

## Recap

### Method of moments to estimate $\theta$

•  $d = \#$  parameters in the model

• set up  $d$  equations:

$$\underbrace{E[X^K]}_{\substack{\uparrow \\ \text{involves } \theta}} = \underbrace{\frac{1}{n} \sum_{i=1}^n X_i^K}_{\substack{\uparrow \\ \text{involves data}}}, \quad K=1, 2, \dots, d$$

and solve for  $\theta$ . solution =  $\hat{\theta}_{MOM}$ .

8

# Method of Maximum Likelihood

$X_1, \dots, X_n \sim \text{i.i.d.}$  and  $\text{identically distributed}$ , follow  $f_\theta(x)$ .

Likelihood function of data: Joint pdf or pmf of sample

data considered as a function of  $\theta$  with data held fixed at the

observed values  $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ .

$$L(\theta) = L(\theta; \underbrace{x_1, x_2, \dots, x_n}_{\text{fixed}}) = \prod_{i=1}^n f_\theta(x_i) = f_\theta(x_1) \cdot f_\theta(x_2) \cdot \dots \cdot f_\theta(x_n)$$

$\uparrow$   $\uparrow$   
indep. indep.

a function of  $\theta$ .

- A function of  $\theta$  — the data are held fixed.

Maximum likelihood estimator (MLE) of  $\theta$ : The value  $\hat{\theta}$  of  $\theta$  that maximizes the likelihood function as a function of  $\theta$ .

- Can think of MLE as the value of  $\theta$  that is “most likely” to have led to the observed data.
- Essentially a calculus problem.

# How to find MLE?

Direct approach: Directly maximize the likelihood function.

**Ex:** Let  $X_1, X_2, \dots, X_n$  represent a random sample from a Uniform  $(0, \theta)$  distribution where  $\theta > 0$ . Find the MLE of  $\theta$ .

Recall:

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta}, & 0 \leq x \leq \theta \\ 0, & \text{o/w.} \end{cases}$$

MLE of  $\theta$ ?

$$\frac{\theta}{2} \uparrow \quad \theta = E[X] = \bar{X} \Rightarrow \theta = 2\bar{X} = \hat{\theta}_{\text{MLE}}.$$

(verify)

$\theta = \text{max. value in the population.}$

MLE of  $\theta$ :

~~$\theta = \text{max}\{x_1, \dots, x_n\}$~~  is a natural estimator of  $\theta$ .

Intuitively:

$\uparrow$   
will see that this is MLE.

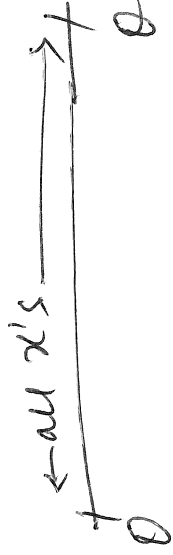
$$I(a \leq x \leq b) = \begin{cases} 1, & a \leq x \leq b \\ 0, & \text{o/w} \end{cases}$$

indicator function

$$L(\theta) = \frac{1}{\theta} I(0 \leq x_1 \leq \theta) \cdot \frac{1}{\theta} I(0 \leq x_2 \leq \theta) \cdot \dots \cdot \frac{1}{\theta} I(0 \leq x_n \leq \theta)$$

$f_{\theta}(x)$

$$f_{\theta}(x_n)$$



$$= \frac{1}{\theta^n}, \text{ provided}$$

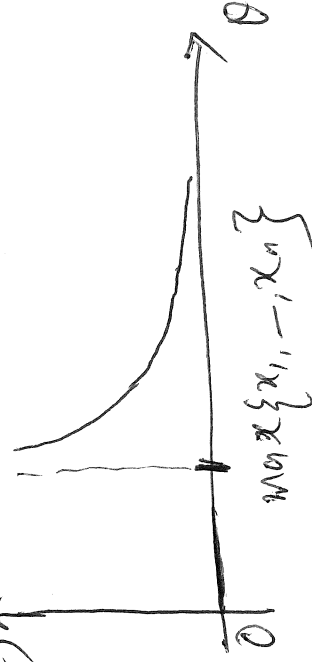
all the  $x_i$ 's are  $\leq \theta$

$$= \frac{1}{\theta^n} \text{ provided}$$

[No need to worry about  $x_i$ 's being negative because by def. all  $x_i$ 's  $\geq 0$ ].

$$= \frac{1}{\theta^n}, \text{ } \max\{x_1, \dots, x_n\} \leq \theta.$$

$$L(\theta)$$



The value of  $\theta$  that maximizes

$$L(\theta) \text{ is } \theta = \max\{x_1, \dots, x_n\}$$

$\hat{\theta}_{MLE}$   
Tdy.

Graph  $\Rightarrow$

Q. Have two estimators for  $\theta$  —  $\hat{\theta}_1$  and  $\hat{\theta}_2$ . Which one will you choose?

A. The one that is better in the sense of having smaller mean-squared error (MSE) defined as

$$MSE[\hat{\theta}] = E[\underbrace{(\hat{\theta} - \theta)^2}_{\uparrow}] \quad \text{--- may depend on } \theta.$$

Avg distance b/w est.  $\hat{\theta}$  and the parameter  $\theta$ .

HW 4 Supplement: ~~Can~~ Use Monte Carlo simulation to compare  $\hat{\theta}_{MSE}$  and  $\hat{\theta}_{MSE}$  of  $\theta$  to figure out which is better.

## Differentiation technique: Maximize the log-likelihood

function  $\ln L(\theta)$  with respect to  $\theta$  instead of  $L(\theta)$  as the former

tends to be easier. The value of  $\theta$  that maximizes  $L(\theta)$  also

maximizes  $\ln L(\theta)$ . (Why?) Since  $\log$  fn. is monotone, the value  $\theta$  that maximizes  $L(\theta)$  is also the value that maximizes  $\log L(\theta)$ .

Step 1: Set up the log-likelihood function.

$$\log L(\theta) = \log \prod_{i=1}^n f_{\theta}(x_i) = \sum_{i=1}^n \log f_{\theta}(x_i) \quad \left\{ \begin{array}{l} \text{the value that} \\ \text{maximizes} \\ \log L(\theta) \end{array} \right.$$

Step 2: Find the *likelihood equation* by partially differentiating  $\ln L(\theta)$  with respect to  $\theta$  and setting the derivative to equal to zero.

$$\boxed{\frac{\partial \log L(\theta)}{\partial \theta} = 0} \rightarrow \text{likelihood equation.}$$

Step 3: Solve the likelihood equation for  $\theta$ . The solution is MLE if it is a point of maxima (no need to verify).

~~Solution of~~  $\hat{\theta}_{MLE}$  = solution of this equation.

**Recall:** Some useful properties of natural log:

$$\lceil \ln \equiv \log \rceil$$

- $\ln(ab) = \ln(a) + \ln(b)$
- $\ln(a^b) = b \ln(a)$
- $\ln(e^a) = a$

**Ex:** Let  $X_1, X_2, \dots, X_n$  represent a random sample from an Exponential ( $\lambda$ ) distribution where  $\lambda > 0$ . Find the MLE of  $\lambda$ .

Recall:  $f_1(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0, \\ 0, & \text{otherwise.} \end{cases} \lambda > 0.$

Assume that all  $x$ 's are positive.

$$\left( \lambda e^{-\lambda x_1} \right) \left( \lambda e^{-\lambda x_2} \right) \dots \left( \lambda e^{-\lambda x_n} \right)$$

step 1: 
$$\begin{aligned} L(\lambda) &= \prod_{i=1}^n f_1(x_i) = \prod_{i=1}^n \lambda e^{-\lambda \left( \sum_{i=1}^n x_i \right)} \\ &= \lambda^n \cdot e^{-\lambda \left( \sum_{i=1}^n x_i \right)} \\ &= \lambda^n e^{-n \bar{x} \lambda} \end{aligned}$$

$$\left[ \bar{x} = \frac{\sum x_i}{n} \right]$$

$$\begin{aligned} \Rightarrow \log L(\lambda) &= \log [ \lambda^n e^{-n\bar{x}\lambda} ] \\ &= n \log(\lambda) - n\bar{x}\lambda. \\ \text{Step 2: } 0 &= \frac{d}{d\lambda} \log L(\lambda) = \frac{n}{\lambda} - n\bar{x} \end{aligned}$$

$$\begin{aligned} \text{Step 3: } \Rightarrow \frac{n}{\lambda} &= n\bar{x} \\ \Rightarrow \lambda &= \left[ \frac{1}{\bar{x}} \right] = \hat{\lambda}_{MLE}. \end{aligned}$$

$$\text{MOM? } \frac{1}{\lambda} = E[X] = \bar{X}$$

$$\Rightarrow \lambda = \left[ \frac{1}{\bar{X}} \right] = \hat{\lambda}_{MOM}$$





## Using R to get MLE

**Ex:** Recall the CPU data — CPU times for  $n = 30$  randomly chosen jobs (in seconds): 70, 36, 43, 69, 82, 48, 34, 62, 35, 15, 59, 139, 46, 37, 42, 30, 55, 56, 36, 82, 38, 89, 54, 25, 35, 24, 22, 9, 56, 19. Graphics suggested that the distribution of these CPU times may be right-skewed. Suppose we *assume* that the parent distribution is Gamma  $(\alpha, \lambda)$ , with both parameters unknown. What are MLE's of these parameters?

Recall:  $f(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, \quad x > 0, \lambda > 0$  ↑ positive

See book: For MLE.

$$\log L(\theta) = \sum_{i=1}^n \log \left[ \frac{\lambda^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\lambda x_i} \right]$$

MLE:

$(\hat{\alpha}, \hat{\lambda})$

— difficult to maximize in closed form.

# We will continue working with the CPU data  
# that we saw earlier

```
cpu <- scan(file="cputime.txt")
```

# Negative of log-likelihood function assuming gamma  
# parent distribution

$(x, \lambda) \rightarrow \text{data}$

```
neg.loglik.fun <- function(par, dat)
```

```
{
```

```
  result <- sum(dgamma(dat, shape=par[1], rate=par[2],  
    log=TRUE))
```

```
  return(-result)
```

```
}
```

# Minimize  $-\log(L)$ , i.e., maximize  $\log(L)$

? optim

Initial values are

close to  
MLE  
see both.

```
ml.est <- optim(par=c(3, 0.1), fn=neg.loglik.fun,
```

↑  
minimizes by default.

```
method = "L-BFGS-B", lower=rep(0,2), hessian=TRUE,  
dat=cpu)
```

*allows  
constraints.*  $\alpha > 0, \lambda > 0.$

```
# > ml.est
```

```
# $par
```

```
# [1] 3.63149628 0.07529459
```

*-  $\alpha, \lambda$*

```
# $value
```

```
# [1] 136.561
```

*= min value of  $\log L(\alpha, \lambda).$*

```
# $counts
```

```
# function gradient
```

```
# 20 20
```

```
# $convergence
```

```
# [1] 0
```

*— good.*

```
# $message
```

```

# [1] "CONVERGENCE: REL_REDUCTION_OF_F <= FACTR*EPSMCH"

# $hessian
      # [,1]      [,2]
# [1,]  9.501374 -398.4584
# [2,] -398.458449 19223.5065
# >

# MLE

# > ml.est$par
# [1] 3.63149628 0.07529459
# >

# their standard errors

# > sqrt(diag(solve(ml.est$hessian)))
# [1] 0.89720941 0.01994668

```

```
# >
```

```
# How well the fitted model represents the data?
```

```
# relative frequency (density) histogram
```

```
hist(cpu, freq=FALSE, xlab="cpu time",  
      ylab="relative frequency",  
      main="histogram vs fitted gamma distribution")
```

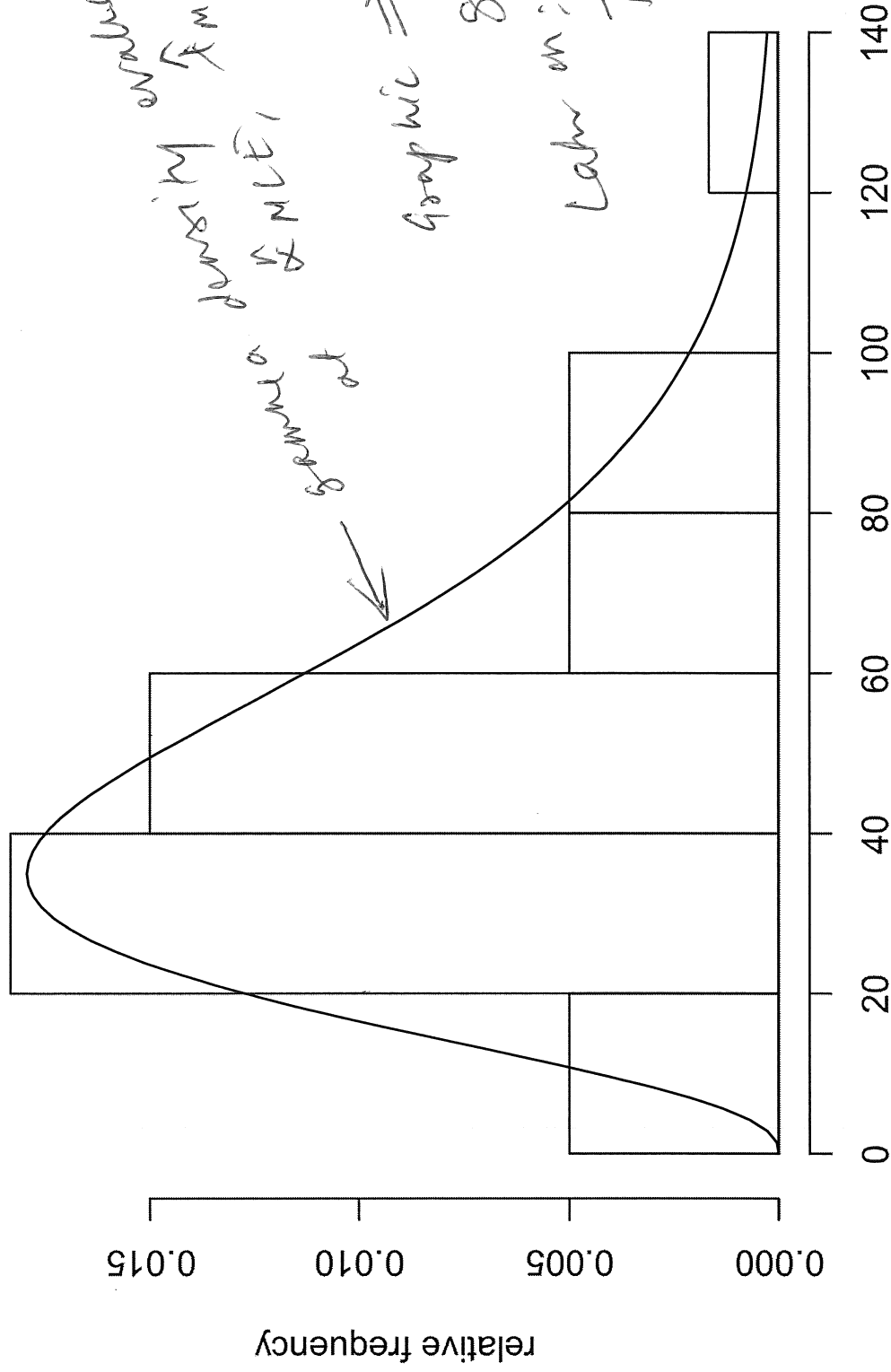
```
# superimpose the fitted density
```

```
gamma.pdf <- function(x, shape=ml.est$par[1],  
                      rate=ml.est$par[2])  
{ dgamma(x, shape=shape, rate=rate) }
```

```
curve(gamma.pdf, from=0, to=140, add=TRUE)
```

add this plot to the  
existing plot

# histogram vs fitted gamma distribution



cpu time