

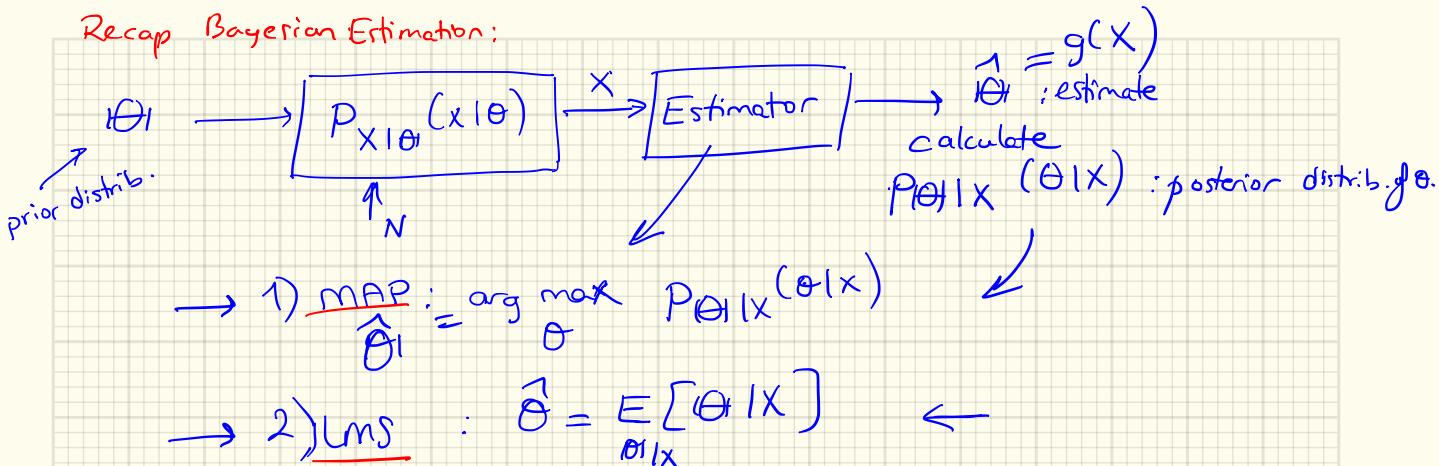
YZV 231E

04.01.2022

Probability Theory & Stats

GU.

Recap Bayesian Estimation:



Some properties of LMS Estimator :

- Estimator $\hat{\theta} = E[\theta|X]$

→ Estimation Error: $\tilde{\theta} \triangleq \hat{\theta} - \theta$

$$E[\tilde{\theta}|X] = E[(\hat{\theta} - \theta)|X] = E[\hat{\theta}|X] - E[\theta|X]$$

given $\hat{\theta}$ is known $\hat{\theta}$

$$- \hat{\theta} = 0$$
$$E[E[\tilde{\theta}|X]] = E[\tilde{\theta}] = 0$$

$$\rightarrow E[\tilde{\theta}|X] = 0$$

→ Estimator $\hat{\theta}$ is UNBIASED.

$$\rightarrow \text{Cov}(\tilde{\theta}, \hat{\theta}) = ? \quad \xrightarrow{\substack{\text{To find this:} \\ E[\tilde{\theta} \cdot \hat{\theta}] - E[\tilde{\theta}] \cdot E[\hat{\theta}]}}$$

$$E[\tilde{\theta} \cdot h(x) | X] = ?$$

$$h(x) \cdot E[\tilde{\theta} | X] = 0$$

$$E[E(\tilde{\theta} h(x) | X)] = E[\tilde{\theta} h(x)] = 0 \quad \text{for any } h(x)$$

We know $\hat{\theta}(x)$: a fn. of X $\therefore E[\tilde{\theta} \cdot \hat{\theta}] = 0$

$$\text{Cov}(\tilde{\theta}, \hat{\theta}) = E[\tilde{\theta} \hat{\theta}] - E[\tilde{\theta}] \cdot E[\hat{\theta}]$$

$$\text{Cov}(\tilde{\theta}, \hat{\theta}) = 0$$

The estimation error is uncorrelated w/ the estimate.

Given X , $h(x)$ is a number,

$$\tilde{\theta} = \hat{\theta} - \theta \rightarrow \theta = \hat{\theta} - \tilde{\theta} ; \hat{\theta} \times \tilde{\theta} \text{ are uncorrelated.}$$

Variance \sim a measure of uncertainty.

$$\text{Var}(\theta) = \underbrace{\text{Var}(\hat{\theta})}_{\substack{\text{uncertainty} \\ \text{in the true r.v.}}} + \underbrace{\text{Var}(\tilde{\theta})}_{\substack{\text{uncertainty} \\ \text{in the estimate}}} + \underbrace{\text{Var}(\epsilon)}_{\substack{\text{uncertainty} \\ \text{in the error.}}}$$

Linear LMS estimator:

Next, consider a simpler estimator of θ , of the form:

Let $g(x) = \hat{\theta} = ax + b$: affine mapping, parameters a, b are unknown.

minimize $E[(\theta - (ax + b))^2]$

Last time

minimize

$E[(\theta - g(x))^2]$: a generic estimator nonlinear.

$$g(x) = E[\theta | X] = \hat{\theta}$$

exercise: Derive

$$E[\hat{\theta}_L^2 - (ax+b)^2 - 2\hat{\theta}_L(ax+b)]$$

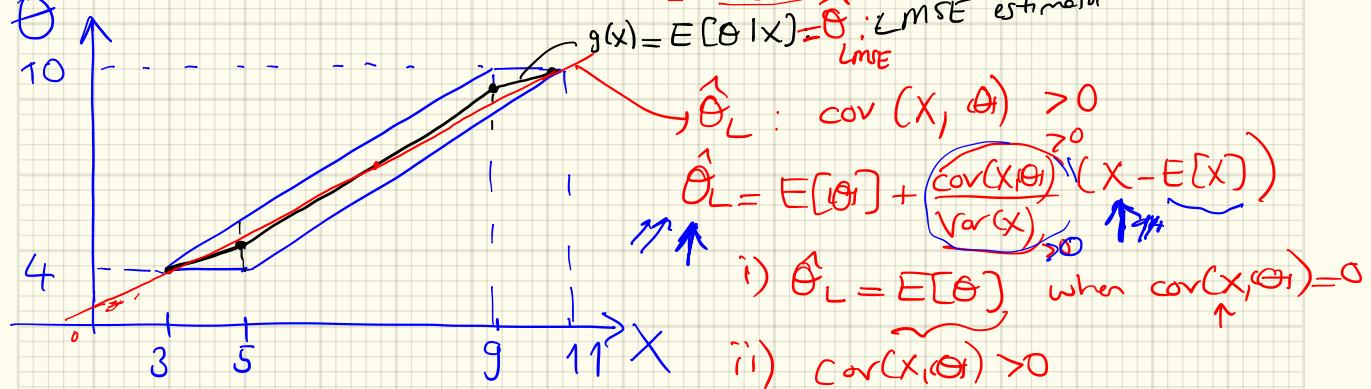
$$= E[\hat{\theta}_L^2] + a^2 E[X^2] + 2ab E[X] + b^2 - 2a E[\hat{\theta}_L X] - 2b E[\hat{\theta}_L]$$

$$\left\{ \begin{array}{l} \frac{\partial}{\partial a} \{ \} = 0 \\ \frac{\partial}{\partial b} \{ \} = 0 \end{array} \right. \rightarrow a, b \quad \checkmark$$

→ "Best" choice of a, b ; "best" linear estimator (in the MSE sense)

$$\hat{\theta}_L = E[\hat{\theta}] + \left[\frac{\text{cov}(X, \hat{\theta})}{\text{Var}(X)} \right] (X - E[X])$$

Recall prev ex.



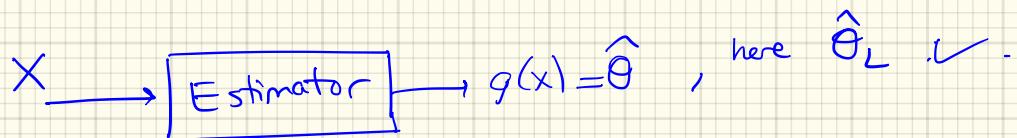
$$\rightarrow \text{What is the MSE? } E[(\hat{\theta}_L - \theta)^2] = (1 - \rho^2) \sigma_\theta^2$$

exercise: derive this result. $\rho = \frac{\text{Cov}(X, \theta)}{\sigma_X \sigma_\theta}$ correlation coeff

$$\boxed{\text{MSE}(\hat{\theta}_L) = (1 - \rho^2) \sigma_\theta^2}$$

variance of the original r.v. θ .

i) $\text{MSE} \uparrow \quad \checkmark \sigma_\theta^2 \uparrow$



ii) When the two r.v.s X, θ are correlated, \rightarrow we improve our estimate,

$\text{MSE} <$

iii) $\underline{\rho = 1} \xrightarrow{\text{maximal correlation case}} \rho^2 = 1 \quad \text{MSE} = 0 \quad \checkmark$

iv) $\rho = 0$: 2 r.v.'s are uncorrelated. $\text{MSE} = \sigma_\theta^2$:
measurements don't help to improve $\hat{\theta}$

Linear LMS w/ Multiple Data: We make several measurements x_1, \dots, x_n

Linear Estimator is of the form:

$$\hat{\theta} = a_1 x_1 + \dots + a_n x_n + b$$

"Optimal" LMSE estimator
 $\hat{\theta} = E[\hat{\theta} | x_1, \dots, x_n]$

→ Find "best" coefficients $a_1, \dots, a_n, b \rightarrow$ "optimal linear LMSE estimator."

Minimize $E[(\hat{\theta} - \theta)^2] = E[(a_1 x_1 + \dots + a_n x_n + b - \theta)^2]$

a_1, a_2, \dots, a_n, b

$$a_1^2 E(x_1^2) + 2a_1 a_2 E[x_1 x_2] + \dots$$

$$\frac{\partial}{\partial a_1} = 0 \quad \frac{\partial}{\partial a_2} = 0 \dots \frac{\partial}{\partial b} = 0$$

we need full posterior distrib. of $\theta | x_1, \dots, x_n$ to calculate $\hat{\theta}$.

But w/ the Linear LMSE estimator; we need to know expectations only;
 means, covariances, variances, ...)

We don't need to know the whole distribution.

— Linear LMS example : $\hat{\theta}$: unknown (random) parameter.

Measurement model: $X_i = \theta + w_i$:

$w_i \sim \mathcal{O}, \sigma_i^2$

prior $\rightarrow \theta \sim \mathcal{N}(\mu, \sigma_0^2)$

Assumption:

θ, w_1, \dots, w_n are independent

X_i : measurement in i^{th} experiment/measurement

$$X_1 = \theta + w_1$$

$$X_2 = \theta + w_2$$

\Rightarrow The form of the "optimal" linear estimator is very nice/neat form.

$$\hat{\theta}_L = \frac{\mu / \sigma_0^2 + \sum_{i=1}^n X_i / \sigma_i^2}{\sum_{i=0}^n 1 / \sigma_i^2}$$

→ sum of all coeffs
in the numerator

a weighted average
of μ, X_1, \dots, X_n

Weights are interesting!

— σ^2 \approx measure of uncertainty, $\frac{1}{\sigma^2} \approx$ reliability.

\star Weights are inversely proportional to the variance of the measurement.

\star The prior mean is treated the same way as the X_i 's.

→ Here, we did not impose any certain shape of the distributions of $w \times \theta$.

* For normal r.v.s; (all normal, $\theta \propto w$):

$$\hat{\theta}_L = E[\Theta | X_1, \dots, X_n]$$

\therefore The optimal (Lmse) estimate & the optimal linear (Lmse) estimate turn out to be the same.

→ If your measurement device measures X^2 or X^3 s instead of X ; want to estimate θ .

→ Choosing X_i in the Linear LMS vs LMS estimate

For LMS estimator

b/c X & X^2 or X^3 have the same info
posterior distribution of $\theta|X$ & $\theta|X^3$

$\theta|X$ or $\theta|(X^3)$ → do not matter
are the same ✓

For Linear LMS estimator

$$\hat{\theta} = aX + b \quad \text{vs} \quad \hat{\theta} = aX^3 + b$$

$$\hat{\theta} = a_1 X + a_2 X^2 + a_3 X^3 + b$$

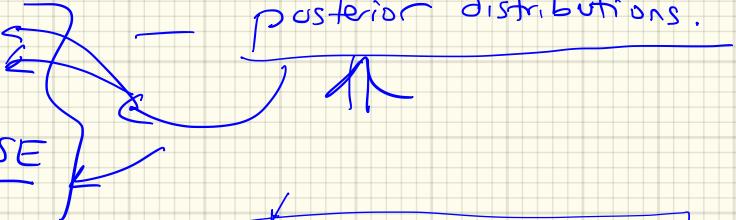
} In Linear LMS estimator,
which features/functions
of observations to choose
from matters.

Bayesian Estimation methods: Use a prior and Bayes' thm to find posterior distributions.

→ MAP

→ MSE

- Linear MSE



- Standard Models :

$$X_i = \theta + w_i;$$

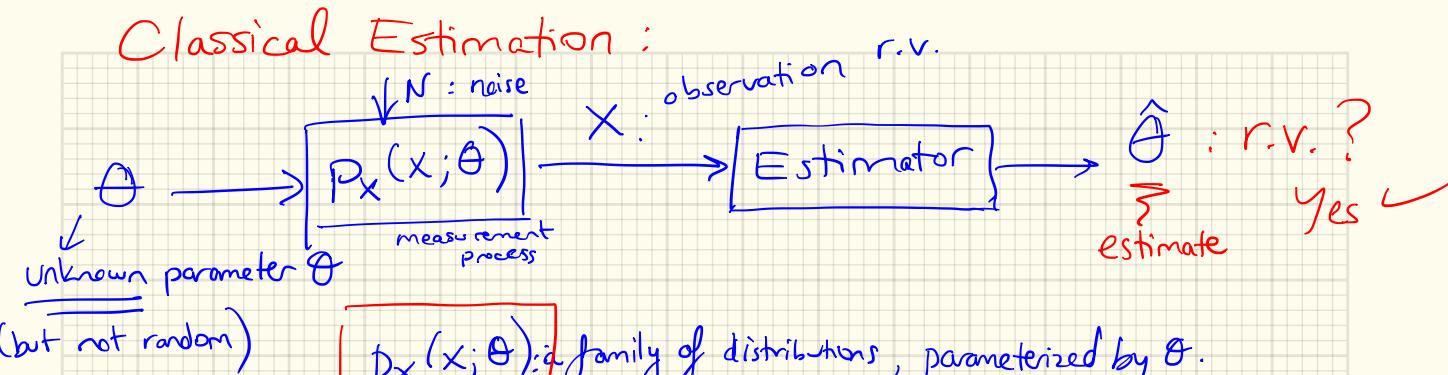
- X_i : Uniform on $[0, \theta]$; Uniform prior on θ

- X_i : Bernoulli(p) : Uniform^(Beta) prior on p .

- X_i : Normal: $\mathcal{N}(\mu, \sigma^2)$: normal prior on θ



Classical Estimation:



$P_X(x; \theta)$: a family of distributions, parameterized by θ .
 ↑; notation → be careful!

Note: not a conditional distrib.: that's why we don't use $P_X(x | \theta)$ as
 in Bayesian setting.

→ Design an estimator $\hat{\theta}$ to keep estimation error $\hat{\theta} - \theta$ small

- $\hat{\theta}(x)$: a fn. of an r.v. - $\hat{\theta}$ is an r.v.

→ Let's check out plausible estimator using $P_X(x; \theta)$



Maximum Likelihood Estimation: (MLE)

- $X \sim p_x(x; \theta)$: a model w/ unknown parameter θ .
- pick θ that "makes the data X most likely".

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} p_x(x; \theta)$$

Recall / compare to the Bayesian MAP estimator

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P_{\theta|X}(\theta|x) = \frac{P_{X|\theta}(x|\theta)}{P_X(x)} P_{\theta}(x)$$

θ is an r.v. \rightarrow we can have a prior $P_{\theta}(\theta)$.

model of the measurement process was $p_{X|\theta}$ - all θ 's are equally likely

★ If the prior P_{θ} is constant, then the ML estimation takes exactly the same form as the Bayesian MAP estimation. \rightarrow

Ex: Let X_1, \dots, X_n ; i.i.d. exponential r.v.'s w/ a certain parameter θ .
 $X_i \sim \theta \exp(-\theta \cdot x_i)$.

joint distrib. $P_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta) = \prod_{i=1}^n \theta \cdot e^{-\theta x_i}$

: prob. of observing a particular instance (X_1, \dots, X_n) given a particular θ value.

MLE: $\max_{\theta} \prod_{i=1}^n \theta \cdot e^{-\theta x_i}$: what's the value of θ that makes the X 's that we observed most likely?

\log  Maximizing this expression \equiv maximizing its logarithm. \rightarrow b/c log is monotonous fm

$$\max_{\theta} (n \log \theta - \theta \sum_{i=1}^n x_i) \quad \Rightarrow \quad \frac{\partial L(\cdot)}{\partial \theta} = 0$$

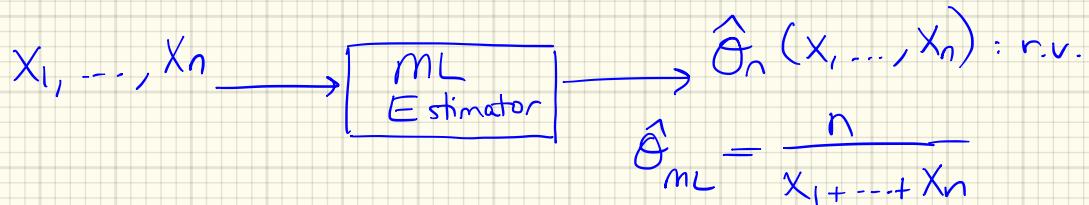
$$\frac{n}{\theta} - \sum_i x_i = 0$$

$$\Rightarrow \hat{\theta}_{ML} = \frac{n}{X_1 + \dots + X_n} := \text{reciprocal of the sample mean of the } X_i's : \frac{X_1 + \dots + X_n}{n}$$

Recall:

\rightarrow If X_i are exponential distrib ; the mean = $\frac{1}{\theta}$:

$\therefore \hat{\theta}_{ML}$ is a reasonable estimate. ✓



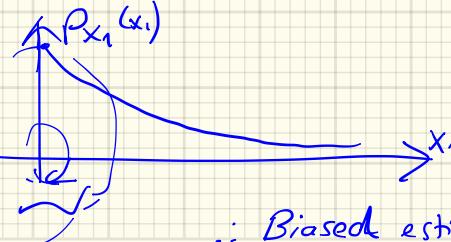
Desirable Properties of Estimators .

① Unbiased : $E[\hat{\theta}_n] = \theta$ ✓

We don't want the estimator to have a systematic error on the +ve or -ve side of the true parameter θ .

Ex : exponential example w/ $n=1$

$$\hat{\theta}_{ML} = \frac{1}{X_1} \rightarrow E[\hat{\theta}_M] = E\left[\frac{1}{X_1}\right] = \infty \neq \theta$$



\therefore Biased estimator.
in this case.

(2) Consistent : $\hat{\theta}_n \xrightarrow{\text{in probability}} \theta$

ex: exponential example: $X_i \sim \exp(\theta)$

$$\left(\frac{X_1 + \dots + X_n}{n} \right) \xrightarrow[\text{in prob.}]{\text{WLLN}} E[X] = \frac{1}{\theta}$$

Sample mean $\longrightarrow \mu$: true mean

$$\hat{\theta} = \frac{n}{X_1 + \dots + X_n} \xrightarrow{\text{prob}} \frac{1}{E[X]} = \theta \quad , \quad \forall \theta$$

$$\hat{\theta} \longrightarrow \theta$$

MLE here is a
consistent estimator

(3) "Small" Mean-Squared Error (MSE)

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= \underbrace{\text{Var}(\hat{\theta} - \theta)}_{\text{w.r.t. } P_{\theta}^1} + \underbrace{(E[\hat{\theta} - \theta])^2}_{\text{Bias} = \hat{\theta} - \theta} \\ &= \text{Var}(\hat{\theta}) + (\text{Bias})^2 \end{aligned}$$

For an (unbiased) estimator $\text{Bias} = 0$

Uncertainty in the estimate
our desire: small $\text{Var}(\hat{\theta})$

\Rightarrow We desire both Small Variance & Small Bias.

$$\text{MSE} = E[(\hat{\theta} - \theta)^2] = \underbrace{\text{Var}(\hat{\theta})}_{\substack{\text{Variance of the} \\ \text{estimator}}} + \underbrace{(\text{Bias})^2}_{\substack{\text{Bias of the estimator.}}}$$

But we have a Bias/Variance Trade-off:

e.g. $X \sim N(\theta, 1)$ say $\hat{\theta} = 100$: constant output Naive estimator

↑
parameter is
the mean

$$\text{Var}(\hat{\theta}) = 0 : \text{smallest variance} \checkmark$$

$$\text{Bias} = \hat{\theta} - \theta : \rightarrow \hat{\theta} : \theta = 0, 1 \dots \text{small}$$

$$\text{Bias} = 100 - \theta$$

$$(\text{Bias})^2 = (100 - \theta)^2$$

Bias \nearrow .

For this naive estimator, MSE is huge!

Conclusion: You may decrease the variance; but \exists a tradeoff
 \rightarrow your bias explodes!!

So to come up w/ low Bias & low Variance in your estimator
→ you have to come up w/ sophisticated estimators.

For parameter estimation (classical estimators): we study:

- 1) MLE : Maximum Likelihood estimation ✓
- 2) Sample Mean Estimator.

Sample Mean Estimator:

X_1, \dots, X_n : i.i.d , mean θ , variance σ^2 .

$$X_i = \underbrace{\theta}_{\text{mean}} + \underbrace{w_i}_{\text{noise}} \quad \left. \begin{array}{l} \hat{\theta} = \frac{X_1 + \dots + X_n}{n} \\ \downarrow \\ \text{Sample mean Estimator} \end{array} \right\}$$

w_i : i.i.d. mean 0, variance σ^2

Properties of the Sample Mean Estimator.

- $E[\hat{\theta}_n] = E\left[\underbrace{\frac{x_1 + \dots + x_n}{n}}_{\text{Unbiased}}\right] = \theta$ (Unbiased) $\Rightarrow \text{Bias} = 0$.
- $\hat{\theta}_n \xrightarrow{\text{WLLN}} \theta$ $\xrightarrow{\text{Consistent}}$
- mSE : $E[(\hat{\theta} - \theta)^2] = \underbrace{\text{Var}(\hat{\theta})}_{\frac{1}{n^2} \cdot n \cdot \sigma^2} + \underbrace{(\text{bias})^2}_{\frac{\sigma^2}{n}} = \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$

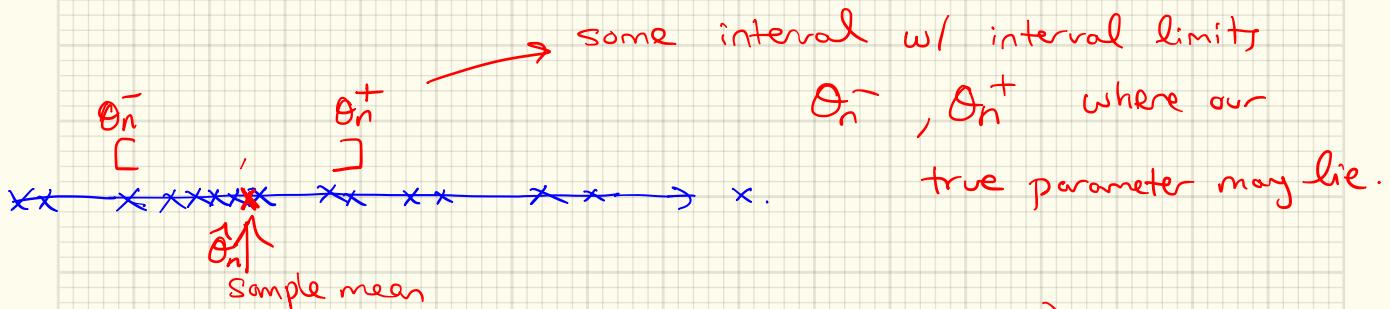
Note: If $X_i \sim \mathcal{N}(0, \sigma^2)$, i.i.d

Sample mean Estimator = ML estimate.
exercise: write $P_X(x; \theta) \Leftarrow \text{maximize w.r.t. } \theta$.

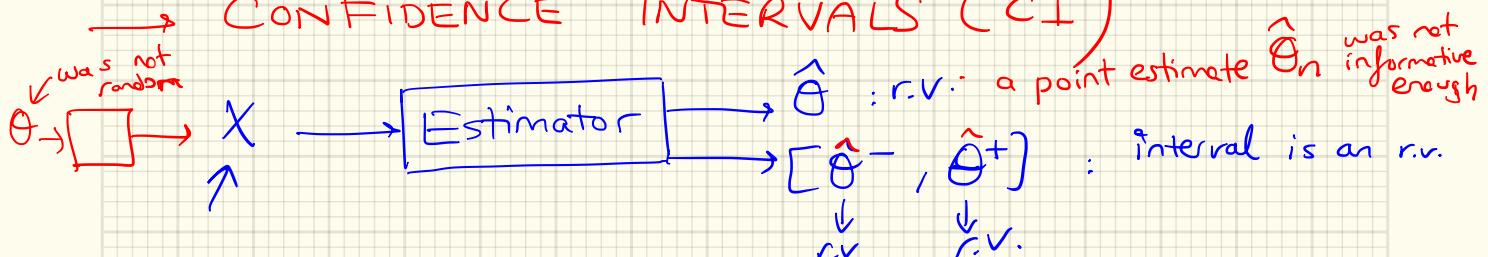
$$\frac{\partial}{\partial \theta} = 0$$

Now, → You report your sample mean : 3.27

Q. How reliable is that number? Can we trust it?



CONFIDENCE INTERVALS (CI)



— We pick an α →

- A $(1-\alpha)$ confidence interval: $[\hat{\theta}_n^- , \hat{\theta}_n^+]$

e.g. $\alpha = 0.05$

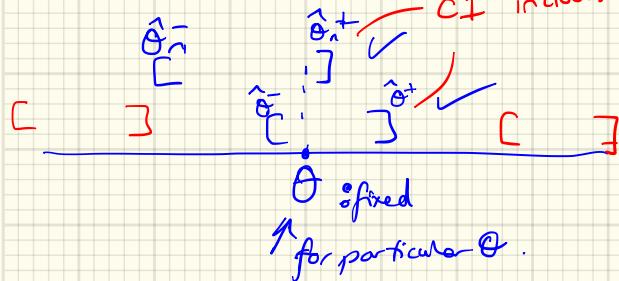
s.t. $P(\hat{\theta}_n^- \leq \theta \leq \hat{\theta}_n^+) \geq 1 - \alpha, \forall \theta$

→ 95% confidence Interval interpretation:

Rather than the statement that $\hat{\theta}$, w/ prob 95% falls in $[\hat{\theta}_n^-, \hat{\theta}_n^+]$,
b/c $\hat{\theta}$: not random,

Say this: w/ prob 95%, the interval falls on the true value of θ

CI includes θ 95% of the time



Q. How do we construct a 95% CI?

. CI in estimation of the mean :

$$\hat{\theta}_n = \frac{X_1 + \dots + X_n}{n} \quad \xrightarrow{\hspace{1cm}} \text{want to calculate}$$

95%
confidence interval

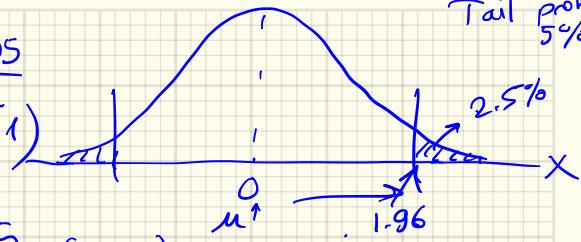
From the normal table:

$$P(x) : \mathcal{N}(0, 1)$$

Tail prob of 5%

$$\Phi(1.96) = 0.975 = 1 - \frac{0.05}{2}$$

Use CLT; standardize $\hat{\theta}_n \sim \mathcal{N}(0, 1)$



$$P\left(\left|\frac{\hat{\theta}_n - \theta}{\sigma/\sqrt{n}}\right| \leq 1.96\right) \approx 0.95 \text{ (CLT)}$$

rewrite

$$P\left(\hat{\theta}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{1.96\sigma}{\sqrt{n}}\right) \approx 0.95$$

construct
the CI

$\hat{\theta}_n^-$: lower end of the CI

$\hat{\theta}_n^+$: upper end

i) as $n \uparrow$; $[\hat{\theta}_n^-, \hat{\theta}_n^+]$: CI \rightarrow

we're more
confident that
our interval captures
the true θ .

ii) If $\sigma \searrow$ (our data has low uncertainty) : CI \searrow

More generally, how to construct the CI?

Let z be st. $\Phi(z) = \left(1 - \frac{\alpha}{2}\right)$ tail prob.

$$P\left(\hat{\theta}_n - \frac{z \cdot \sigma}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{z \cdot \sigma}{\sqrt{n}}\right) \approx 1 - \alpha$$

Here

- n is given
- $\sigma = ?$

For Unknown σ ; options:

1) Use an upper bound on σ ;

$$\text{e.g. } X_i \sim \text{Bernoulli} \quad \sigma \leq \frac{1}{2}$$

} more conservative
CI estimates
ie larger.

2) Estimate σ from the data \rightarrow use heuristic estimates

$$\text{e.g. } X_i \sim \text{Bernoulli}(\theta) \rightarrow \sigma = \sqrt{\theta(1-\theta)}$$

$\downarrow \hat{\theta} \downarrow \downarrow \downarrow \sigma \approx \sqrt{\hat{\theta}(1-\hat{\theta})}$

estimate
for the
standard
deviation

3) Use a generic estimate of the variance : Sample Variance

$$\sigma^2 = E[(X_i - \theta)^2]$$

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 \xrightarrow{\text{WLN}} \sigma^2$$

✓ Good if
we knew
 θ .

→ We don't know θ (the mean)

→ Insert the sample mean estimate

$$\tilde{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\theta}_n)^2 \xrightarrow{\text{WLN}} \sigma^2$$

$$E[\tilde{\sigma}_n^2] = \sigma^2 \quad : \text{unbiased}$$

Now use $\tilde{\sigma}_n = \sqrt{\tilde{\sigma}_n^2}$ in constructing your Confidence Interval limits