

COMPARING TWO ESTIMATORS

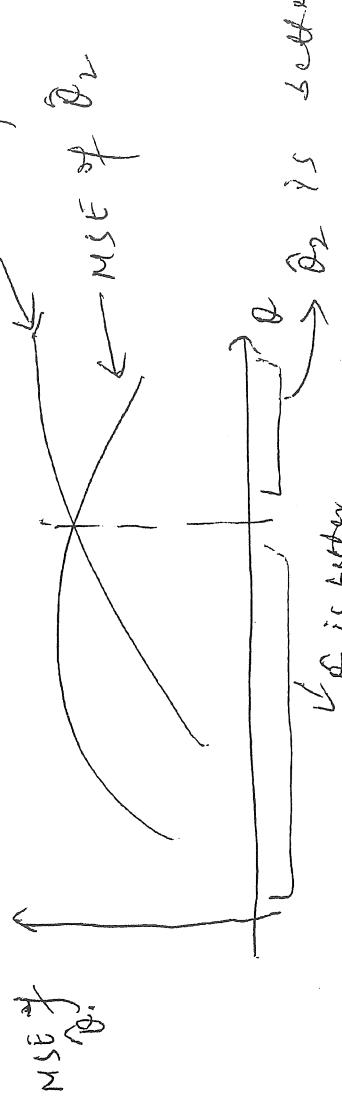
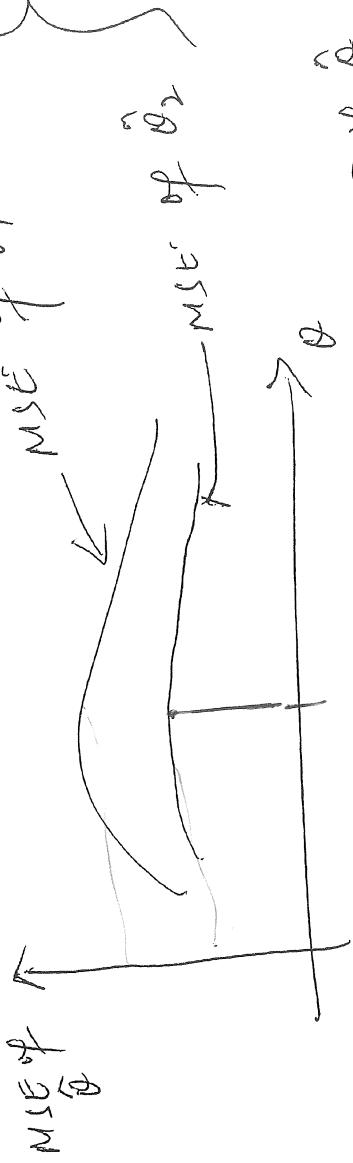
Mean squared error (MSE) of $\hat{\theta}$ if θ :

$$E[(\hat{\theta} - \theta)^2]$$

$$\text{Note: } \text{MSE}_{\hat{\theta}} = \frac{\text{var}[\hat{\theta}] + \text{Bias}_{\hat{\theta}}^2}{n}$$

$$E[(\hat{\theta} - \theta)^2] \rightarrow \text{depends on } \theta.$$

Average squared-distance from $\hat{\theta}$ to θ .



Note:

$$\begin{aligned} \text{MSE}_{\theta}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ &= E\left[\left(\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta\right)^2\right] \\ &= \text{Var}_{\theta}(\hat{\theta}) + [\text{Bias}_{\theta}(\hat{\theta})]^2, \end{aligned}$$

$$\text{where Bias}_{\theta} = E[\hat{\theta}] - \theta$$

$$\text{then } \text{MSE}_{\theta}(\hat{\theta}) = \text{Var}_{\theta}(\hat{\theta}) + \text{Bias}_{\theta}(\hat{\theta}).$$

- If $\hat{\theta}$ is unbiased, $\text{MSE}_{\theta}(\hat{\theta}) = \text{Var}_{\theta}(\hat{\theta})$.

Large sample properties of MLE $\hat{\theta}$ of θ

→ set of all possible values of X = sample space → certain conditions when n is large,

Result: Assume that $\{x : f_\theta(x) > 0\}$ is free of θ . Then, under

$$\hat{\theta} \approx N(\theta, \hat{I}^{-1}), \text{ where } \hat{I} = -\frac{\partial^2 \log\{L(\theta)\}}{\partial \theta^2} \Big|_{\theta=\hat{\theta}}$$

\hat{I} involves n

MLE

"Hessian information"
"Observed information"

Case 1: θ is scalar. Then, $\widehat{SE}(\hat{\theta}) \approx \sqrt{\hat{I}^{-1}}$

Case 2: θ is a vector, say, $\theta = (\theta_1, \dots, \theta_d)$. In this case, \hat{I} is a $d \times d$ matrix. Here $\widehat{SE}(\hat{\theta}_j) \approx (j\text{-th diagonal element of } \hat{I}^{-1})^{1/2}$.

Properties of MLE:

- Consistent; asymptotically unbiased
- Asymptotically normal
- Optimal if the assumed model holds
- Not a good choice if the assumed model does not hold

Recall: $SE(\bar{x}) = \frac{s}{\sqrt{n}}$

$$SE(\hat{\theta}) = \frac{s}{\sqrt{I}}$$

$\hat{\theta}$ is unbiased in \mathbb{T}

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 (α, λ), with both parameters unknown. What are MLE's of these parameters?

$$f_{\alpha, \lambda}(x) = \begin{cases} \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, & x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

~~cross product this~~

cannot multiply

$$\log \left[L(\alpha, \lambda) \right] = \log \left[\prod_{i=1}^n f_{\alpha, \lambda}(x_i) \right]$$

maximize this numerically.

Optim

to know

optimization methods

```
# 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
neg.loglik.fun <- function(par, dat)
{
  result <- sum(dgamma(dat, shape=par[1], rate=par[2]),
log=TRUE)
  return(-result)
}
```

Optimization do by gradient minimization.

log-likelihood

```
# Minimize -log(L), i.e., maximize log(L) of moment estimates
# approximating the true values
ml.est <- optim(par=c(3, 0.1),
                 fn=neg.loglik.fun,
                 method="Brent")
```

```

method = "L-BFGS-B",
dat=cpu) allows taking path up of constraints lower=rep(0,2), hessian=TRUE,
# > m1.est
# $par
# [1] 3.63149628 0.07529459 → meets
# $value
# [1] 136.561 min of -log L(x̂) ⇒ -136.561 is
# $counts
# function gradient
# 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

```

QQQ
12 / 14

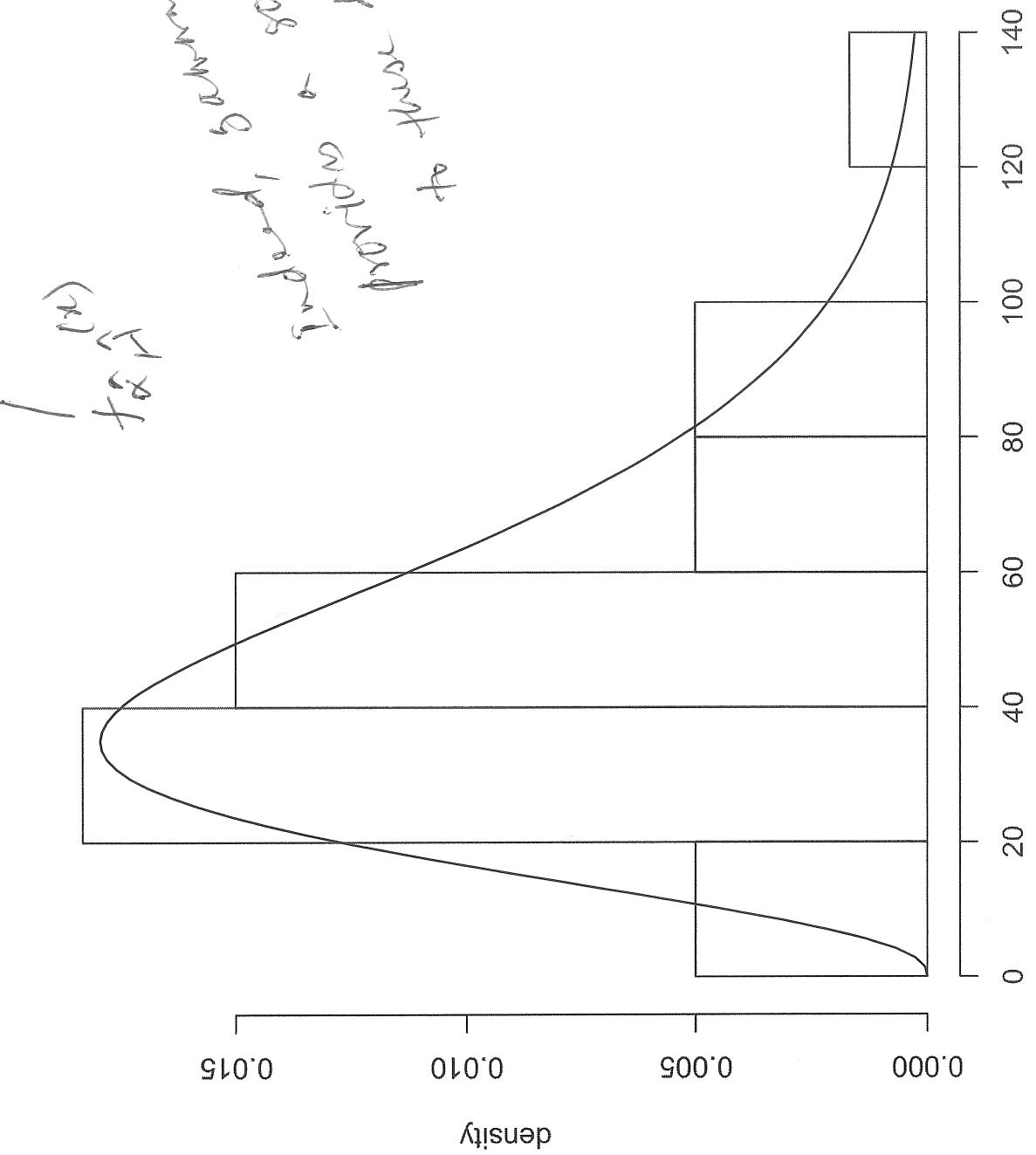
$\hat{\beta}$
 $\hat{\sigma}^2$
 \hat{f}
 $\hat{\epsilon}$
 $\hat{\alpha}$

>  with gamma curve ~ will see later in.

How well the fitted model represents the data?

```
# density histogram  
  
hist(cpu, freq=FALSE, xlab="cpu time" ,  
      ylab="density" ,  
      main="histogram vs fitted gamma distribution")  
  
# superimpose the fitted density  
  
gamma.pdf <- function(x, shape=m1.est$par[1] ,  
                      rate=m1.est$par[2])  
{ dgamma(x, shape=shape, rate=rate) }  
  
curve(gamma.pdf, from=0, to=140, add=TRUE)
```

histogram vs fitted gamma distribution



Pitfalls of Numerical techniques

- Function needs to be well-behaved
- Roots need to be captured in an interval
- Multiple roots,
- Convergence to local minima
- flat likelihood in region of maximum
- Negative estimate for variance.
- and others.

So far: Estimating θ by $\hat{\theta}$ — point estimation

Confidence intervals (Section 9.2)

Set up: Same as before, i.e.,

Population: $X \sim f_{\theta}(x)$

$\theta = \text{unknown}$

X_1, X_2, \dots, X_n (random sample)

? = $P[\sum_i (\hat{\theta} - \theta)^2 = 0] = 0$

cont. r.v., which is typically the case

Motivation: Estimator $\hat{\theta}$ is a single number that gives a plausible value of the unknown θ . But rarely the two will be equal. So, often it is preferable to give an interval of plausible values — a **confidence interval (CI)**, which contains the unknown θ with a specified high probability.

Definition: An interval $[L, U]$ is a $100(1 - \alpha)\%$ CI for θ if $L = L(\underline{X}_1, \dots, \underline{X}_n)$ and $U = U(\underline{X}_1, \dots, \underline{X}_n)$ are such that

$$P(L \leq \theta \leq U) = 1 - \alpha.$$

"Coverage probability"

random fixed random

$\Rightarrow P[\theta \text{ is captured by random interval } [L, U]] = P[\theta \in [L, U]]$

- L and U are *random*, so the CI is *random*.
- Parameter θ is not random — it is unknown but fixed.
- $(1 - \alpha)$ = *confidence coefficient* or *confidence level*.
- In practice, $(1 - \alpha) = 0.90$ or 0.95 (most common) or 0.99 .

A general method for constructing CI for θ

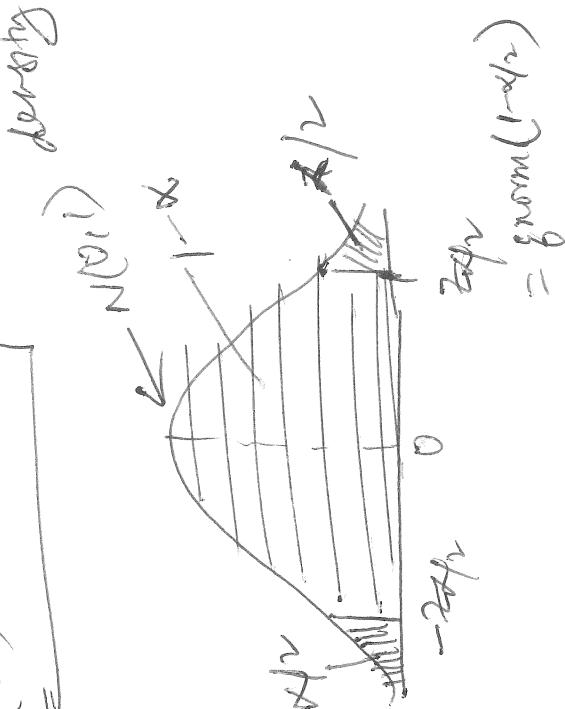
Step 1: Find an estimator $\hat{\theta}$ of θ that has a normal distribution with *known* variance, i.e., $\hat{\theta} \sim N(\hat{\theta}, \text{var}(\hat{\theta}))$.

Step 2: Standardize $\hat{\theta}$ to get Z , where

$$\text{"pivot"} \rightarrow Z = \frac{\hat{\theta} - \theta}{\text{SE}(\hat{\theta})} \sim N(0, 1)$$

A final note:
and $1 - \alpha = P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2})$
is called
normal
dist.
completely known.

$$\text{SD}(\hat{\theta}) = \text{SE}(\hat{\theta})$$



$$= \text{norm}(1 - \alpha)$$

Thus, the $100(1 - \alpha)\%$ CI is:
see next page

Have $2\pi\rho$ such that

$$1-\alpha = P[-2\pi\rho \leq Z \leq 2\pi\rho] \\ = P[-2\pi\rho \leq \frac{\theta - \hat{\theta}}{SE(\hat{\theta})} \leq 2\pi\rho] =$$

$$\text{Plug in } \theta \text{ diff. of} \\ \text{var diff. of} \\ \sum P[-2\pi\rho SE(\hat{\theta}) - \hat{\theta} \leq -\theta \leq \hat{\theta} + 2\pi\rho SE(\hat{\theta}) - \hat{\theta}]$$

Multiply by $SE(\hat{\theta})$

and expand $\hat{\theta}$

$$P[\frac{\hat{\theta} - 2\pi\rho SE(\hat{\theta})}{SE(\hat{\theta})} \leq \theta \leq \frac{\hat{\theta} + 2\pi\rho SE(\hat{\theta})}{SE(\hat{\theta})}]$$

Multiply by -1

$$\Rightarrow \text{From def., } 100(1-\alpha)\% CI \text{ for } \theta : \hat{\theta} \pm 2\pi\rho SE(\hat{\theta}) \\ = [L, U]$$

A general method, continued

Note 1: If the distribution of $\hat{\theta}$ is approximately normal, then the CI is also approximate.

Note 2: In case of MLE, if n is large, then often $\hat{\theta}$ approximately follows a $N(\hat{\theta}, \hat{I}^{-1})$ distribution. In this case, an approximate $100(1 - \alpha)\%$ CI for θ is:

$$\hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{1}{\hat{I}(\hat{\theta})}} \approx \hat{\theta} \pm z_{\alpha/2} \sqrt{\frac{1}{E(\hat{\theta})}}.$$

Note 3: If the distribution is not normal or n is not large, we can use the method of bootstrap to construct a CI (later).