

YZV 231E

03.01.2022

Probability Theory & Stats

GU.

Decap: Limit Theorems:

we have a large # i.i.d. r.v.s.

e.g. polling ; 1000 people

X_1, X_2, \dots, X_n
values

— $M_n = \frac{X_1 + \dots + X_n}{n}$

↙
sample mean

: estimate of the expected value

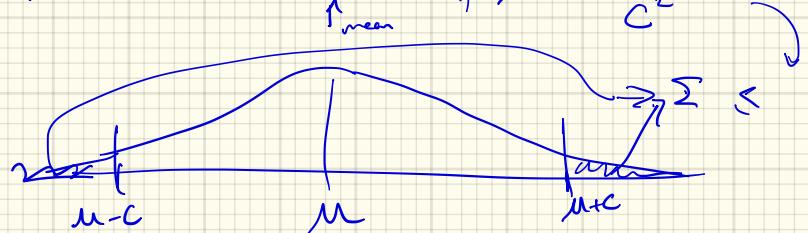
— M_n : is a random variable b/c the sample you collected is random.
Is M_n representative of the ^{the} expected value ?

$$M_n \xrightarrow{\quad} E[X] = \mu \text{ in probability}$$

— Convergence in probability : A seq. Y_1, Y_2, \dots, Y_n of r.v.s converges to a number a in prob. if as $n \rightarrow \infty$

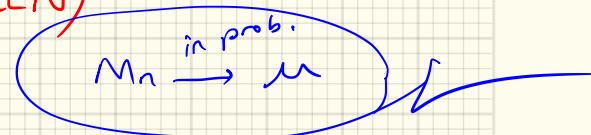
$$P(|Y_n - a| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$$

$$-\text{Chebyshoff Ineq: } P(|X - \mu| \geq c) \leq \frac{\sigma_x^2}{c^2}$$



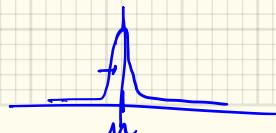
(Weak Law of Large Numbers (WLLN))

$$M_n = \frac{x_1 + \dots + x_n}{n}$$



$$\rightarrow E[M_n] = \underbrace{E[X_1] + \dots + E[X_n]}_{n} = \frac{n \cdot \mu}{n} = (\mu) \rightarrow \text{true mean of the population}$$

$$\rightarrow \text{Var}(M_n) = \frac{1}{n^2} n \cdot \underbrace{\text{Var}(X_i)}_{\sigma^2} = \left(\frac{\sigma_x^2}{n} \right) \xrightarrow{n \rightarrow \infty} 0$$



Back to Ex: Polling: p : fraction of population that prefers ...

$$X_i = \begin{cases} 1 & , \text{ if yes (logically)} \\ 0 & , \text{ otherwise} \end{cases}$$

Bernoulli r.v.

$$\bar{M}_n = \frac{X_1 + \dots + X_n}{n} : \text{ prediction of the fraction } p.$$

Sample mean

$$\text{Goal: } P(|\bar{M}_n - p| > \underbrace{0.01}_{\text{accuracy}}) \leq \underbrace{0.05}_{\gamma - \text{confidence}}$$

Given

Specs w/ 95% confidence of $\leq 1\%$ error

i) Collect data, calculate sample mean \bar{M}_n & check whether you satisfy the specs.

ii) Calculate the sample size n so that the specs are satisfied

$$\rightarrow \text{Chebyshev Ineq: } P(|\bar{M}_n - p| > 0.01) \leq \frac{\sigma_{\bar{M}_n}^2}{(0.01)^2} \leq 0.05 \quad \sigma_{\bar{M}_n}^2 = ?$$

$$\underline{\sigma_x^2} = p(1-p) \rightarrow \arg \max_p \sigma_x^2(p) \rightarrow p = \frac{1}{2} \quad \sigma_x^2 \leq \frac{1}{4}$$

$$P(|M_n - p| \geq 0.01) \leq \frac{\sigma_x^2}{n \cdot (10^{-4})} \leq \frac{1}{n \cdot 4 \cdot 10^{-4}} \leq 0.05$$

\approx

$$\Rightarrow n \geq \frac{1}{4 \cdot 10^{-5} \cdot 5 \cdot 10^{-2}}$$

$n = 50$ K samples :

fine on Twitter, but not practical for "in person" polling.

$\rightarrow \therefore$ play w/ specs :

e.g. $P(|M_n - p| > (0.03)) \leq 0.05$

accuracy less conservative

Chebyshev

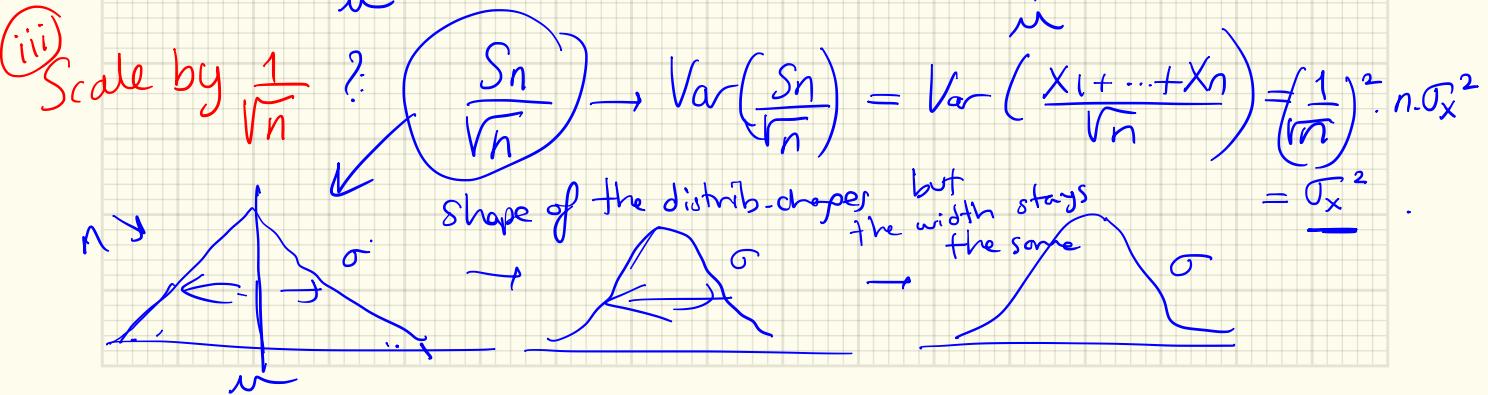
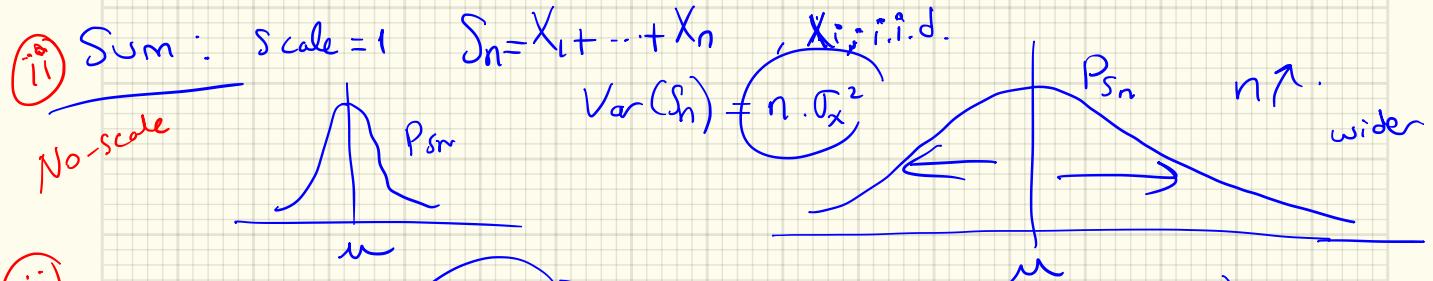
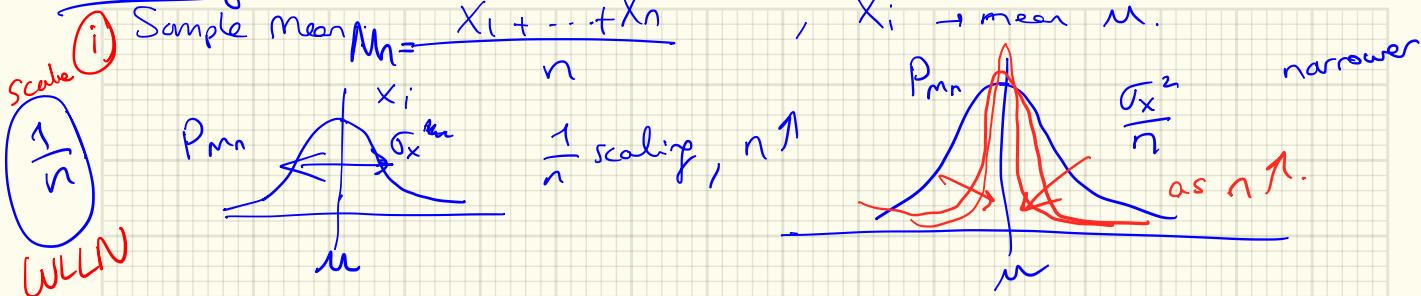
Saves a factor of ~ 10

$$(0.03)^2$$

$$3^2 \cdot 10^{-4}$$

$\rightarrow n \sim 5000$ people.

Scaling:



Different scalings of M_n ; X_1, \dots, X_n i.i.d w/ μ, σ^2

3 variants of the sum:

- 1) $S_n = X_1 + \dots + X_n$: variance $n\sigma^2$
- 2) $M_n = \frac{S_n}{n}$: variance $\frac{\sigma^2}{n}$
- 3) $\frac{S_n}{\sqrt{n}}$: constant variance σ^2

Central Limit Theorem (CLT)

Standardize: $S_n = X_1 + \dots + X_n$ by

$$Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}} = \frac{S_n - n \cdot E[X]}{\sqrt{n} \sigma}$$

$\downarrow r.v.$ $\underbrace{\sigma_{S_n}}_{\text{: standard deviation}}$ $\uparrow \sqrt{n}\sigma$

\Rightarrow

This is STANDARDIZATION of an r.v.

$$\left. \begin{array}{l} E[Z_n] = 0 \\ \text{Var}(Z_n) = 1 \end{array} \right\} Z_n \xrightarrow{\text{in distribution}} Z \sim \mathcal{N}(0, 1)$$

Let Z be a standard normal r.v.

Theorem (CLT) : For every c ,

$$P(Z_n \leq c) \xrightarrow{\text{Cdf of } Z_n \text{ r.v.}} P(Z \leq c) \xrightarrow{\text{standard normal cdf}}$$

$\hat{\Phi}(c) = P(Z \leq c)$: standard normal CDF is available from normal tables.

$$\frac{X_1, \dots, X_n}{\downarrow \sum} \xrightarrow{\text{i.i.d.: } X_i: \text{any distribution}} \xrightarrow{n \uparrow} \text{Normal Distribution}$$

$$Z_n = \frac{S_n - nE[X]}{\sqrt{n} \cdot \sigma}$$

$$\rightarrow S_n = \sigma\sqrt{n} Z_n + nE[X]$$

When n is large $\stackrel{\uparrow}{\sim} \text{Normal.} \therefore S_n$ is $\sim \text{Normal}$

b/c this is a linear transform of Z_n .

Ex: Let $X_i \sim N(0,1)$ i.i.d. Define $Y_N = \sum_{i=1}^N X_i^2$
 Approximate Y_N distrib. by a Gaussian. \rightarrow justified by CLT.
 b/c X_i i.i.d $\rightarrow X_i^2$ i.i.d. $\text{Var}(Y_N) = N \cdot \text{Var}(X^2)$

CLT says: $\tilde{Y}_N = \frac{Y_N - E[Y_N]}{\sqrt{Y_N}} = \frac{Y_N - N \cdot E[X^2]}{\sqrt{N \cdot \text{Var}(X^2)}}$ $\xrightarrow[\text{CLT.}]{\sim} N(0,1)$

$$E[X^2] = \text{Var}(X) = 1$$

$$\begin{aligned} \text{Var}(X^2) &= \underbrace{E[X^4]}_{= 3} - \underbrace{(E[X^2])^2}_{= 1} = 2 \\ &\quad (\text{check}) \end{aligned}$$

$$Y_N = \sum_{i=1}^N X_i^2$$

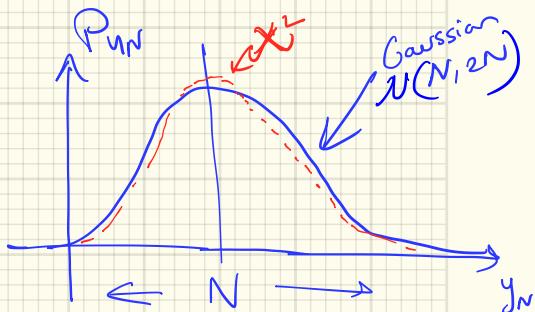
$$\rightarrow \tilde{Y}_N = \frac{Y_N - N}{\sqrt{2N}} \underset{n \rightarrow \infty}{\approx} N(0, 1) \quad \text{by CLT.}$$

$$\rightarrow \text{Approx. } Y_N = \sqrt{2N} \tilde{Y}_N + N \rightarrow Y_N \sim N(N, 2N)$$

\tilde{Y}_N Gaussian \checkmark Gaussian w/ mean 0 var 1

$$E(Y_N) = E(\sqrt{2N} \tilde{Y}_N + N)$$

$$\text{Var}(Y_N) = 2N \cdot 1^0 + 0$$



$Y_N : X^2$ chi-squared.

as $N \uparrow$
approx. by a

Gaussian becomes better standard dev. $\sqrt{2N}$

$Y_N \sim \chi_N^2$ exact

Ex: Pollster's problem using CLT.

- p : fraction of the population that prefers " - - - "

- X_i : i^{th} randomly selected person , $X_i = \begin{cases} 1 & \text{if yes} \\ 0 & \text{otherwise} \end{cases}$

$M_n = \frac{X_1 + \dots + X_n}{n}$: this is the estimate for the fraction of the population that preferred . . .

→ We define 2 specifications for the poll : $\equiv 2$ parameters

$$\Rightarrow P(|M_n - p| \geq 0.01) \leq 0.05 \rightarrow \text{confidence}$$

$\stackrel{\text{accuracy}}{\approx}$ Specs → Want probability 95% that our estimate M_n is within 1% of the true value p .

Event of interest : $|M_n - p| > 0.01$

Standardize . $\rightarrow \left| \frac{X_1 + \dots + X_n - np}{\sqrt{n}} \right| \geq 0.01$

$\stackrel{\text{some event}}{\approx} \left| \frac{(X_1 + \dots + X_n - np)}{\sqrt{n} \sigma} \right| \geq \frac{0.01 \sqrt{n}}{\sigma}$

standardized r.v $\rightarrow Z_n \stackrel{\text{?}}{=} -$

Now: $P(|M_n - p| > 0.01) = P(|Z_n| > \frac{0.01\sqrt{n}}{\sigma})$

Using CLT $\rightarrow \approx P(|Z| > 0.01\sqrt{n}/\sigma)$

where Z is a standard normal r.v.

Q: What is σ ? We know an upper bound

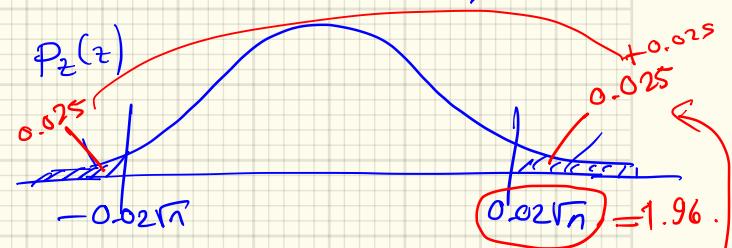
$$\begin{aligned} \sigma^2 &\leq \frac{1}{4} \\ \sigma &\leq \frac{1}{2} \end{aligned}$$

$$P(|Z| > \frac{0.01\sqrt{n}}{\sigma}) \leq P(|Z| > 0.02\sqrt{n}) \quad \frac{1}{\sigma} > 2$$

i) Given $n = 10,000$

$$\begin{aligned} P(|Z| > \frac{0.02\sqrt{10^4}}{2}) &= P(|Z| > 2) \\ &= 2 P(Z > 2) \\ &= 2(1 - P(Z \leq 2)) \\ &\stackrel{\Phi(2)}{=} 2(1 - 0.9772) \\ &\stackrel{\text{use Normal table}}{\downarrow} \\ &= 2(1 - 0.9772) \end{aligned}$$

$$\begin{aligned} &= 0.0456 \rightarrow \text{prob of error} = 4.6\% \\ &< 5\% \end{aligned}$$



ii) Find n s.t. $P(|Z| > 0.02\sqrt{n}) = 0.05$

$$\begin{aligned} \Phi(-1.96) &= 0.975 \\ 1.96 &\text{ from the table} \\ 0.02\sqrt{n} &= 1.96 \Rightarrow n = 9604 \end{aligned}$$

→ W/ $n=9604$ people in the poll
our probability of error is 0.05. ✓

→ We use CLT in different ways → approximate distributions.

Apply to Binomial:

$X_i: \text{Bernoulli}(p)$, indep. → Fix p , $0 < p < 1$

$S_n = X_1 + \dots + X_n$; Binomial(n, p)

$\xrightarrow{\text{Binomial}}$ Mean np , Variance $= np(1-p)$ ✓

CDF of $\frac{S_n - np}{\sqrt{np(1-p)}}$ $\xrightarrow{n \uparrow}$ Standard Normal Distrib.

Check whether this approx. is
good: $n=36$, $p=0.5$; find $P(S_n \leq 21)$

$$\begin{aligned} \text{CLT:} \\ \text{mean} &= 18 = np \\ \text{var} &= 9 = np(1-p) \end{aligned}$$

standardize

$$\frac{S_n - 18}{\sqrt{9}} \leq \frac{21 - 18}{3} = 1$$

std normal

$$\frac{Z_n}{3} \leq 1 \approx Z \leq 1$$

$$\rightarrow \Phi(1) = P(Z \leq 1) = 0.843 \rightarrow \text{an approximate answer.}$$

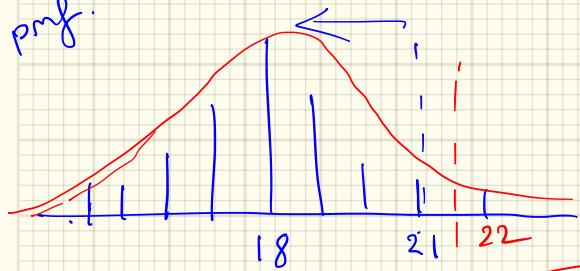
from the table

OK. but b/c
 S_n : a discrete
 r.v

Exact answer :

$$\sum_{k=0}^{21} \binom{36}{k} \left(\frac{1}{2}\right)^{36} = 0.8785$$

pmf.



w/ CLT : we treated S_n as normal

$$P(S_n \leq 21) = P(S_n < 22)$$

\therefore use $P(S_n \leq 21.5)$

$$\frac{S_n - 18}{3} \leq \frac{21.5 - 18}{3} = 1.17$$

De-Moivre - Laplace CLT.

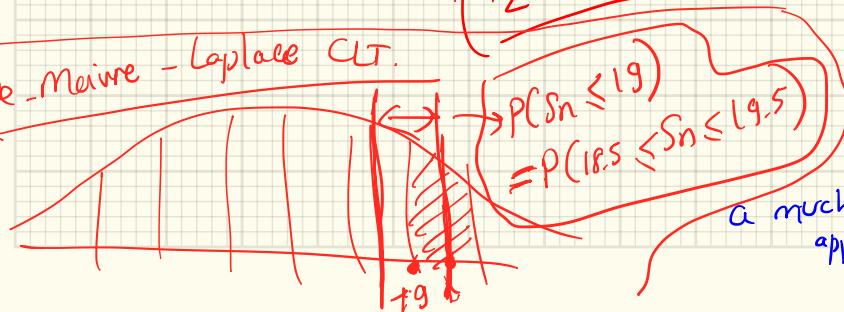
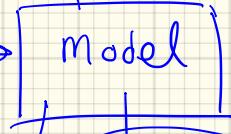
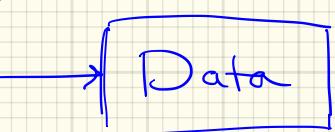
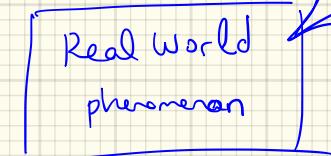


Table $P(Z \leq 1.17) = 0.879$

a much better approximation

STATISTICAL INFERENCE :



e.g. customer arrivals

estimate
Electricity usage
Watts

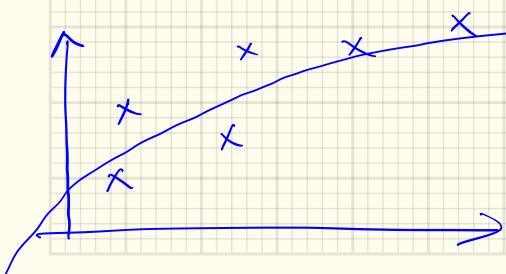
→ e.g. Poisson r.v.
→ Come up w/ a probability model
→ then calculate its parameters.

Real world

Estimation Problems



Hypothesis Testing Problems



1) Bayesian Statistical Estimation

2) Classical Statistical Estimation

This is what we do w/ probability theory.

e.g.
here:
the arrival rate
 λ .

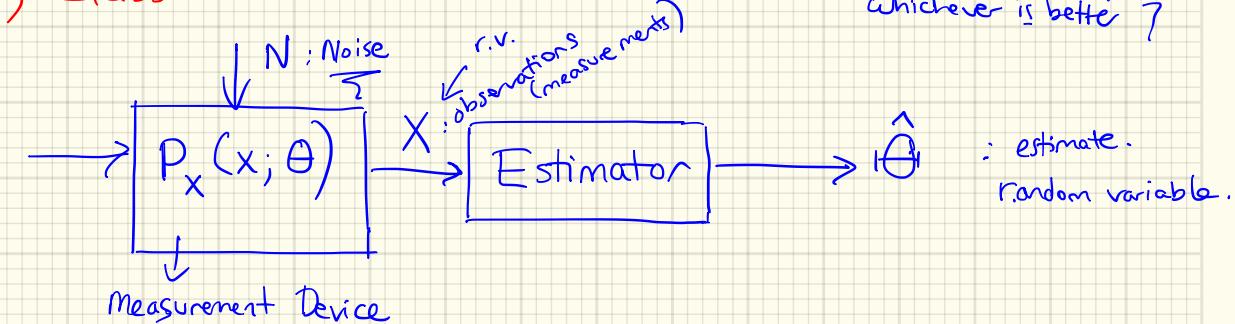
Estimation :

1) Classical Statistics

Classical vs Bayesian Statistics

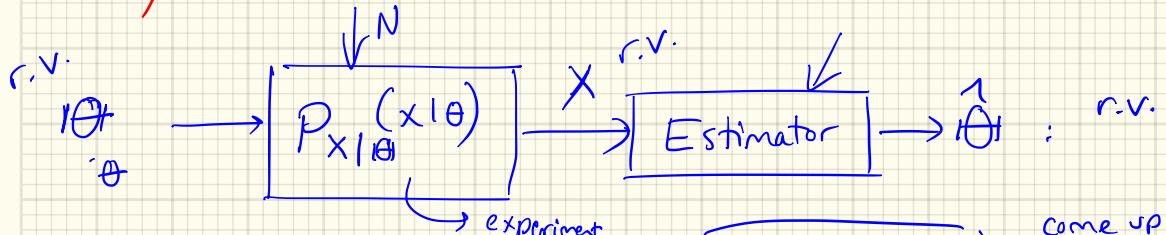
↓
ongoing debate.
↓ whichever is better ?

Unknown quantity
we're trying
to estimate

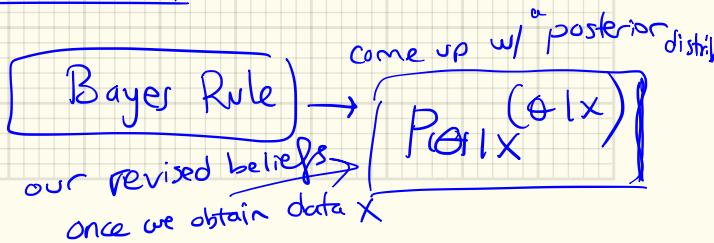


θ : Unknown parameter (not an r.v.)

2) Bayesian Statistics : Use priors on θ



$P_\theta(\theta)$: prior distribution.
our initial beliefs about θ
before the experiment

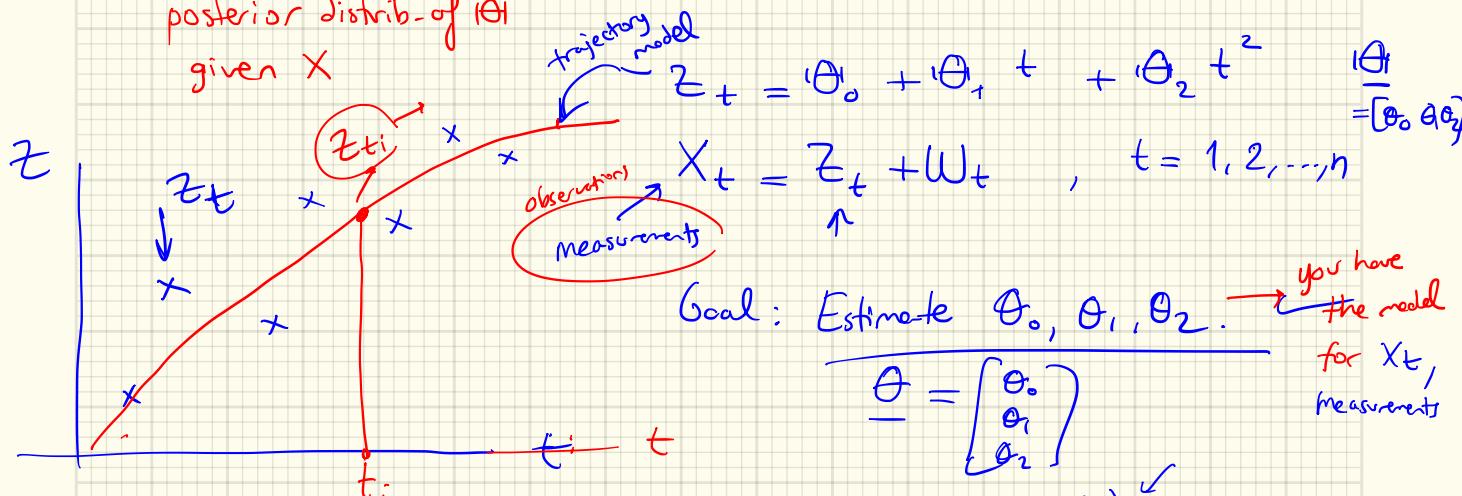


Bayesian inference: Use Bayes' rule conditional model of the experimental process prior

$$P_{\Theta|X}(\theta|x) = \frac{P_X(x|\theta) - P_\theta(\theta)}{P_X(x)}$$

we want to calculate posterior distrib. of θ
given X

P : pmfs
or
pdfs



continuous r.v.s θ_i .

comes w/ $P_{\theta_0, \theta_1, \theta_2 | X_1, \dots, X_n}(\theta_0, \theta_1, \theta_2 | X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n | \underline{\theta}) P_{\underline{\theta}}(\underline{\theta})}{P(X_1, \dots, X_n)}$

ex: Coin w/ unknown parameter $\theta \xrightarrow{\text{eg. } \epsilon_p}$

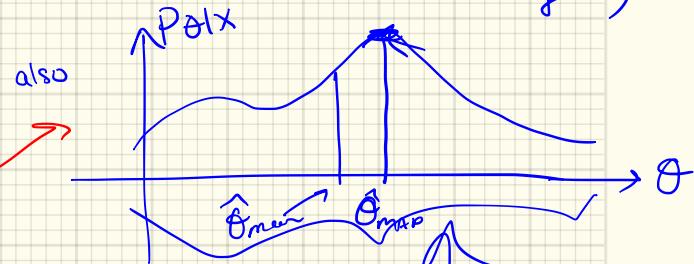
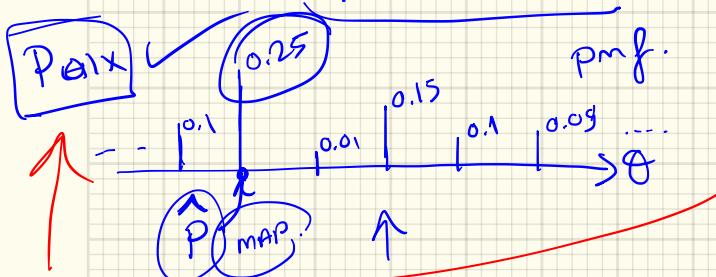
→ In classical statistics: we flip the coin many times x_1, \dots, x_n form s_n .

$$\hat{\theta}_n = \frac{s_n}{n} \xrightarrow{\text{LLN}} \theta \quad ?$$

check this w/ LLN.

Now: Bayesian approach: Assume a prior on θ (e.g. uniform prior if you don't know anything -)

Goal: $p(\theta|x) = ?$ use Bayes rule to find!



Output of Bayesian Inference: $p(\theta|x)$

→ Interested in a single answer/output



Bayesian
setting

r.v.

Θ

$P_{\Theta}(e)$
Initial belief

generates data

$P(x|\theta)$

conditional prob.
distrib,
depends on θ

r.v.

X

observe X ,

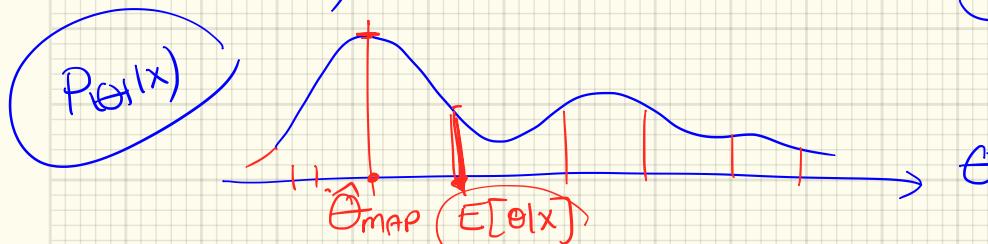
r.v.

$\hat{\theta} = g(x)$

Estimator $g(\cdot)$

Given X ; what is $\theta|x$:

$P_{\theta|x}(\theta|x)$
revised belief



Q: If we want to report a single number for θ

$P(\theta|x)$
complete answer
to a Bayesian inference
problem.
Point Estimate of θ

①

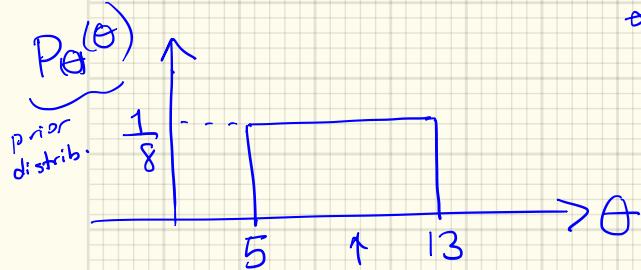
MAP (Maximum a posteriori) Estimate

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|x)$$

②

LMS (Least-Mean-Squared) Error Estimate :

Least Mean Squares Estimation : (LMS)



* Estimation in the absence of information
(no X observed)

→ Goal: Come up w/ a point estimate.

using LMS criterion

point estimate c : minimize

: make a small error;
squaring the error →
we penalize large errors
even more.

$$\min E \left[(\theta - c)^2 \right]$$

This is called the Least-Mean Squares Estimation

Optimal estimate $c = ?$

$$E[\theta^2] - 2E[\theta] \cdot c + c^2$$

$$\frac{d}{dc} = -2E[\theta] + 2c = 0 \rightarrow c = E[\theta]$$

min.
by taking
 $\frac{\partial}{\partial c}$ = 0

$$\boxed{c = E[\theta]}$$

; above ex. $c = 9$. $\underline{= E[\theta]}$

To judge how good is our estimate)

"optimal"
mean-squared
error

$$E \left[\underbrace{(\theta - E[\theta])^2}_{\text{least-mean-squared error}} \right] = \text{Var}(\theta)$$

→ Expectation is the "best" estimate if you're minimizing the least-mean-squared error.

→ Next : we have data X , how will this estimate change ?

LMS estimation of θ based on $\underline{\underline{X}}$: Given X , we are in a conditional universe

$$\min. E[(\theta - c)^2 | X = x]$$

is minimized by $c = E[\theta | X = x]$

now, you use conditional expectation



$$g(\cdot) \rightarrow g(x) = \hat{\Theta}$$

Let $\hat{\Theta} = g(x)$ minimize

What other forms of estimator?

$$\min_{x, \theta} E \left[(\Theta - g(x))^2 \right] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\Theta - g(x))^2 p_{X|\Theta}(x, \theta) dx d\theta$$

$$\min = \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} (\Theta - g(x))^2 p_{\Theta|X}(\theta|x) d\theta \right) \frac{p_X(x)}{p_X(x) > 0} dx.$$

We can minimize $\int_{-\infty}^{\infty} (\Theta - c)^2 p_{\Theta|X}(\theta|x) d\theta = f(c)$

$$\frac{\partial f}{\partial c} = 0$$

$$2 \int_{-\infty}^c (\Theta - c) p_{\Theta|X}(\theta|x) d\theta = 0$$

$$\int \Theta p_{\Theta|X}(\theta|x) d\theta = c \cdot \int p_{\Theta|X}(\theta|x) d\theta \Rightarrow$$

$= 1$.

c unknown.

$$\rightarrow C = E[\theta|x] = \int \theta \cdot P_{\theta|x}(\theta|x) d\theta$$

$$= g(x) = \overbrace{\theta|x}$$

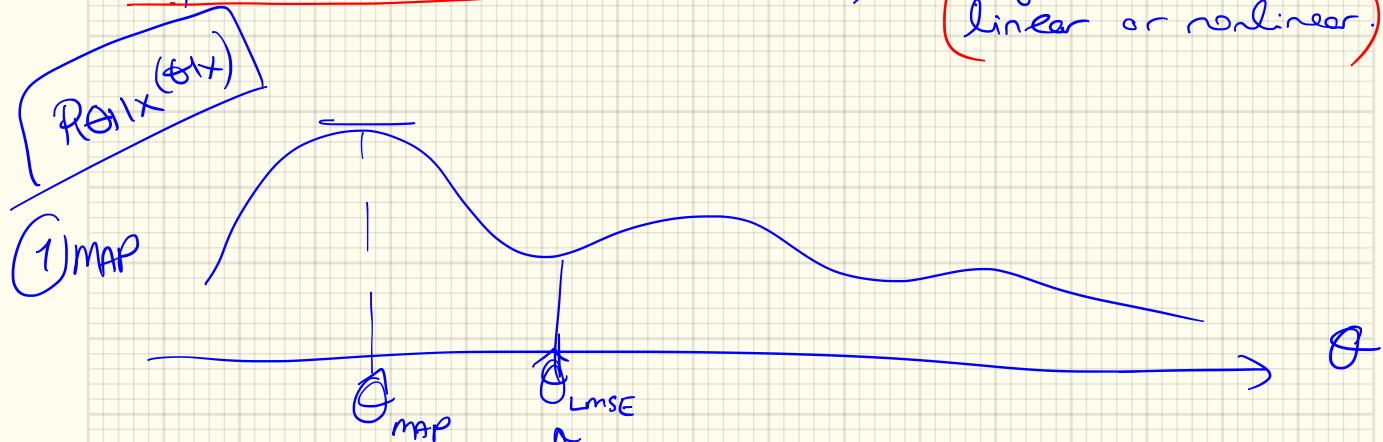
$$\Rightarrow \therefore \boxed{\hat{\theta}(x) = E[\theta|x]}$$

$$= g(x)$$

"Optimal" estimator (LmSE sense)

in the mSE sense.
(LmSE) sense

among all estimators
(linear or nonlinear.)



② $E[\theta|x]$ = conditional mean of the posterior density LmSE estimator

Bayesian estimation w/ several measurements

Unknown r.v. Θ

Observe values of r.v.s

X_1, \dots, X_n

"Best" estimator

LMS E:

MAP: $\arg \max P_{\Theta}(x_1, \dots, x_n)$

calculate, using Bayes' rule:

$$P_{\Theta}(x_1, \dots, x_n)$$

$$E[\Theta | x_1, \dots, x_n]$$

Issues:

- Have to come up w/ a reasonable prior \sim distrib. for Θ .
- Have to calculate the posterior distrib $P_{\Theta|x}$ & its expectations:

Those Calculations can be hard \rightarrow even intractable!!

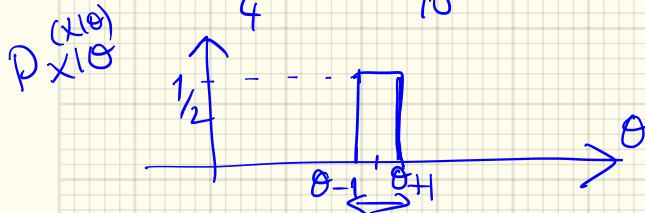
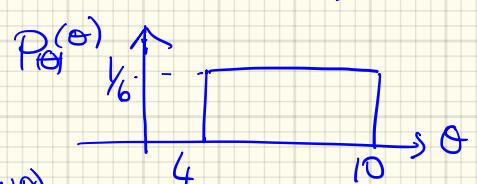
w/ multi-dimensional integrals ..

$$\int \dots \int$$

Ex: Suppose $\theta \sim U[4, 10]$; $\theta \rightarrow$ measurement X

Observation model:

$$X = \theta + U$$



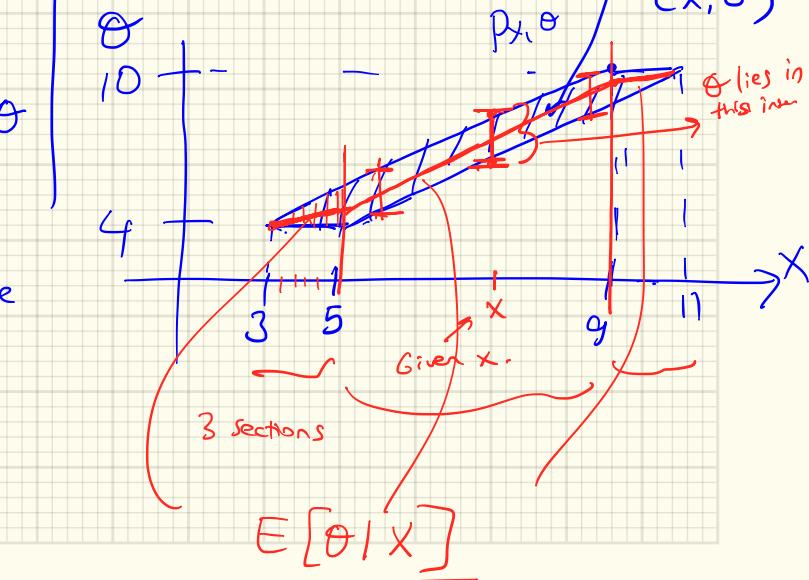
④ A plausible estimator.
"optimal" estimator in LMS-sense
 $E[\theta|x]$; P_theta|x

construct the joint density $P_{x,\theta}$

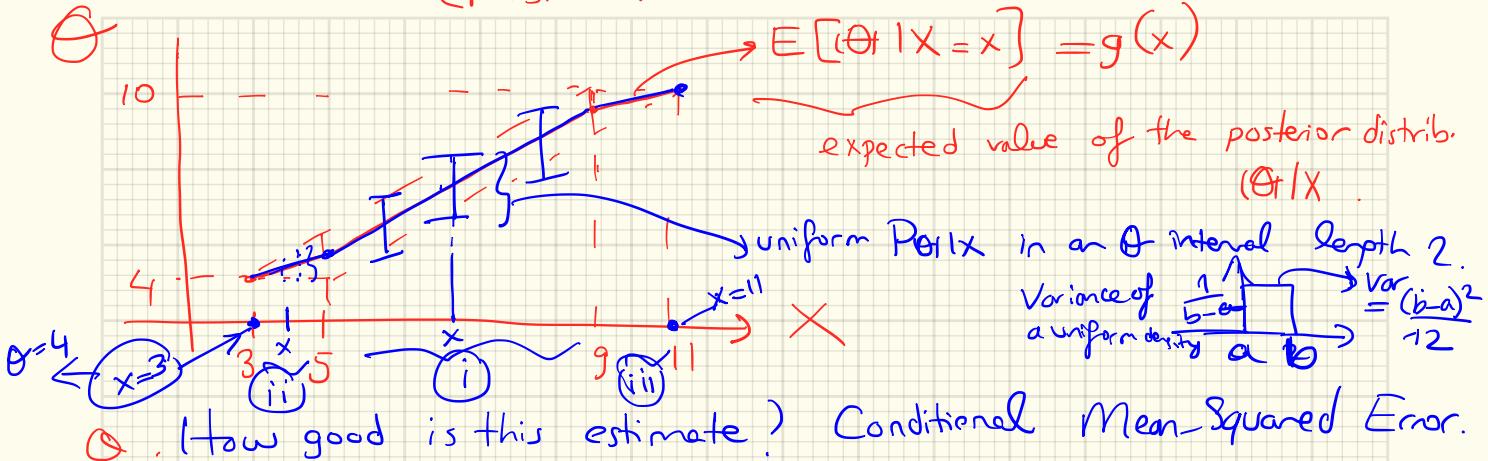
$$P_{x,\theta} = P_\theta \cdot P_{x|\theta}$$

$$= \frac{1}{6} \cdot \frac{1}{2} = \frac{1}{12}$$

on the support (x, θ)



LMS-estimator:



Q. How good is this estimate? Conditioned Mean-Squared Error.

$$E[(\theta - E[\theta|X])^2 | X=x]$$

true value of θ estimate we made

: variance of the posterior distrib.

