

## CHAPTER 3

### Estimators

- Statistical model = family  $P_\theta$

Data generated from one of the distributions in the model

Parametric model :  $\{P_\theta : \theta \in \Theta\}$   
if  $\Theta \subseteq \mathbb{R}^d$

Estimation = process of determining the parameter  $\theta$  that gives the best fitting model.

- Often we want to estimate a function  $g(\theta)$  given the data  $x$

ESTIMATOR or STATISTIC for  $g(\theta)$

is a random vector  $T(x)$   
that depends ONLY on

$$\underline{x} = (X_1, \dots, X_n)$$

ESTIMATE

is the observed value  $T(x)$

Both are indicated also by  $\hat{\theta}$

METHODS of Estimation

- maximum likelihood
- method of moments
- Bayes method

## 3.2 Mean Square Error

A good estimator must be "close" to the estimand  $g(\theta)$

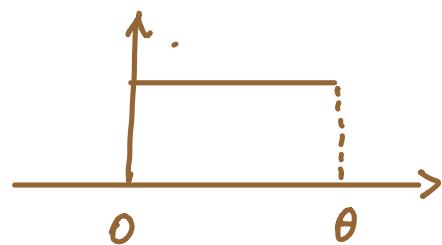
To measure the distance we consider the distribution of the distance

$$\| T(X) - g(\theta) \|^2$$

under the assumption the  $\theta$  is the true parameter

### Example 3.2.

Let  $X = (X_1, \dots, X_n) \stackrel{iid}{\sim} \text{Unif}(0, \theta)$



$$E(X_i) = \frac{\theta}{2}$$

Idea : Estimate  $\frac{\theta}{2}$  by  $\bar{X}_m = \frac{\sum X_i}{n}$

or Estimate  $\theta$  by  $2\bar{X}_m$

Justified by the Law of large numbers

$$\bar{X}_n \xrightarrow{P} \mu = \frac{\theta}{2} \quad \text{for } n \rightarrow \infty$$

For instance,

$$x = (3.03, 2.7, 7, 1.59, 5.04, 5.92, \\ 9.82, 1.11, 4.26, 6.96) \quad n = 10$$

$$2\bar{x} = 9.49$$

Note: this underestimates  $\theta$

because  $9.49 < x_7 = 9.82$

Idea 2: take  $\hat{\theta} = X_{(n)} = \max_i x_i$

But, again  $X_{(n)} < \theta$ .

We could correct  $X_{(n)} \rightarrow \frac{n+2}{n+1} X_{(n)}$

$\Rightarrow$  which estimator is best?

IDEA: the best estimator is the one with a sampling distribution most concentrated on  $\theta$  "in the long run"; or that is best on average

See Notebook

Note: our simulation shows that if  $\theta = 1$   $T_1 = \frac{n+2}{n+1} X_{(n)}$  is better than  $T_2 = 2 \bar{X}_n$

Mean square error

of an estimator  $T$  for  $g(\theta)$  is

$$\text{MSE}(\theta; T) = E_\theta \|T - g(\theta)\|^2$$

this is a function of  $\theta$ .

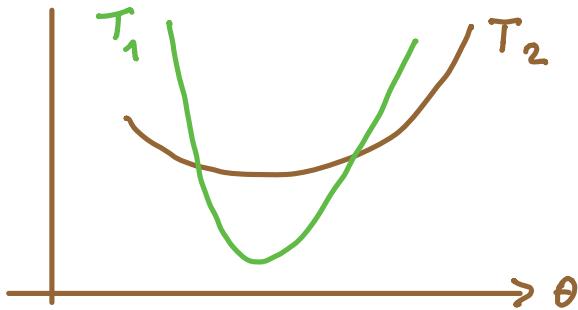
Note : the subscripted  $E_{\theta}$  is essential. This means that we compute the expected square deviation of  $T$  from  $g(\theta)$  under the assumption that  $\theta$  is the true value of the parameter.

We prefer an estimator with MSE small for ALL values of  $\theta$ .

If:  $MSE[\theta, T_1] \leq MSE[\theta, T_2]$  for all  $\theta$  with a strict inequality for at least one value of  $\theta$ ,

we prefer  $T_1$  and we say that  $T_2$  is inadmissible

- Sometimes this does not happen



Result: the MSE can be decomposed

$$\text{MSE}(\theta, T) = \text{var}_{\theta} T + [E_{\theta} T - g(\theta)]^2$$

↓                              ↓  
 Variance                      bias  
 of the estimator

Proof:  $E_{\theta} (T - \theta)^2 = E_{\theta} U^2$   
 (scalar)

We know that

$$\text{var}_{\theta}(U) = E_{\theta} U^2 - (E_{\theta} U)^2$$

$$\therefore E_{\theta} U^2 = \text{var}_{\theta} U + (E_{\theta} U)^2$$

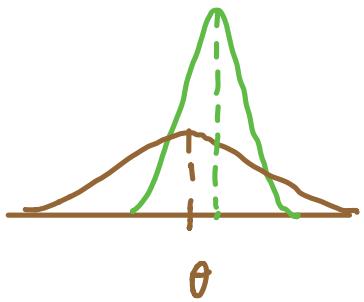
Thus:

$$\begin{aligned} E_{\theta} (T - \theta)^2 &= \text{var}_{\theta} (T - \theta) + [E_{\theta} (T - \theta)]^2 \\ &= \text{var}_{\theta} T + [E_{\theta} T - \theta]^2 \end{aligned}$$

### Unbiased estimator

$T$  is unbiased for  $g(\theta)$  if

$$E_{\theta} T = g(\theta) \text{ for all } \theta \in \mathbb{R}$$



While unbiased estimators look very desirable, there is a trade-off between variance and bias.

### Standard error

The standard deviation of an estimator  $T$ :  $\sigma_{\theta}(T) = \sqrt{\text{var}_{\theta} T}$  is called the standard error

The standard error gives an idea of the quality of an estimate.

Example:  $(x_1, \dots, x_n) \stackrel{iid}{\sim} N(\theta, \sigma_0^2)$

$\sigma_0^2$  = Known variance

$\bar{X}$  is an unbiased estimator of  $\theta$  (the mean) because

$$E_{\theta}(\bar{X}) = \frac{1}{n} E_{\theta} \sum_{i=1}^n X_i = \frac{1}{n} \sum_i E X_i = \frac{n\theta}{n} = \theta$$

The standard error is

$$\begin{aligned}\sigma_{\theta}(\bar{X}) &= \sqrt{\text{var}_{\theta} \bar{X}} = \sqrt{\frac{1}{n^2} \text{var}(\sum_i X_i)} \\ &= \sqrt{\frac{\sigma_0^2}{n}} = \frac{\sigma_0}{\sqrt{n}}\end{aligned}$$

This is the variability of  $\bar{X}$  around  $\theta$  in repeated sampling.

### Example 3.6

1)  $(X_1, \dots, X_n) \stackrel{iid}{\sim} U(0, \theta)$

$2\bar{X}$  is unbiased for  $\theta$

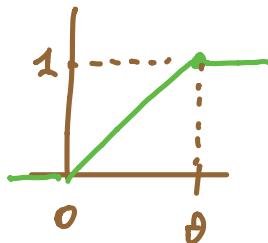
$$\begin{aligned}\rightarrow E_\theta(2\bar{X}) &= 2E_\theta(\bar{X}) = \frac{2}{n} \sum_i E_\theta X_i \\ &= \frac{2}{n} \sum_{i=1}^n \frac{\theta}{2} = \theta.\end{aligned}$$

$$\begin{aligned}\rightarrow \text{MSE}(\theta, 2\bar{X}) &= \text{var}_\theta(2\bar{X}) = 4 \text{var}_\theta(\bar{X}) \\ &= 4 \cdot \frac{\text{var}(X_1)}{n}\end{aligned}$$

$$\text{As } \text{var} X_1 = \frac{\theta^2}{12} = \frac{4\theta^2}{12n} = \frac{\theta^2}{3n}.$$

2) Is  $X_{(n)}$  unbiased? The density of  $X_{(n)}$  can be found in this way

$$\begin{aligned}F_{(n)}(x) &= P(X_{(n)} \leq x) \\ &= P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x)\end{aligned}$$



$$\begin{aligned}
 &= \prod_{i=1}^n P(X_i \leq x) \\
 &= \prod_{i=1}^n \frac{x}{\theta} = \frac{x^n}{\theta^n}
 \end{aligned}$$

So that

$$f_{(n)}(x) = \frac{d}{dx} \frac{x^n}{\theta^n} = \frac{n x^{n-1}}{\theta^n}.$$

The Expected value is

$$\begin{aligned}
 \int_0^\theta x \cdot \frac{n x^{n-1}}{\theta^n} dx &= \frac{n}{\theta^n} \cdot \frac{x^{n+1}}{n+1} \Big|_0^\theta \\
 &= \frac{n}{n+1} \theta.
 \end{aligned}$$

→ So  $X_{(n)}$  is biased

→ The MSE is  $\text{var}_\theta X_{(n)} + \left(\frac{n}{n+1} \theta - \theta\right)^2$

The variance is

$$\text{var}_\theta(X_{(n)}) = E_\theta X_{(n)}^2 - [E_\theta X_{(n)}]^2$$

$$= \int_0^\theta x^2 \cdot \frac{n x^{n-1}}{\theta^n} dx - \left[ \frac{n \theta}{n+1} \right]^2$$

$$= \frac{n}{\theta^n} \int_0^\theta x^{n+1} dx - \frac{n^2 \theta^2}{(n+1)^2}$$

$$= \frac{n}{\theta^n} \frac{x^{n+2}}{n+2} \Big|_0^\theta - \frac{n^2 \theta^2}{(n+1)^2}$$

$$= \frac{n \theta^2}{n+2} - \frac{n^2 \theta^2}{(n+1)^2}$$

$$= \theta^2 \frac{n}{(n+2)(n+1)^2} .$$

And finally :

$$\begin{aligned} \text{MSE}(\theta, X_{(n)}) &= \theta^2 \frac{n}{(n+2)(n+1)^2} + \left( \frac{\theta n}{n+1} - \theta \right)^2 \\ &= \frac{2\theta^2}{(n+2)(n+1)}. \end{aligned}$$

Compare

$$\text{MSE}(\theta, 2\bar{X})$$

$$\frac{\theta^2}{3n}$$

$$\text{MSE}(\theta, X_{(n)})$$

$$\frac{2\theta^2}{(n+2)(n+1)}$$

See notebook

Other estimators

$$\rightarrow \frac{n+1}{n} X_{(n)} \quad \text{unbiased}$$

$$\rightarrow \frac{n+2}{n+1} X_{(n)} \quad \text{better}$$

$$\text{MSE} = \frac{\theta^2}{(n+1)^2}$$

## Example 3.7

In general  $(X_1, \dots, X_n) \stackrel{iid}{\sim} F$

with mean and variance

$$\mu = E_F(X_i) \quad \sigma^2 = \text{var}_F(X_i)$$

: a nonparametric model

→ Two estimators of  $\mu$  and  $\sigma^2$

$$\bar{X} = \frac{1}{n} \sum_i X_i \quad S_x^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X})^2$$

→  $\bar{X}$  is unbiased for  $\mu$

$$\rightarrow \text{MSE}(F; \bar{X}) = \text{var}_F(\bar{X})$$

$$= \frac{1}{n^2} \sum_i \text{var}_F(X_i) = \frac{\sigma^2}{n}$$

The precision of the estimator  $\bar{X}$  increases by a factor of  $\sqrt{n}$

→ The sample variance is unbiased for  $\sigma^2$ . Proof:

$$s_x^2 = \frac{1}{n-1} \left[ \sum_i (x_i - \mu)^2 - n (\bar{x} - \mu)^2 \right]$$

→ ESERCIZIO

$$E_F(s_x^2) = \frac{1}{n-1} \left[ \sum_i E_F(x_i - \mu)^2 - n E_F(\bar{x} - \mu)^2 \right]$$

$$= \frac{1}{n-1} \left[ \sum_{i=1}^n \sigma^2 - n \text{var}_F(\bar{X}) \right]$$

$$= \frac{1}{n-1} \left( n\sigma^2 - n \frac{\sigma^2}{n} \right) = \sigma^2.$$

Note that unbiasedness is not preserved under transformation

$$E_F(s_x) \neq \sigma$$

$$E_F(\bar{X}^2) \neq \mu^2$$

### 3.3. Maximum likelihood

Is the most common method for finding estimators.

- Is based on a function called the likelihood function

Given the data  $\underline{x} = (x_1, \dots, x_n)$

and a parametric model  $P_\theta \leftrightarrow p_\theta$

the likelihood is the density

$$p_\theta(\underline{x})$$

of the observation, as a function of  $\theta \in \mathbb{R}$ .

### Example 3.9 - Binomial

Data after Tossing a biased coin 10 times:

( H T T H T T H T T T )

so we get 3 heads.

Let  $x = 3$  and the model is

$$X \sim \text{Bin}(10, p)$$

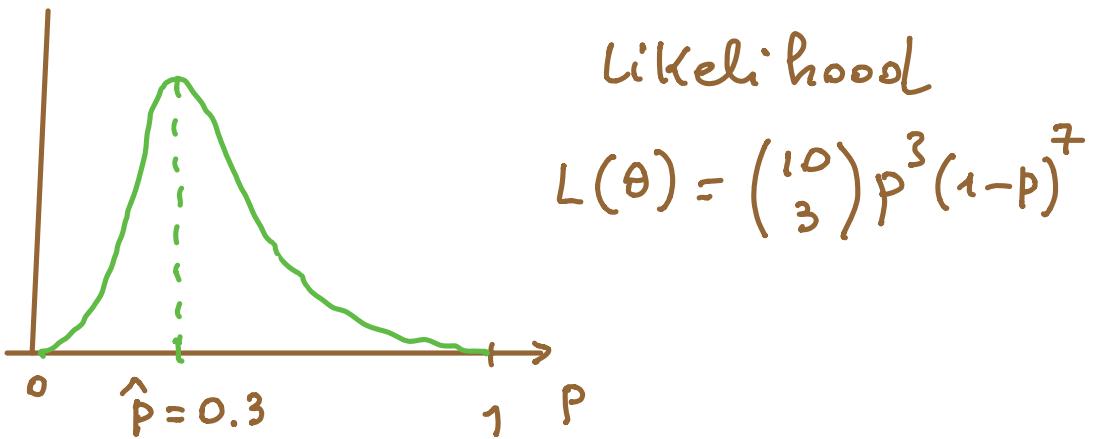
with  $p$  = prob. of success unknown.

What's the probability of getting 3 heads on 10 tosses?

$$P_p(x=3) = \binom{10}{3} p^3 (1-p)^7$$

is a function of  $p$

See the notebook



$\hat{p} = 0.3$  is the max probability  
of getting 3 heads

If  $\underline{X}$  is a random vector with  
density\*  $p_\theta(\underline{x})$ ,  $\theta \in \Theta$   
the likelihood function is the  
function

$$\theta \mapsto L(\theta; \underline{x}) = p_\theta(\underline{x})$$

for  $\underline{x}$  fixed at the observed value

\* we talk of density both for   
continuous  
discrete

## i.i.d Sample

If  $\underline{x} = (x_1, \dots, x_n) \stackrel{i.i.d}{\sim} p_\theta(x)$

the likelihood is

$$L(\theta; \underline{x}) = \prod_{i=1}^n p_\theta(x_i)$$

Here  $p_\theta(x_i)$  is

the marginal density of  $x_i$ .

Maximum likelihood estimate (MLE)

Is the value  $\hat{\theta}(\underline{x}) \in \Theta$  that maximizes  $L(\theta; \underline{x})$ .

The maximum likelihood is an intuitively reasonable principle for finding estimators.  
MLE are not necessarily the best ones.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta; \underline{x}) = \underset{\theta}{\operatorname{argmax}} \log L(\theta; \underline{x})$$

↓  
likelihood      log likelihood

If  $L(\theta)$  is differentiable in  $\Theta \subset \mathbb{R}^K$   
 and it takes the maximum in an  
 interior point of  $\Theta$  then

$$\frac{\partial}{\partial \theta_j} \log L(\underline{\theta}; \underline{x}) \Big|_{\underline{\theta} = \hat{\theta}} = 0 \quad j = 1, \dots, K$$

System of likelihood equations

NB - cannot be solved always explicitly

→ Check the form of the likelihood

To verify if a solution is the max.

- If the solution is the maximum the 2nd derivative is negative
- All the eigenvalues of the Hessian are negative

$$\text{Score function} = \frac{\partial}{\partial \theta_j} \log L(\underline{\theta}, \underline{x})$$

an i.i.d  
the likelihood equation for sample

$$\sum_{i=1}^n \frac{\partial}{\partial \theta_j} \log P_{\theta}(x_i) = 0 \quad j=1, \dots, k$$

### Example 3.12 - Exponential

$$(x_1, \dots, x_n) \stackrel{\text{i.i.d}}{\sim} \text{Exp}(\lambda)$$

$$\lambda = \text{rate} \quad E_{\lambda}(x_i) = \frac{1}{\lambda} > 0$$

$$\begin{aligned}
 \frac{d}{d\lambda} \log p_\theta(x_i) &= \frac{d}{d\lambda} \log (\lambda e^{-\lambda x_i}) \\
 &= \frac{d}{d\lambda} (\log \lambda - \lambda x_i) \\
 &= \frac{1}{\lambda} - x_i
 \end{aligned}$$

Likelihood equations

$$\sum_{i=1}^n \left( \frac{1}{\lambda} - x_i \right) = \frac{n}{\lambda} - \sum x_i = 0$$

$$1 \text{ solution} = \frac{1}{\bar{x}}$$

2nd derivative

$$\frac{d}{d\lambda} \left( \frac{1}{\lambda} - x_i \right) = -\frac{1}{\lambda^2}$$

$$\text{So } \frac{d^2}{d\lambda^2} \log L(\underline{x}; \lambda) = -\frac{n}{\lambda^2} < 0, \forall \lambda$$

$$\text{so } \hat{\lambda}_{ML} = \frac{1}{\bar{x}} . \quad //$$

### Example 3.13 - Binomial .

$x = \# \text{ successes}$        $n = \# \text{ trials}$

$$L(p; x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\log L(p; x) = \ell(p)$$

$$\log \binom{n}{x} + x \log p + (n-x) \log (1-p)$$

If  $0 < x < n$ ,

$$\lim_{p \rightarrow 0} \ell(p) = -\infty \quad \lim_{p \rightarrow 1} \ell(p) = -\infty$$

so the maximum is in  $(0, 1)$

Likelihood equation

$$\frac{d}{dp} \log L(p; x) = \frac{x}{p} - \frac{n-x}{1-p} = 0$$

$$\text{Single solution } \hat{p} = \frac{x}{n} = \frac{\# \text{ successes}}{\# \text{ trials}}$$

If

$$x=0, \ell(p) = n \log(1-p) \Rightarrow \hat{p} = 0$$

$$x=n, \ell(p) = n \log p \Rightarrow \hat{p} = 1$$

Both can be written as  $\hat{p} = \frac{x}{n}$ .

Example 3.14 - Normal distribution

$$(x_1, \dots, x_n) \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$$

$$(\mu, \sigma^2) \mapsto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x_i - \mu)^2\right]$$

$$\Theta = \mathbb{R} \times (0, \infty)$$

log-likelihood

$$(\mu, \sigma^2) \mapsto -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

## Likelihood equations

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\mu = \frac{\sum x_i}{n} = \bar{x}$$

$$\sigma^2 = \frac{\sum_i (x_i - \bar{x})^2}{n}$$

for  $\mu = \bar{x}$  likelihood has a maximum  
for every  $\sigma > 0$

$$\text{Hessian} = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \mu^2} & \frac{\partial^2 \ell}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 \ell}{\partial (\sigma^2)^2} \end{pmatrix}$$

$$H = \begin{pmatrix} -\frac{n}{\sigma^2} & -\frac{1}{\sigma^4} \sum_i (x_i - \mu) \\ \frac{1}{\sigma^4} \sum_i (x_i - \mu) & \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_i (x_i - \mu)^2 \end{pmatrix}$$

Substitute  $\mu \rightarrow \bar{x}$  and  $\sigma^2 \rightarrow \hat{\sigma}^2$

$$H = \begin{pmatrix} \frac{n}{\hat{\sigma}^2} & 0 \\ 0 & \frac{n}{2\hat{\sigma}^4} - \frac{n\hat{\sigma}^2}{\hat{\sigma}^6} \end{pmatrix}.$$

$$= \begin{pmatrix} -\frac{n}{\hat{\sigma}^2} & 0 \\ 0 & -\frac{n}{2\hat{\sigma}^4} \end{pmatrix}.$$

Both eigenvalues are negative so

$(\bar{x}, \frac{n-1}{n} S_x^2)$  is the MLE //

## Notes

1) ... MLE can be taken outside  
the interior of  $\Theta$

- 2) ... likelihood equations do not hold
- 3) ... likelihood not everywhere differentiable
- 4) ... MLE has several local maxima.

## Example 3.15 - Uniform

$(x_1, \dots, x_n) \stackrel{\text{iid}}{\sim}$  Uniform on  $[0, \theta]$

$\underline{x} = (x_1, \dots, x_n) \rightarrow x_1 \leq \theta, x_2 \leq \theta, \dots, x_n \leq \theta$

$$\Rightarrow \boxed{x_{(n)} \leq \theta}$$

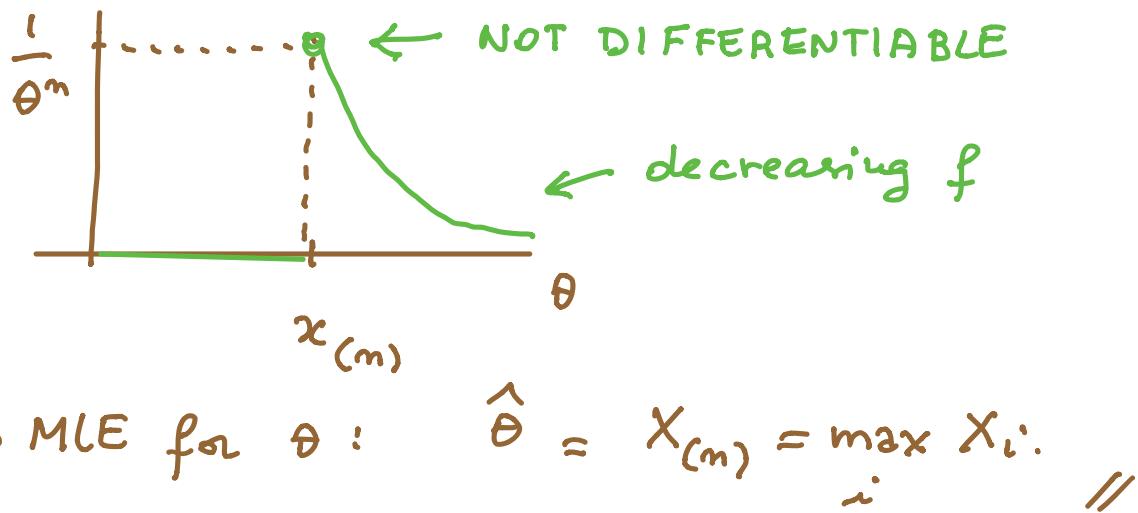
$$L(\theta; \underline{x}) = \prod_{i=1}^n p_\theta(x_i)$$

$$\text{where } p_\theta(x_i) = \frac{1}{\theta} \cdot 1_{0 \leq x_i \leq \theta}$$

Therefore:

$$L(\theta; \tilde{x}) = \left(\frac{1}{\theta}\right)^n \cdot 1_{x_1 > 0} \cdot 1_{x_{(n)} \leq \theta}$$

$$= \begin{cases} \left(\frac{1}{\theta}\right)^n & \text{if } \theta \geq x_{(n)} \\ 0 & \text{if } \theta < x_{(n)} \end{cases}$$



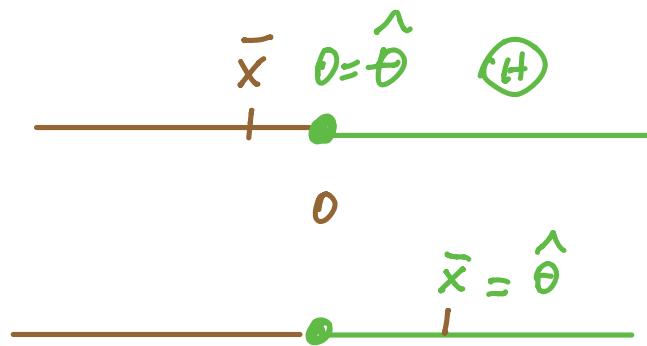
### 3.16 - Normal distrib. with restrictions

$$(X_1, \dots, X_n) \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2=1)$$

under the restriction  $\mu \geq 0$

$$\Theta = [0, \infty)$$

Given  $\underline{x} = (x_1, \dots, x_n)$  it can happen that the unrestricted MLE  $\hat{\theta} = \bar{x}$  is negative while  $\mu > 0$ . So this is not the restricted MLE



$$\text{MLE} = \max(0, \bar{x}).$$

The statistical model  
the MLE } depend on the  
form of  $\Theta$ .

## Equivariance of the MLE

What's the MLE of  $g(\theta)$ ?

Let  $g: \Theta \rightarrow H$  a bijective fun.

We parametrize the model with

$$\eta = g(\theta), \quad \eta \in H$$

$\hat{\theta}$  is MLE of  $\theta \Rightarrow \hat{\eta} = g(\hat{\theta})$  is MLE of  $\eta$

For any function  $g$

Def.

$g(\hat{\theta})$  is the MLE of  $g(\theta)$

### Example 3.17 Exponential (again)

$(X_1, \dots, X_n) \stackrel{i.i.d}{\sim} \text{Exp}(\lambda)$

MLE of  $\mu = \frac{1}{\lambda}$

$$\hat{\lambda}_{ML} = \frac{1}{\bar{x}} \Rightarrow \hat{\mu} = \bar{x}. \quad //$$

### Example 3.18 - Gamma

$(X_1, \dots, X_n) \stackrel{i.i.d}{\sim} \text{Gamma}(\alpha, \lambda)$

$$P_{\alpha, \lambda}(x) = \frac{x^{\alpha-1}}{\Gamma(\alpha)} \lambda^\alpha e^{-\lambda x} \quad \begin{matrix} \lambda = \text{inverse} \\ \text{scale} \end{matrix}$$

$\alpha = \text{shape}$



Gamma function

$$\Gamma(\alpha) = \int_0^\infty s^{\alpha-1} e^{-s} ds$$

$$\log P_{\alpha, \lambda}(x_i) = (\alpha - 1) \log x_i + \alpha \log \lambda -$$

$$- \lambda x_i - \log \Gamma(\alpha)$$

$$l(\alpha, \lambda) = (\alpha - 1) \sum_i \log x_i + n \alpha \log \lambda -$$

$$- \lambda \sum_i x_i - n \log \Gamma(\alpha)$$

Parameter

$$\theta = (\alpha, \lambda) \in \Theta$$

$$\Theta = [0, \infty) \times [0, \infty)$$

$$\frac{\partial l}{\partial \alpha} = \sum_i \log x_i + n \log \lambda - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} = 0$$

$$\frac{\partial l}{\partial \lambda} = \frac{n \alpha}{\lambda} - \sum_i x_i = 0$$

From the 2nd equation

$$\lambda = \frac{\alpha}{\bar{x}}$$

Substituting into the 1st eq.  
we solve this equation

$$\sum_i i \log x_i + n \log\left(\frac{\alpha}{\bar{x}}\right) - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} = 0$$

numerically.

---

See the notebook

---

### Example 3.19 - Counting bacteria

With Petri dishes  
it is possible to  
see by naked eye  
colonies of bacteria.



Petri dish  
(capsula di Petri)

Assumption : # colonies  $\sim$  Poisson( $\mu$ )

$\frac{\mu}{\text{1 cc}}$   
1 ml

of contaminated  
water.

Dilute the water in 100 buckets  
so that

$$X_1, \dots, X_{100} \stackrel{iid}{\sim} \text{Poisson}(\mu/100)$$

and record

$$Y_1, \dots, Y_{100} \quad \text{where } Y_i = \begin{cases} 0 & \text{if no colonies} \\ 1 & \text{otherwise} \end{cases}$$

Estimate  $\mu$  from  $(Y_1, \dots, Y_{100})$

$$Y_i \stackrel{iid}{\sim} \text{Bernoulli}(P = P(X_i = 1))$$

$$\text{so that } p = 1 - \frac{e^{-\lambda} \lambda^0}{0!} = 1 - e^{-\frac{\mu}{100}}.$$

$$\text{Therefore, } \hat{p} = \frac{\sum Y_i}{100}$$

$$\text{and } \left(1 - e^{-\frac{\mu}{100}}\right) = \hat{p}$$

$$\hat{\mu} = -100 \log(1 - \sum Y_i/100)$$

## 3.4 – Method of moments estimator

---

- Simpler alternative to maximum likelihood
- Requires only the theoretical form of the moments.
- Is based on imposing the matching of theoretical and sample moments

$$\begin{array}{ccc} & \swarrow & \downarrow \\ E_\theta(x^j) & & \overline{x^j} = \frac{1}{n} \sum_{i=1}^n x_i^j \end{array}$$

Method of moments estimator  
is the value  $\hat{\theta}$  where

$$E_{\hat{\theta}}(X^j) = \bar{X}^j$$

take  $j$   
as small  
as possible

The MME for  $g(\theta)$  is taken to be  $g(\hat{\theta})$

### Example 3.27 - Exponential

$$x_i \sim \lambda e^{-\lambda x} \quad \lambda > 0$$

$$E(x_i) = \frac{1}{\lambda} \Rightarrow \frac{1}{\lambda} = \bar{x} \Rightarrow \hat{\lambda}_{MM} = \frac{1}{\bar{x}}.$$

$$\hat{\lambda}_{MM} = \hat{\lambda}_{ML}.$$

### Example 3.28 - Uniform

$$x_i \sim \frac{1}{\theta}$$

$$E(x_i) = \theta/2 \quad \bar{x} = \theta/2 \Rightarrow \hat{\theta}_{MM} = 2\bar{x}$$

$$\hat{\theta}_{MM} \neq \hat{\theta}_{ML}.$$

### Example 3.29 - Normal

$$x_i \sim N(\mu, \sigma^2)$$

$$E(x_i) = \mu \quad E(x_i^2) = \sigma^2 + \mu^2$$

$$\begin{cases} \mu = \bar{x} \\ \sigma^2 + \mu^2 = \bar{x^2} \end{cases}$$

$$\hat{\mu}_{MM} = \bar{x} \quad \hat{\sigma}_{MM}^2 = \bar{x^2} - \bar{x}^2 \\ = \frac{1}{n} \sum_i (x_i - \bar{x})^2.$$

$$\hat{\sigma}_{MM}^2 = \hat{\sigma}_{ML}^2.$$

//

### Example 3.30 - Gamma distribution

$$x_i \sim \text{Gamma}(\alpha, \lambda)$$

$$E(x_i) = \frac{\alpha}{\lambda} \quad E(x_i^2) = \frac{\alpha}{\lambda^2} + \frac{\alpha^2}{\lambda^2} \\ = \frac{\alpha}{\lambda^2} (1 + \alpha)$$

$$\begin{cases} \bar{X} = \frac{\alpha}{\lambda} \\ \bar{X^2} = \frac{\alpha}{\lambda^2} (1+\alpha) \end{cases} \rightarrow \begin{cases} \hat{\lambda}_{MM} = \frac{(\bar{X})^2}{\bar{X^2} - (\bar{X})^2} \\ \hat{\alpha}_{MM} = \frac{\bar{X}}{\bar{X^2} - (\bar{X})^2} \end{cases}$$

are different from MLE and in  
closed form

### 3.5 Bayes estimators

A different approach to inference

- $\theta$  is fixed and unknown
- $\theta$  is a random variable described by a density  $\pi(\theta)$  representing the information prior to observing the data.

I) The Bayesian approach starts by defining the prior distribution  $\pi(\theta)$  on in addition to the statistical model.

2) the Bayesian approach uses data and Bayes formula to transform the prior into a posterior  $\pi(\theta | \underline{x})$

Bayes risk of  $T(x)$  given  $\pi(\theta)$

is the weighted average of the MSE

$$R(\pi, T) = \int E_{\theta} (T - g(\theta))^2 \pi(\theta) d\theta$$

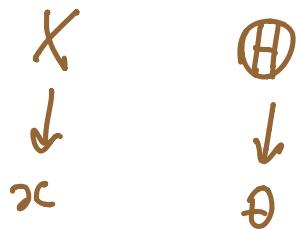
with weight  $\pi(\theta)$ .

— Bayes estimator —

is the estimator  $T$  that minimizes  $R(\pi, T)$  over all estimators.

## Bayes formula

Discrete form.



$$P_{\theta}(x) = P(x | \Theta = \theta) \text{ likelihood}$$
$$\pi(\theta) \text{ prior}$$

Posterior

$$P(\Theta = \theta | X = x) =$$
$$= \frac{P(X = x | \Theta = \theta) P(\Theta = \theta)}{\sum_{\theta} P(X = x | \Theta = \theta) P(\Theta = \theta)}$$

## Continuous form

### Posterior

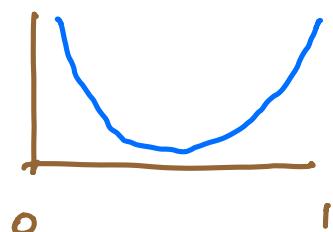
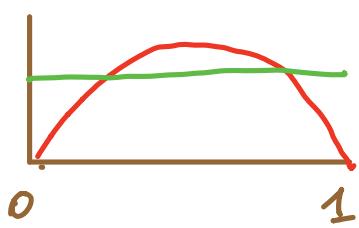
$$P_{\text{Posterior}}(\theta | x) = \frac{P(x | \theta) \pi(\theta)}{\int P(x | \theta) \pi(\theta)}$$

- The denominator is a normalizing constant.

### Example 3.38 - Binomial

$$X \sim \text{Bin}(n, p)$$

Prior on  $p \in (0, 1)$

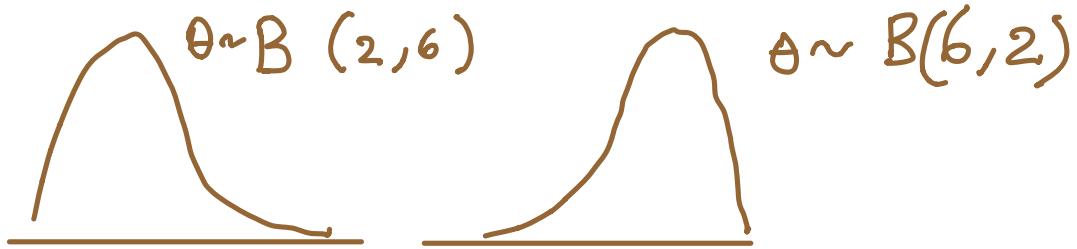


A useful family of priors is the Beta

## Beta density

See the notebook

Some properties  $\theta \sim \text{Beta}(\alpha, \beta)$



$$E(\theta) = \frac{\alpha}{\alpha + \beta}$$

$$\pi(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

where  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$

$$= \int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta$$

The nice thing of the Beta  
is that  
The posterior is also Beta!

$$p(\theta | x) = \frac{\binom{n}{x} \theta^x (1-\theta)^{n-x} \pi(\theta)}{\int_0^1 \binom{n}{x} \theta^x (1-\theta)^{n-x} \pi(\theta) d\theta}$$

$$\text{Numerator} = \frac{\theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}{C(x, \alpha, \beta)}$$

Normalization constant ↗

$$\text{must be } = \int_0^1 \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta$$

Thus  $\Theta | x=x \sim \text{Beta}(\alpha+x, \beta+n-x)$

## Example / Fig. 3.7

see → Notebook 3

### Comments

•	Prior	Model	Posterior
	Beta	Binomial	Beta

The Beta family is said conjugate for the Binomial

- The Bayes estimator is

$T(x) = \frac{X + \alpha}{n + \alpha + \beta}$  Expected value of the posterior. (Theorem 3.36)

$$= \frac{X + \alpha}{n + \alpha + \beta}$$

- For each prior  $B(\alpha, \beta)$  we have a Bayes estimator.
- The MLE  $X/n$  is Not a Bayes est.

## Frequentist evaluation of the Bayes estimator

Calculation of the

$$MSE(\theta; T_{\alpha, \beta}) =$$

$$\text{var}_{\theta} \left( \frac{x + \alpha}{n + \alpha + \beta} \right) + \left[ E_{\theta} \left( \frac{x + \alpha}{n + \alpha + \beta} \right) - \theta \right]^2$$

$$= \frac{\text{Var}_\theta(x)}{(n + \alpha + \beta)^2} + \left[ \frac{E_{\theta}(x) + \alpha}{n + \alpha + \beta} - \theta \right]^2$$

$$= \frac{n\theta(1-\theta)}{(n + \alpha + \beta)^2} + \left( \frac{n\theta + \alpha}{n + \alpha + \beta} - \theta \right)^2$$

$\rightarrow \underline{\text{See Notebook 3}}$

Theorem 3.36 → See Th. 3.33

The Bayes estimate for  $g(\theta)$  with respect to the prior  $\pi(\theta)$  is

$$\begin{aligned} T(x) &= E_{\Theta | X=x} (g(\theta)) \\ &= \int g(\theta) p(\theta | x) d\theta. \end{aligned}$$

Proof : First step : rewrite the Bayes risk  $E(MSE)$

$$MSE(\theta; T) =$$

$$= \begin{cases} E_\theta (T(x) - g(\theta))^2 & \text{(usual notation)} \\ E[(T(x) - g(\Theta))^2 | \Theta = \theta] & \text{(Bayesian notation)} \end{cases}$$

This is a conditional expectation

Reminder on cond. expectation

$$1) \quad E(z|y=y) = \int z p(z|y) dz.$$

this is a function of  $x$ .

$$2) \quad E(z) = \int E(z|y=y) p(y) dy$$

The Bayes risk is

$$R = \int E \left[ \underbrace{(T(x) - g(\theta))^2}_{z} \mid \underbrace{\oplus = \theta}_{y} \right] \pi(\theta) d\theta$$

$$= E(z) = E[(T(x) - g(\oplus))^2]$$

$$= \int E \left[ \underbrace{(T(x) - g(\oplus))^2}_{z} \mid \underbrace{X=x}_{y} \right] p_x(x) dx$$

To minimize the Risk is equivalent  
to minimize the integrand wrt  $x$ :

For every  $x$

$$\min_t E[(t - g(\theta))^2 | X=x] p_x(x)$$

$$\min_t E[(t - g(\theta))^2 | X=x] \quad (*)$$

We know that  $E(t - Y)^2 \rightarrow \min_t$

for  $t = E(Y)$ .

In this case  $Y \sim g(\theta) | X=x$

So that the Bayes estimate is

$$(*) \Leftrightarrow t = E(g(\theta) | X=x)$$

$$= \int g(\theta) p(y|x) d\theta.$$

//

## Interpretation of a Bayes estimate

The Bayes estimate is a weighted average of

$$\hat{\theta}_{ML} \quad \text{and} \quad E(\theta) = \theta_0$$

the Maximum likelihood

The mean of the prior

Binomial model

$$X \sim \text{Bin}(n, \theta)$$

$$\theta \sim \text{Beta}(\alpha_0, \beta_0)$$

$$\theta | X \sim \text{Beta}(x + \alpha_0, n - x + \beta_0)$$

$$\hat{\theta}_B = w \hat{\theta}_{ML} + (1-w) \theta_0$$

with  $w = n / (n + \alpha_0 + \beta_0)$ .

Proof:

$$\frac{n}{n + \alpha_0 + \beta_0} \cdot \frac{x}{n} + \frac{\alpha_0 + \beta_0}{n + \alpha_0 + \beta_0} \cdot \frac{\alpha_0}{\alpha_0 + \beta_0}$$

||

$$\hat{\theta}_B = \frac{x + \alpha_0}{n + \alpha_0 + \beta_0}$$

## Example 3.37 Exponential

$(x_1, \dots, x_n) \stackrel{iid}{\sim} \text{Exp}(\theta)$     $\theta$  rate

Assume that  
the prior is

$$\boxed{\theta \sim \text{Exp}(\lambda)} \quad \downarrow \text{known}$$

The posterior is

$$P(\theta | x) \propto p(x | \theta) \pi(\theta)$$

$$\hookrightarrow \prod_{i=1}^n \theta e^{-\theta x_i} \cdot \lambda e^{-\lambda \theta}$$

$$= \theta^{n-1} e^{-\theta \sum x_i} \lambda^n e^{-\lambda \theta} = \theta^n \lambda e^{-\theta(\sum x_i + \lambda)}$$

$$P(\theta | x) = \frac{\theta^n \lambda e^{-\theta(\sum x_i + \lambda)}}{C(x, \lambda)}$$

This is a Gamma density

$$\theta | x \sim \text{Gamma}(n+1, \lambda + \sum x_i)$$

The Bayes estimate is the expected value of this distribution.

Thus the Bayes estimate of  $\theta$  is :

$$T_\lambda(x) = \frac{n+1}{\lambda + \sum x_i};$$

→ The Bayes estimate of  $\theta^2$  is

The 2nd moment of the posterior  
 $=$  variance + mean<sup>2</sup>  
 $= \frac{n+1}{(\lambda + \sum x_i)^2} + \frac{(n+1)^2}{(\lambda + \sum x_i)^2}$   
 $= \frac{(n+1)(n+2)}{(\lambda + \sum x_i)^2}$ .

Note : • The Bayes estimator is not equivariant.

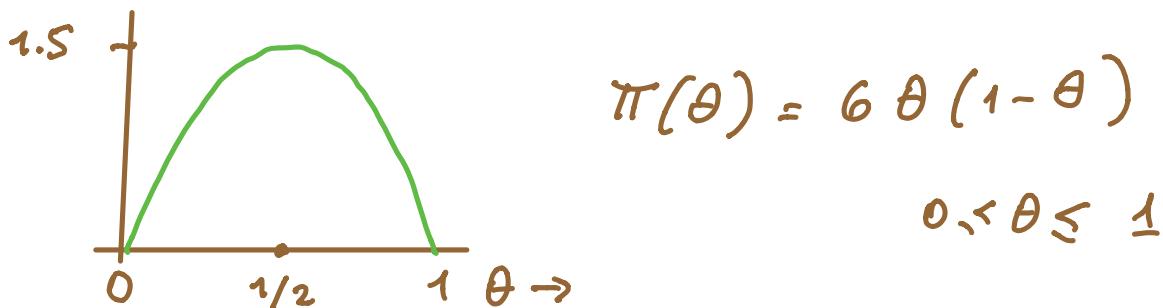
• For large  $n$   $T_\lambda(x) \rightarrow \frac{1}{\bar{x}} = \hat{\theta}_{ML}$

### Example 3.39 Geometric.

Geometric = # Trials until success

$$P_\theta(x) = (1-\theta)^{x-1} \theta \quad x = 1, 2, 3 \dots$$

Let  $\theta \sim \text{Beta}(\alpha=2, \beta=2)$ .



Posterior

$$\begin{aligned} P(\theta | \underline{x}) &= \frac{\prod_{i=1}^n (1-\theta)^{x_i-1} \theta}{C(\underline{x})} \cdot 6\theta(1-\theta) \\ &= \frac{\theta^{m+1} (1-\theta)^{n(\bar{x}-1)+1}}{C'(\underline{x})} \end{aligned}$$

This is a Beta  $(n+2, m(\bar{x}-1)+2)$

And the Bayes estimator is

$$T(x) = \frac{n+2}{m\bar{X}+4} .$$