{\partial^2}{\partial \theta^2} f_\theta(x) dx ?$$

✓ Answer: 0.0

### Solution:

Since we know  $\int_{-\infty}^{\infty} f_\theta(x) dx = 1$ , this implies that

$$\frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_\theta(x) dx = \frac{\partial}{\partial \theta} 1 = 0.$$

Since we are allowed to interchange the integral and derivative, this implies that

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_\theta(x) dx = 0.$$

Similarly for the second derivative,

$$\int_{-\infty}^{\infty} \frac{\partial^2}{\partial \theta^2} f_\theta(x) dx = \frac{\partial^2}{\partial \theta^2} \int_{-\infty}^{\infty} f_\theta(x) dx = \frac{\partial^2}{\partial \theta^2} 1 = 0.$$

**Remark:** If  $f$  is "nice enough," analytically speaking, then we can rigorously justify interchanging the integral and derivative.

## Deriving a Useful Formula for the Fisher Information II

1/1 point (graded)

As before, let  $f_\theta$  denote the pdf (probability density function) of the continuous distribution  $\mathbf{P}_\theta$ . By definition,

$$\ell(\theta) = \ln L_1(X, \theta) = \ln f_\theta(X)$$

where  $X \sim \mathbf{P}_\theta$ . Differentiating, we see

$$\ell'(\theta) = \frac{\partial}{\partial \theta} \ln f_\theta(X) = \frac{\frac{\partial}{\partial \theta} f_\theta(X)}{f_\theta(X)}.$$

What is

$$\mathbb{E}[\ell'(\theta)] = \mathbb{E}\left[\frac{\frac{\partial}{\partial \theta} f_\theta(X)}{f_\theta(X)}\right]?$$

0

✓ Answer: 0.0

(Note that  $X \sim \mathbf{P}_\theta$ .)

**Solution:**

Observe that

$$\mathbb{E}\left[\frac{\frac{\partial}{\partial \theta} f_\theta(X)}{f_\theta(X)}\right] = \int_{-\infty}^{\infty} \left( \frac{\frac{\partial}{\partial \theta} f_\theta(x)}{f_\theta(x)} \right) f_\theta(x) dx = \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_\theta(x) dx = 0,$$

by the computation in the previous question.

### Deriving a Useful Formula for the Fisher Information III

0/1 point (graded)

As before, let  $f_\theta$  denote the pdf (probability density function) of the continuous distribution  $\mathbf{P}_\theta$ . By definition,

$$\ell'(\theta) = \ln L_1(X, \theta) = \ln f_\theta(X)$$

where  $X \sim \mathbf{P}_\theta$ .

Using the previous question, which of the following are equal to  $\text{Var}(\ell'(\theta)) = \text{Var}\left(\frac{\partial}{\partial \theta} \ln f_\theta(X)\right)$ ? (Choose all that apply. )

$\mathcal{I}(\theta)$  ✓

$\mathbb{E}[\ell'(\theta)]$

$\mathbb{E}[(\ell'(\theta))^2]$  ✓

$\int_{-\infty}^{\infty} \frac{(\frac{\partial}{\partial \theta} f_\theta(x))^2}{f_\theta(x)} dx$  ✓

✗

#### Solution:

We consider the choices in order.

- By definition,  $\mathcal{I}(\theta) = \text{Var}(\ell'(\theta))$ , so the first answer choice  $\mathcal{I}(\theta)$  is correct.
- By the previous question,  $\mathbb{E}[\ell'(\theta)] = 0$ , so this answer choice is incorrect.
- By definition of variance,

$$\text{Var}(\ell'(\theta)) = \mathbb{E}[\ell'(\theta)^2] - \mathbb{E}[\ell'(\theta)]^2,$$

and  $\mathbb{E}[\ell'(\theta)] = 0$ , by the previous question. Hence,  $\mathbb{E}[(\ell'(\theta))^2] = \text{Var}(\ell'(\theta))$ , and so the answer choice  $\mathbb{E}[(\ell'(\theta))^2]$  is correct.

- The last choice  $\int_{-\infty}^{\infty} \frac{(\frac{\partial}{\partial \theta} f_\theta(x))^2}{f_\theta(x)} dx$  is correct because, using the previous bullet,

$$\begin{aligned} \text{Var}(\ell'(\theta)) &= \mathbb{E}[(\ell'(\theta))^2] \\ &= \mathbb{E}\left[\left(\frac{\frac{\partial}{\partial \theta} f_\theta(X)}{f_\theta(X)}\right)^2\right] \\ &= \int_{-\infty}^{\infty} \frac{(\frac{\partial}{\partial \theta} f_\theta(x))^2}{f_\theta(x)} dx. \end{aligned}$$

**Remark:** A convenient way to compute the Fisher information is to use the fourth answer choice, which gives the useful formula

$$\mathcal{I}(\theta) = \int_{-\infty}^{\infty} \frac{(\frac{\partial}{\partial \theta} f_\theta(x))^2}{f_\theta(x)} dx.$$

lecture version of above  
assuming that we can take the derivative of the integral  
 $f(x)$  is a pdf so integrates to 1

$X$  is continuous with pdf  $f_\theta(x) = L_1(x, \theta)$

$$\int f_\theta(x) dx = 1 \quad \frac{\partial}{\partial \theta} \int f_\theta(x) dx = 0$$

also

$$\boxed{\begin{aligned} \int \frac{\partial}{\partial \theta} L_1(x, \theta) dx &= 0 \quad (1) \\ \int \frac{\partial^2}{\partial \theta^2} L_1(x, \theta) dx &= 0 \quad (2) \end{aligned}}$$

$$\left. \begin{aligned} l(\theta) &= \sum \log L_1(x, \theta) \\ &= \frac{\sum L_1(x, \theta)}{L_1(x, \theta)} \end{aligned} \right\}$$

$$\mathbb{E}[\ell'(\theta)] = \int \frac{\frac{\partial}{\partial \theta} L_1(x; \theta)}{L_1(x; \theta)} L_1(x; \theta) dx = \mathbb{E}(b_{y^{(1)}})$$

$$\begin{aligned}\mathbb{V}_{\theta}[\ell'(\theta)] &= \mathbb{E}[(\ell'(\theta))^2] \\ &= \int \left( \frac{\frac{\partial}{\partial \theta} L_1(x; \theta)}{L_1(x; \theta)} \right)^2 L_1(x; \theta) dx \\ &= \int \left( \frac{\frac{\partial}{\partial \theta} L_1(x; \theta)}{L_1(x; \theta)} \right)^2 dx.\end{aligned}$$

$$\ell''(\theta) = \frac{\frac{\partial^2}{\partial \theta^2} L_1(x; \theta) L_1(x; \theta) - \left( \frac{\partial}{\partial \theta} L_1(x; \theta) \right)^2}{(L_1(x; \theta))^2}$$

$$\begin{aligned}-\mathbb{E}[\ell''(\theta)] &= - \int \frac{\frac{\partial^2}{\partial \theta^2} L_1(x; \theta) L_1(x; \theta) - \left( \frac{\partial}{\partial \theta} L_1(x; \theta) \right)^2}{(L_1(x; \theta))^2} dx \\ &= - \int \frac{\frac{\partial^2}{\partial \theta^2} L_1(x; \theta) dx}{L_1(x; \theta)^2} + \mathbb{V}_{\theta}(\ell'(\theta))\end{aligned}$$

$$= O(b_{y^{(2)}})$$

$$\Rightarrow -\mathbb{E}[\ell''(\theta)] = \mathbb{V}_{\theta}(\ell'(\theta))$$

### Definition of Fisher Information

Let  $\theta \in \Theta \subset \mathbb{R}^d$  and let  $(E, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$  be a statistical model. Let  $f_\theta(\mathbf{x})$  be the pdf of the distribution  $\mathbf{P}_\theta$ . Then, the Fisher information of the statistical model is

$$I(\theta) = \text{Cov}(\nabla \ell(\theta)) = -\mathbb{E}[\mathbf{H}\ell(\theta)],$$

where  $\ell(\theta) = \ln f_\theta(\mathbf{x})$ .

The definition when the distribution has a pmf  $p_\theta(\mathbf{x})$  is also the same, with the expectation taken with respect to the pmf.

slide 37

this way of writing the PMF for Bernoulli is a clever way to say that when  $x$  is 1 the pmf is  $x$  and when  $x$  is 0 it is  $1-p$  (instead of writing the  $1(x=1)$  odd function)

$$\begin{aligned} \text{pmf} \Rightarrow f_p(x) &= p^x (1-p)^{1-x} \\ \ell'(p) &= \frac{x}{p} - \frac{1-x}{1-p} & \ell''(p) &= -\frac{x}{p^2} + \frac{1-x}{(1-p)^2} \\ \text{Var}(\ell'(p)) &= \text{Var}\left(x\left(\frac{1}{p} + \frac{1}{1-p}\right)\right) = \frac{1}{p(1-p)} \end{aligned}$$

$$\begin{aligned} \text{Var}\left(\frac{x}{p(1-p)}\right) &= \frac{p(1-p)}{p^2(1-p)^2} = \frac{1}{p^2(1-p)} \\ \mathbb{E}[\ell''(p)] &= -\frac{\mathbb{E}[x]}{p^2} - \frac{1-\mathbb{E}[x]}{(1-p)^2} \\ &= -\frac{1}{p} - \frac{1}{1-p} = -\frac{1}{p(1-p)} \end{aligned}$$

## Fisher Information of the Binomial Random Variable

1/1 point (graded)

Let  $X$  be distributed according to the binomial distribution of  $n$  trials and parameter  $p \in (0, 1)$ . Compute the Fisher information  $\mathcal{I}(p)$ .

*Hint:* Follow the methodology presented for the Bernoulli random variable in the above video.

$\mathcal{I}(p):$

n/(p*(1-p))	✓ Answer: n/(p*(1-p))
-------------	-----------------------

$\frac{n}{p \cdot (1-p)}$
---------------------------

pmf	$\binom{n}{k} p^k q^{n-k}$
-----	----------------------------

### Solution:

The logarithm of the pmf of a binomial random variable  $X$ , treated as a random function, can be written as

$$\ell(p) \triangleq \ln \left( \binom{n}{X} \right) + X \ln p + (n - X) \ln (1 - p), \quad X \in \{0, 1, \dots, n\}.$$

The derivative of  $\ell(p)$  with respect to  $p$  is

$$\ell'(p) = \frac{X}{p} - \frac{n - X}{1 - p},$$

which means the second derivative is

$$\ell''(p) = -\frac{X}{p^2} - \frac{n - X}{(1 - p)^2}.$$

The Fisher information  $\mathcal{I}(p)$ , therefore, is

$$\begin{aligned} \mathcal{I}(p) &= -\mathbb{E} [\ell''(p)] = \mathbb{E} \left[ \frac{X}{p^2} + \frac{n - X}{(1 - p)^2} \right] \\ &= \frac{np}{p^2} + \frac{n - np}{(1 - p)^2} \\ &= \frac{n}{p(1 - p)}. \end{aligned}$$

## Fisher Information of a Bernoulli-Like Random Variable

0/1 point (graded)

Consider the following experiment: You take a coin that lands a head (H) with probability  $0 < p < 1$  and you toss it twice. Define  $X$  as the following random variable:

$$X = \begin{cases} 1 & \text{if outcome is HH} \\ 0 & \text{otherwise} \end{cases}$$

Compute the Fisher information  $\mathcal{I}(p)$ .

$\mathcal{I}(p):$

1/(p*(1-p^2))	✗ Answer: 4/(1-p^2)
---------------	---------------------

$\frac{1}{p \cdot (1-p^2)}$
-----------------------------

**Solution:**

Following the Bernoulli and binomial examples,

$$\ell(p) \triangleq 2X \ln p + (1-X) \ln(1-p^2), \quad X \in \{0, 1\}.$$

The derivative of  $\ell(p)$  with respect to  $p$  is

$$\ell'(p) = \frac{2X}{p} - 2p \cdot \frac{1-X}{1-p^2},$$

which means the second derivative is

$$\ell''(p) = -\frac{2X}{p^2} - 2 \cdot \frac{(1-X)}{1-p^2} - 4p^2 \cdot \frac{1-X}{(1-p^2)^2}.$$

The Fisher information  $\mathcal{I}(p)$ , therefore, is

$$\begin{aligned} \mathcal{I}(p) &= -\mathbb{E}[\ell''(p)] = \mathbb{E}\left[\frac{2X}{p^2} + 2 \cdot \frac{(1-X)}{1-p^2} + 4p^2 \cdot \frac{1-X}{(1-p^2)^2}\right] \\ &= \frac{2p^2}{p^2} + \frac{2(1-p^2)}{(1-p^2)} + 4p^2 \cdot \frac{1-p^2}{(1-p^2)^2} \\ &= 4 + \frac{4p^2}{1-p^2} \\ &= \frac{4}{1-p^2} \end{aligned}$$

### Fisher Information of a Modified Gaussian Random Vector

0/4 points (graded)

Let  $\mathbf{X}$  be a Gaussian random vector with **independent** components  $X^{(i)} \sim \mathcal{N}(\alpha + \beta t_i, 1)$  for  $i = 1, \dots, d$ , where  $t_i$  are known constants and  $\alpha$  and  $\beta$  are unknown parameters.

Compute the Fisher information matrix  $\mathcal{I}(\theta)$  using the formula  $\mathcal{I}(\theta) = -\mathbb{E}[\mathbf{H}\ell(\theta)]$ .

Use **S\_1** for  $\sum_{i=1}^d t_i$  and **S\_2** for  $\sum_{i=1}^d t_i^2$ .

$$\mathcal{I}(\theta)_{1,1} = \boxed{0} \quad \text{✖ Answer: } d + 0*S_1 + 0*S_2 \quad \mathcal{I}(\theta)_{1,2} = \boxed{0} \quad \text{✖ Answer: } 0*d + S_1 + 0*S_2$$

$$\mathcal{I}(\theta)_{2,1} = \boxed{0} \quad \text{✖ Answer: } 0*d + S_1 + 0*S_2 \quad \mathcal{I}(\theta)_{2,2} = \boxed{0} \quad \text{✖ Answer: } 0*d + 0*S_1 + S_2$$

*Hint:* Let  $\theta = [\alpha \ \beta]^T$  denote the parameters of the statistical model.  $\ell(\theta)$  is a real-valued function of  $\theta$  as given by the joint pdf at any fixed  $\mathbf{x}$ .

**Solution:**

Let  $\theta = [\alpha \ \beta]^T$  denote the parameters of the statistical model. The joint distribution of the random vector  $\mathbf{X}$  is equal to the product of the marginal distributions of its components,  $(2\pi^{-1/2} e^{-(x^{(i)} - \alpha - \beta t_i)^2})^d$ , by independence. Therefore, the random vector  $\mathbf{X}$  has the density

$$f_\theta(\mathbf{x}) = (2\pi)^{-\frac{d}{2}} e^{-\frac{1}{2} \sum_{i=1}^d (x^{(i)} - \alpha - \beta t_i)^2}, \quad \text{where } \mathbf{x} = [x^{(1)} \ x^{(2)} \ \dots \ x^{(d)}]^T \in \mathbb{R}^d.$$

Taking the log of the pdf yields (written as a random function)

$$\ell(\theta) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \left[ \sum_{i=1}^d \left( (X^{(i)} - \beta t_i)^2 - 2\alpha(X^{(i)} - \beta t_i) + \alpha^2 \right) \right]$$

Therefore,

$$\nabla \ell(\theta) = \begin{bmatrix} \sum_{i=1}^d (X^{(i)} - \beta t_i - \alpha) \\ \sum_{i=1}^d (t_i X^{(i)} - \beta t_i^2 - \alpha t_i) \end{bmatrix},$$

from which we can obtain the hessian

$$\mathbf{H}\ell(\theta) = \begin{bmatrix} \sum_{i=1}^d (-1) & \sum_{i=1}^d (-t_i) \\ \sum_{i=1}^d (-t_i) & \sum_{i=1}^d (-t_i^2) \end{bmatrix}.$$

Therefore,

$$\mathcal{I}(\theta) = -\mathbb{E}[\mathbf{H}\ell(\theta)] = \begin{bmatrix} d & \sum_{i=1}^d t_i \\ \sum_{i=1}^d t_i & \sum_{i=1}^d t_i^2 \end{bmatrix},$$

where the expectation is taken with respect to the pdf of the random vector  $\mathbf{X}$ . Since none of the entries of the hessian contained any  $X^{(i)}$ , the expectation was simply the hessian matrix itself.

## A Geometric View on the Fisher Information

1/1 point (graded)

Let  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  denote a statistical model. Recall that the MLE (maximum likelihood estimator) for one observation maximizes the log-likelihood for one observation, which is the random variable  $\ell(\theta) = \ln L_1(X, \theta)$  where  $X \sim \mathbf{P}_\theta$ . Suppose we observe  $X_1 = x_1$ , and now consider the graph of the function  $\theta \mapsto \ln L_1(x_1, \theta)$ .

What does the Fisher information  $\mathcal{I}(\theta)$  represent?

**Hint:** Use the definition  $\mathcal{I}(\theta) = -\mathbb{E}[\ell''(\theta)]$ .

It gives you an approximation for the true parameter  $\theta^*$ .

It tells you the average slope of the function  $\theta \mapsto \ln L(x_1, \theta)$  is.

It tells you, on average, how curved the function  $\theta \mapsto \ln L(x_1, \theta)$  is.

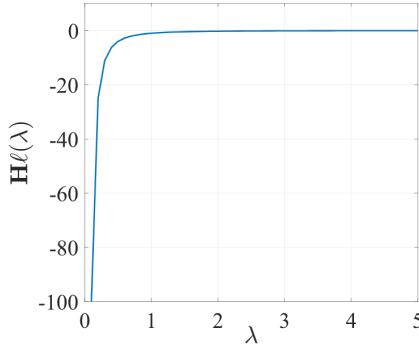


**Solution:**

The third choice "It tells you, on average, how curved the function  $\theta \mapsto \ln L(x_1, \theta)$  is." is correct. Recall that the Fisher information is also equal to the expected second-derivative of the log-likelihood:  $\mathcal{I}(\theta) = -\mathbb{E}[\ell''(\theta)]$ . Since the second-derivative measures concavity/convexity (how curved a function is at a particular point),  $\mathcal{I}(\theta)$  measures the *average* curvature of the function  $\theta \mapsto \ell(\theta) = \ln L_1(x_1, \theta)$ .

**Remark:** It turns out that the Fisher information tells how curved (on average) the log-likelihood  $\ln L_n(x_1, \dots, x_n, \theta)$  for several samples  $X_1 = x_1, \dots, X_n = x_n$  is. In particular,  $\mathcal{I}(\theta^*)$  tells how curved (on average) the log-likelihood is near the true parameter. As a rule of thumb, if the Fisher information  $\mathcal{I}(\theta^*)$  is large, then we expect the MLE to give a good estimate for  $\theta^*$ .

Consider the exponential statistical model with  $f_\lambda(x) = \lambda e^{-\lambda x}$ ,  $x \geq 0$  and  $\lambda \in (0, \infty)$ . The second derivative of  $\ell(\lambda)$ , which is  $\mathbf{H}\ell(\lambda) = \ell''(\lambda)$  (you will compute later in a homework exercise), does not depend upon  $x$  and is shown in the following figure.



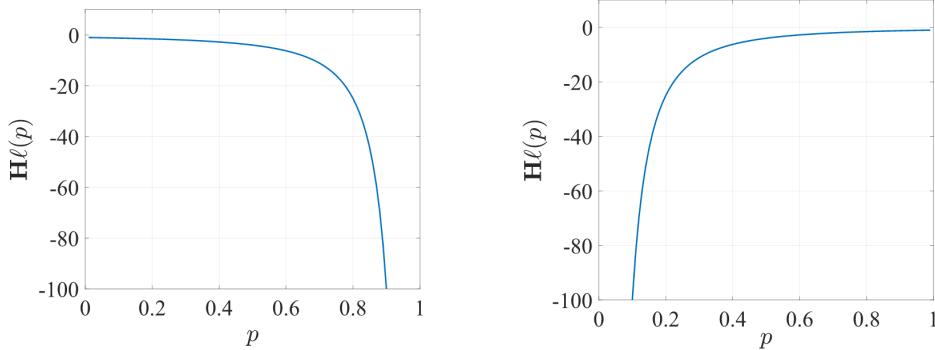
$\ell''(\lambda)$  for the exponential statistical model with parameter  $\lambda$  for all  $X$ .

The Fisher information,  $\mathcal{I}(\theta) = -\mathbb{E}[\mathbf{H}\ell(\theta)]$ , captures the negative of the expected curvature of  $\ell(\theta)$ . For example, for the exponential statistical model, the expected curvature of  $\ell(\lambda)$  is  $\ell''(\lambda)$  itself and this is shown in the figure above. The Fisher information in this case is always positive. The fact that  $\ell''(\lambda)$  is negative for all  $x$  also means that the log-likelihood function  $\ell(\lambda)$  is a concave function of  $\lambda$  for all  $x$ .

If we consider the Bernoulli statistical model with parameter  $p \in (0, 1)$ , we derived in a lecture video (also can be seen in the slides) that

$$\ell''(p) = -\frac{X}{p^2} - \frac{1-X}{(1-p)^2}, \quad X \in \{0, 1\}.$$

Here, we see that  $\ell''(p)$  is a random function that depends upon  $X$ . The following two figures show  $\ell''(p)$  for  $X = 0$  and  $X = 1$ .



$\ell''(p)$  for the Bernoulli statiscal model with parameter  $p \in (0, 1)$  for  $X = 0$  (left) and  $X = 1$  (right).

The **asymptotic normality of the ML estimator**, which will be discussed in the upcoming video, depends upon the Fisher information. For a one-parameter model (like the exponential and Bernoulli), the asymptotic normality result will say something along the lines of following: that the asymptotic variance of the ML estimator is inversely proportional to the value of Fisher information at the true parameter  $\theta^*$  of the statistical model. This means that if the value of Fisher information at  $\theta^*$  is high, then the asymptotic variance of the ML estimator for the statistical model will be low.

## Weak Consistency of the MLE

1/1 point (graded)

Let  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  denote a statistical model associated to a statistical experiment  $X_1, \dots, X_n \stackrel{iid}{\sim} \mathbf{P}_{\theta^*}$  for some true parameter  $\theta^*$  that we would like to estimate. You construct the maximum likelihood estimator  $\hat{\theta}_n^{MLE}$  for  $\theta^*$ . Which of the following conditions is **not** necessary for the MLE  $\hat{\theta}_n^{MLE}$  to converge to  $\theta^*$  in probability?

- The model  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  is identified. (Recall that the parameter  $\theta$  is identified if the map  $\theta \mapsto \mathbf{P}_\theta$  is injective.)
- For all  $\theta \in \Theta$ , the support of  $\mathbf{P}_\theta$  does not depend on  $\theta$ .
- The MLE  $\hat{\theta}_n^{MLE}$  is given by the sample average.
- The Fisher information  $I(\theta)$  is non-zero in an interval containing true parameter  $\theta^*$ . (Note that this is what it means for a  $1 \times 1$  matrix, a scalar, to be invertible.)



### Solution:

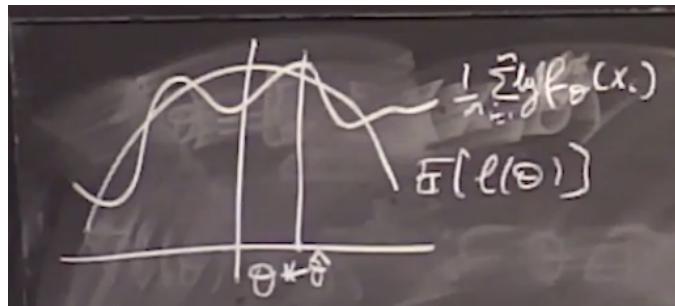
The power of the theorem for the convergence of the MLE is that it applies even in situations where the MLE is *not* the sample average. Hence, the third choice, "The MLE  $\hat{\theta}_n^{MLE}$  is given by the sample average.", is correct, as it is not an assumption needed for the theorem statement. On the other hand, the first, second, and fourth choices are all hypotheses in the theorem statement regarding the convergence of the MLE. Hence, "The model  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  is identified.", "For all  $\theta \in \Theta$ , the support of  $\mathbf{P}_\theta$  does not depend on  $\theta$ .", and "The Fisher information  $I(\theta)$  is invertible in an interval containing the true parameter  $\theta^*$ " are all incorrect responses.

continuation of slide 38

what is the Fisher Information for

using an estimator (squiggle line)

want closeness in y axis to translate to closeness in x axis



$$\frac{\partial}{\partial \theta} \sum_{i=1}^n l_i(\theta) = \sum_{i=1}^n l'_i(\theta) = 0 \quad \text{because } \hat{\theta} \text{ is maximum}$$

$$E[l'(\theta^*)] = 0$$

equal to 0 because  $\theta^*$  is the maximum of the expected log likelihood

1st order taylor expansion of  $l(\hat{\theta})$  around  $l(\theta^*)$

$$0 = \sum_{i=1}^n l_i'(\hat{\theta}) \simeq \sum_{i=1}^n [l_i'(\theta^*) + (\hat{\theta} - \theta^*) \underline{l_i''(\theta^*)}]$$

$$= \sum_{i=1}^n [(l_i'(\theta^*) - E[l_i'(\theta^*)]) + (\hat{\theta} - \bar{\theta}^*) \underline{l_i''(\theta^*)}]$$

now have sums of random variables minus their expectations  
so the central limit theorem comes in

$$\text{By CLT } \frac{1}{\sqrt{n}} \sum_{i=1}^n (l_i'(\theta^*) - E[l_i'(\theta^*)]) \xrightarrow{n \rightarrow \infty} N(0, \underbrace{\text{Var}(l_i'(\theta))}_{I(\theta)})$$

$$0 \simeq N(0, I(\theta^*)) + (\hat{\theta} - \theta^*) \frac{1}{\sqrt{n}} \sum_{i=1}^n l_i''(\theta^*)$$

$$0 \simeq N(0, I(\theta^*)) + \sqrt{n} (\hat{\theta} - \theta^*) \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n l_i''(\theta^*)}_{\xrightarrow{n \rightarrow \infty} -I(\theta^*)}$$

swapped the n around so now we know what the term on the right converges in probability to (Fisher information)

$$\Rightarrow \sqrt{n} (\hat{\theta} - \theta^*) \sim N(0, \frac{1}{I(\theta^*)})$$

which is the formula at the bottom of slide 38 (1/Fisher information)

Asymptotic = approaches but never touches the line

## Fisher Information and Asymptotic Normality of the MLE

1/1 point (graded)

Consider the statistical model  $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$  associated to the statistical experiment  $X_1, \dots, X_n \stackrel{iid}{\sim} \mathbf{P}_{\theta^*}$ , where  $\theta^*$  is the true parameter. Assume that the conditions of the theorem for the convergence of the MLE hold. Which of the following statements about the Fisher information  $\mathcal{I}(\theta)$  is true?

The Fisher information  $\mathcal{I}(\theta^*)$  at the true parameter gives a good approximation for  $\theta^*$ .

The Fisher information  $\mathcal{I}(\theta^*)$  at the true parameter determines the asymptotic mean of the random variable  $\hat{\theta}_n^{\text{MLE}}$ .

The Fisher information  $\mathcal{I}(\theta^*)$  at the true parameter determines the asymptotic variance of the random variable  $\hat{\theta}_n^{\text{MLE}}$ .



### Solution:

As stated in the theorem,

$$\sqrt{n}(\hat{\theta}_n^{\text{MLE}} - \theta^*)$$

converges to a normal random variable  $\mathcal{N}(0, \mathcal{I}(\theta^*)^{-1})$ . Hence, the Fisher information determines the asymptotic variance, and so the third choice is correct.

## Asymptotic Normality of the MLE

0/1 point (graded)

Consider the statistical model  $(\{0, 1\}, \{\text{Ber}(\theta)\}_{\theta \in (0, 1)})$ . Let  $\ell(\theta)$  denote the **log-likelihood of one observation** of this model. You observe samples  $X_1, \dots, X_n \sim \text{Ber}(\theta^*)$  and construct the MLE  $\hat{\theta}_n^{\text{MLE}}$  for  $\theta^*$ . By the theorem for the convergence of the MLE (you are allowed to assume that all necessary conditions for this theorem hold), this implies that

$$\sqrt{n}(\hat{\theta}_n^{\text{MLE}} - \theta^*) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, \sigma^2)$$

for some constant  $\sigma^2$  that depends on  $\theta^*$ . The quantity  $\sigma^2$  is referred to as the **asymptotic variance**. Use the theorem for the convergence of the MLE to find the expression for  $\sigma^2$ .

What is  $\sigma^2$ ? Express your answer in terms of  $T := \theta^*$ .

Type **T** for  $T$ , using the variable  $T$  to stand for  $\theta^*$ .

$$\sigma^2 = \boxed{1/T}$$

✖ Answer:  $T*(1-T)$

$\frac{1}{T}$

**Solution:**

We have that for this model the Fisher information is  $I(\theta) = \frac{1}{\theta} + \frac{1}{(1-\theta)} = \frac{1}{\theta(1-\theta)}$ . Applying the theorem for the convergence of the MLE,

$$\sqrt{n}(\hat{\theta}_n^{\text{MLE}} - \theta^*) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, I(\theta^*)^{-1})$$

Hence,

$$\sigma^2 = I(\theta^*)^{-1} = \theta^*(1 - \theta^*).$$

**Remark:** Alternatively, the asymptotic variance can be computed directly from the MLE, which is given, in the Bernoulli case, by the sample mean  $\bar{X}_n$ .

## Estimating Moments

0/1 point (graded)

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \mathbf{P}_{\theta^*}$  be integer-valued random variables and let  $(\mathbb{Z}, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$  be the associated statistical model. Let  $p_\theta$  denote the pmf of  $\mathbf{P}_\theta$ . Assume that for all  $\theta \in \Theta$ , the **k-th moment**

$$m_k(\theta) := \mathbb{E}[X^k] = \sum_{x \in \mathbb{Z}} x^k p_\theta(x)$$

exists for all  $k \geq 1$ . Use the law of large numbers to fill in the formula so that  $\hat{m}_k(\theta)$  is a consistent estimator for  $m_k(\theta)$ .

(Type  $\mathbf{X}_i$  for  $X_i$ .)

$$\hat{m}_K(\theta) = \frac{1}{n} \sum_{i=1}^n \boxed{k * \mathbf{X}_i}$$

✖ Answer:  $\mathbf{X}_i^k$

STANDARD NOTATION

**Solution:**

The weak law of large numbers implies that

$$\frac{1}{n} \sum_{i=1}^n X_i^k \rightarrow \mathbb{E}[X^k],$$

where the convergence is in probability.

## Mapping Parameters to Moments I

2/2 points (graded)

Let  $(\mathbb{R}, \{N(\mu, \sigma^2)\}_{\mu \in \mathbb{R}, \sigma > 0})$  be the statistical model of a normal random variable  $X$ . Let

$$m_k(\mu, \sigma) = \mathbb{E}[X^k]$$

denote the  $k$ -th moment of  $X$ . Let  $\psi : \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}^2$  be defined by  $\psi(\mu, \sigma) = (m_1(\mu, \sigma), m_2(\mu, \sigma))$ . (Since we have two parameters of interest,  $\mu$  and  $\sigma$ , it makes sense to work with the first two moments. The hope is that the two moments will uniquely determine the parameters of interest  $\mu$  and  $\sigma$ .)

Express  $m_1(\mu, \sigma)$  and  $m_2(\mu, \sigma)$  in terms of  $\mu$  and  $\sigma$ .

$$m_1(\mu, \sigma) = \boxed{\text{mu}} \quad \checkmark \text{ Answer: mu}$$

$\mu$

$$m_2(\mu, \sigma) = \boxed{\text{mu}^2 + \sigma^2} \quad \checkmark \text{ Answer: mu}^2 + \sigma^2$$

$\mu^2 + \sigma^2$

**STANDARD NOTATION**

**Solution:**

Note that

$$\begin{aligned} m_1(\mu, \sigma) &= \mathbb{E}[X] = \mu \\ m_2(\mu, \sigma) &= \mathbb{E}[X^2] = (\mathbb{E}[X])^2 + (\mathbb{E}[X^2] - (\mathbb{E}[X])^2) = \mu^2 + \sigma^2. \end{aligned}$$

Hence,  $\psi(\mu, \sigma) = (\mu, \mu^2 + \sigma^2)$ .

## Mapping Parameters to Moments II

2/3 points (graded)

Let

$$\psi : \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}^2 \\ (\mu, \sigma) \mapsto (m_1(\mu, \sigma), m_2(\mu, \sigma)).$$

denote the moments map considered in the previous problem, where  $m_k(\mu, \sigma)$  denotes the  $k$ -th moment of the distribution  $N(\mu, \sigma^2)$ .

Is  $\psi$  one-to-one on the domain  $\mathbb{R} \times (0, \infty)$ ? (Equivalently, given the outputs  $m_1$  and  $m_2$ , can we use them to uniquely reconstruct  $\mu \in \mathbb{R}$  and  $\sigma > 0$ ?)

Yes

No



If  $\psi$  is one-to-one on the given domain and  $\psi(\mu, \sigma) = (m_1, m_2)$ , what is  $\mu$  expressed in terms of  $m_1$  and  $m_2$ ? (If  $\psi$  is not one-to-one, enter 0.)

$$\mu = \boxed{m_1}$$

✓ Answer:  $m_1$

If  $\psi$  is one-to-one on the given domain and  $\psi(\mu, \sigma) = (m_1, m_2)$ , what is  $\sigma$  expressed in terms of  $m_1$  and  $m_2$ ? (If  $\psi$  is not one-to-one, enter 0.)

$$\sigma = \boxed{m_2 - (m_1)^2}$$

✗ Answer:  $\sqrt{m_2 - m_1^2}$

### Solution:

Note that

$$m_1(\mu, \sigma) = \mathbb{E}[X] = \mu \\ m_2(\mu, \sigma) = \mathbb{E}[X^2] = (\mathbb{E}[X])^2 + (\mathbb{E}[X^2] - (\mathbb{E}[X])^2) = \mu^2 + \sigma^2.$$

Hence,  $\psi(\mu, \sigma) = (\mu, \mu^2 + \sigma^2)$ . This function is one-to-one on the domain  $\mathbb{R} \times (0, \infty)$ . Since  $m_1(\mu, \sigma) = \mu$ , we can reconstruct the first parameter directly from the first moment:  $\mu = m_1$ .

Next, since we know  $m_2(\mu, \sigma) = \sigma^2 + \mu^2$ , we can back-solve for  $\sigma$ :

$$\sigma = \sqrt{m_2 - \mu^2} = \sqrt{m_2 - m_1^2}.$$

Above we took the positive square-root because we have insisted a priori that  $\sigma > 0$ .

**Remark:** Assuming that  $m_1$  and  $m_2$  are one of the outputs of  $\psi$ , we have essentially shown how to construct  $\psi^{-1}(m_1, m_2)$ . In general, computing inverses is a computationally difficult problem, but in this particular example, the function  $\psi$  is simple enough that it is possible to invert by hand.

## Method of Moments Concept Question I

2/2 points (graded)

Let  $(E, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$  denote a statistical model associated to a statistical experiment  $X_1, \dots, X_n \stackrel{iid}{\sim} \mathbf{P}_{\theta^*}$  where  $\theta^* \in \Theta$  is the true parameter. Assume that  $\Theta \subset \mathbb{R}^d$  for some  $d \geq 1$ . Let  $m_k(\theta) := \mathbb{E}[X^k]$  where  $X \sim \mathbf{P}_\theta$ .  $m_k(\theta)$  is referred to as the  **$k$ -th moment of  $\mathbf{P}_\theta$** . Also define the moments map:

$$\begin{aligned}\psi : \Theta &\rightarrow \mathbb{R}^d \\ \theta &\mapsto (m_1(\theta), m_2(\theta), \dots, m_d(\theta)).\end{aligned}$$

Assume that  $\psi$  is one-to-one (and hence, invertible).

Which of the following is equal to  $\theta^*$ ?

$(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi^{-1}(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi^{-1}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2, \dots, \frac{1}{n} \sum_{i=1}^n X_i^d\right)$



Which of the following is the method of moments estimator for  $\theta^*$ ?

$(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi^{-1}(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$

$\psi^{-1}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2, \dots, \frac{1}{n} \sum_{i=1}^n X_i^d\right)$



### Solution:

Observe that  $\psi(\theta^*) = (m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$  by definition of  $\psi$ . Since  $\psi$  is invertible, then we know that  $\psi^{-1}(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*)) = \theta^*$ . Hence,  $\psi^{-1}(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$  is the correct response to the first question.

The remaining choices are incorrect.

- $(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$  is the list of moments of  $\mathbf{P}_{\theta^*}$ , not the parameter  $\theta^*$  itself.
- $\psi(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$  is incorrect, as this is really  $\psi^2(\theta^*)$ , which is not necessarily  $\theta^*$ .
- This is the method of moments estimator, not the true parameter  $\theta^*$ .

The method of moments estimator is given by

$$\hat{\theta}_n^{\text{MM}} = \psi^{-1}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2, \dots, \frac{1}{n} \sum_{i=1}^n X_i^d\right),$$

**Remark:** The above expression is consistent with the procedure we followed in the previous problems that we used to construct the method of moments estimator for a Gaussian statistical model with unknown mean and variance. Namely, we expressed the true parameters in terms of the true moments, and then plugged in the sample means into that expression. Informally, since we expect the sample means to give a good approximation for the true moments, plugging in the sample moments into the expression for the true parameters (in terms of the moments) should also give a good approximation for the true parameters. This is the strategy of the method of moments, and in general, the strategy of replacing expectations with averages is a recurring theme in statistics and in this course.

## Applying the Method of Moments to a Gaussian Statistical Model

1/2 points (graded)

We let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu^*, (\sigma^*)^2)$  and consider the associated statistical model  $(\mathbb{R}, \{N(\mu, \sigma^2)\}_{\mu \in \mathbb{R}, \sigma > 0})$ . Let

$$\begin{aligned}\psi : \mathbb{R} \times (0, \infty) &\rightarrow \mathbb{R}^2 \\ (\mu, \sigma) &\mapsto (m_1(\mu, \sigma), m_2(\mu, \sigma)).\end{aligned}$$

denote the moments map considered in the previous problem, where  $m_k(\mu, \sigma)$  denotes the  $k$ -th moment of the distribution  $N(\mu, \sigma^2)$ .

To answer the next question, you should recall:

1. your result on writing  $\mu$  and  $\sigma$  in terms of  $m_1$  and  $m_2$  obtained in the previous problem, and
2. the estimators  $\widehat{m_1}$  and  $\widehat{m_2}$  (the sample moments) for the true moments  $m_1$  and  $m_2$ .

Suppose we observe the data-set  $X_1 = 0.5, X_2 = 1.8, X_3 = -2.3, X_4 = 0.9$ .

What is the method of moments estimator  $\widehat{\mu}^{\text{MM}}$  for  $\mu^*$  evaluated on this data-set? (You are encouraged to use whatever computational tools may be helpful.)

$$\widehat{\mu}^{\text{MM}}(0.5, 1.8, -2.3, 0.9) = \boxed{0.225} \quad \checkmark \text{ Answer: } 0.225$$

What is the method of moments estimator  $\widehat{\sigma}^{\text{MM}}$  for  $\sigma^*$  evaluated on this data-set? (You are encouraged to use whatever computational tools may be helpful.)

$$\widehat{\sigma}^{\text{MM}}(0.5, 1.8, -2.3, 0.9) = \boxed{0} \quad \times \text{ Answer: } 1.532$$

**Solution:**

Since we computed  $\psi^{-1}$  explicitly in the previous problem, we can apply the method of moments to estimate the true parameters  $\mu^*$  and  $\sigma^*$ . Let

$$\widehat{m}_1 = \frac{1}{n} \sum_{i=1}^n X_i, \quad \widehat{m}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$$

denote the first and second sample moments, respectively. Then the **method of moments estimator** is defined to be

$$(\widehat{\mu}_n^{\text{MM}}, \widehat{\sigma}_n^{\text{MM}}) := \psi^{-1} (\widehat{m}_1, \widehat{m}_2) = (\widehat{m}_1, \sqrt{\widehat{m}_2 - \widehat{m}_1^2}).$$

Using a calculator (or other computational software) we can compute,

$$\widehat{m}_1 (0.5, 1.8, -2.3, 0.9) = \frac{0.5 + 1.8 - 2.3 + 0.9}{4} \approx 0.225$$

and

$$\widehat{m}_2 (0.5, 1.8, -2.3, 0.9) = \frac{(0.5)^2 + (1.8)^2 + (-2.3)^2 + (0.9)^2}{4} \approx 2.3975.$$

Applying the method of moments,

$$\widehat{\mu}_n^{\text{MM}} (0.5, 1.8, -2.3, 0.9) = \widehat{m}_1 (0.5, 1.8, -2.3, 0.9) \approx 0.225.$$

and

$$\widehat{\sigma}_n^{\text{MM}} (0.5, 1.8, -2.3, 0.9) = \sqrt{\widehat{m}_2 - (\widehat{m}_1)^2} \approx \sqrt{2.3975 - (0.225)^2} \approx 1.532.$$