

12.12.2022

YZV 231E

Probability Theory & Stats

Week 12

Gü.

# Recap: Limit Theorems

**WLLN:** Large # i.i.d. r.v.s

Convergence in prob.  

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

$X_1, \dots, X_n$  : like sampling from a population.

$(X_1, X_2, \dots, X_n) \xrightarrow{n \rightarrow \infty}$   
 ↳ "close" to  $a \pm \varepsilon$

r.v.  $M_n$   
 sample mean

$$M_n = \frac{X_1 + X_2 + \dots + X_n}{n} : \text{estimate of the expected value.}$$

WLLN:  $M_n \xrightarrow[\text{in probability.}]{\text{convergence}} E[X] = \mu$

Recall Tchebycheff  $P(|X - \mu| \geq c) \leq \frac{\sigma_x^2}{c^2}$   
 $\mu, \sigma^2$

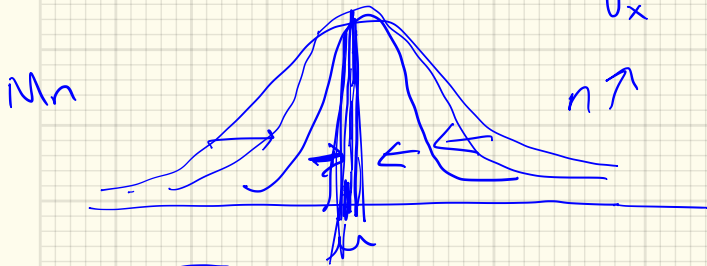


$$\rightarrow M_n = \left( \frac{X_1 + X_2 + \dots + X_n}{n} \right)$$

$$E[M_n] = \frac{E[X_1] + \dots + E[X_n]}{n} = \frac{n \cdot \mu}{n} = \mu$$

μ: true mean of the population

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



**Pollipo ex:**  $\underline{p}$ : fraction of population that prefer mt.

$$\underline{M_n = \frac{X_1 + \dots + X_n}{n}}$$

prediction of the fraction  $p$ .

Bernoulli:  $X_i = \begin{cases} 1, & \text{if yes (prefer)} \\ 0, & \text{if no (does not prefer)} \end{cases}$

Goal:  $P(|\underbrace{M_n}_{\text{true fraction } n} - p| \geq \underbrace{0.01}_{\text{accuracy}}) \leq \underbrace{0.05}_{\text{confidence}} \equiv 95\%$

1 - confidence

Types of Polling Problems : Specify Specs  
w/ 95% confidence of  $\leq 1\%$  error

- i) Collect data ( $n$ ) , calculate sample mean  $M_n$   
& check whether we satisfy the specs
- ii) Calculate the sample size  $n$  so that the specs are satisfied.

---

Check the polling ex. from last time :

→  $n = 50K$  people <sup>/samples</sup> were needed to satisfy the specs.  
Not practical.

→ 1-2 K , samples → you need to change your specs.

# ★ Different scalings of $X_1 + X_2 + \dots + X_n$ , $X_i$ i.i.d. w/ $\mu, \sigma^2$ .

i) Sample Mean r.v.  $\text{Scale} = \frac{1}{n}$   $\rightarrow M_n = \frac{X_1 + \dots + X_n}{n}$  (sum r.v.)

$E[M_n] = \mu$

$\text{Var}(M_n) = \frac{\sigma_x^2}{n}$  as  $n \uparrow$  gets narrower

ii)  $\text{Scale} = 1$  Sum r.v.  $S_n = X_1 + \dots + X_n$ ;  $X_i$  i.i.d.

$E[S_n] = n \cdot \mu$ ,  $\text{Var}(S_n) = n \cdot \sigma_x^2$

$n \uparrow$  narrower distrib. wider wider flatter out in the limit

density shifts

iii)  $\text{Scale} = \frac{1}{\sqrt{n}}$   $\frac{S_n}{\sqrt{n}} \rightarrow \text{Var}\left(\frac{S_n}{\sqrt{n}}\right) = \text{Var}\left(\frac{X_1 + \dots + X_n}{\sqrt{n}}\right) = \left(\frac{1}{\sqrt{n}}\right)^2 n \cdot \sigma_x^2 = \sigma_x^2$

Shape of the distrib. may change

constant variance. distrib. width stays the same!

Different scalings of  $S_n$  :  $X_1, X_2, \dots, X_n$  i.i.d. w/  $\mu$ ,  $\sigma^2$ .  
mean  $\mu$ , var  $\sigma^2$ .

3 variants : 1)  $S_n = X_1 + \dots + X_n$  : Variance  $n\sigma^2 \rightarrow \infty$ .

2)  $M_n = \frac{S_n}{n}$  : Variance  $\frac{\sigma^2}{n} \rightarrow 0$ .

3)  $\frac{S_n}{\sqrt{n}}$  : Constant Variance  $\sigma^2$  ✓

STANDARDIZE  $S_n = X_1 + \dots + X_n$  by

$$Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}} = \frac{S_n - n \cdot E[X]}{\sqrt{n\sigma_x^2}} = \frac{S_n - nE[X]}{\sqrt{n} \cdot \sigma}$$

$\sigma_{S_n}$  : standard deviation  
 $\sqrt{n\sigma_x^2} = \sqrt{n} \cdot \sigma_x = \sqrt{n} \cdot \sigma$

This is called standardization of an r.v.  $\rightarrow$  Zero mean, unit variance.

now

$$\left\{ \begin{array}{l} E[Z_n] = 0 \\ \text{Var}(Z_n) = 1 \end{array} \right.$$

Compare  $Z_n$  to a standard Normal r.v.

$$\underbrace{Z_n}_{\text{cdf of } Z_n} \xrightarrow{\text{(in distribution)}} \underbrace{Z}_{\text{cdf}} \sim \mathcal{N}(0, 1)$$

Let  $Z$  be a standard normal r.v.

# CENTRAL LIMIT THEOREM (CLT) : statement about CDFs.

We have a large # i.i.d. r.v.s  $X_i$ 's ( $X_i$ 's could be any r.v., any distribution!)  
→ Sum them ( $S_n$ ) → Standardize them to  $Z_n$   
 $Z_n = \frac{S_n - E(S_n)}{\sqrt{n}\sigma}$  n.m.

For every  $c$ ,

$$P(Z_n \leq c) \rightarrow P(Z \leq c) = \Phi(c)$$

cdf of  $Z_n$

useful b/c we have

standard Normal  
≡ Gaussian w/ 0 mean & unit var  
CDF of std normal available from tables.

Given

$$X_1, X_2, \dots, X_n$$

$$\rightarrow \text{normalized } \sum X_i$$

$X_i$ 's w/ mean  $\mu$  var:  $\sigma^2$  finite.

(Distrib. CDF)

Normal Distrib. CDF

$$Z_n = \frac{S_n - nE(X)}{\sqrt{n}\sigma}$$

$$\rightarrow S_n = \sqrt{n}\sigma Z_n + nE(X)$$

is Normal r.v.

$Z_n \sim$  Normal when  $n$  is large

Exercise

Check why  $\bar{S}_n$  is also a normal r.v.

(b/c  $S_n \leftarrow$  linear transformation of  $Z_n$ )

we assume pretend  $Z_n$  is Normal r.v. when n is large enough.

- CLT is a limit theorem:  $n \rightarrow \infty$ .

- In practice, maybe  $n=30 \rightarrow$  gives accurate results.

Ex: Let  $X_i \sim N(0,1)$  i.i.d.

Approximate  $Y_N$  distrib. by a Gaussian.

Define

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

Q. Is this justified?  
Yes

b/c  $X_i$  are i.i.d. (w/ finite mean & var)

$$\text{Var}(Y_N) = N \text{Var}(X_i^2)$$

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



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

$$\text{Var}(X^2) = E[X^4] - (E[X^2])^2 = 3 - 1 = 2.$$

$\underbrace{X^2 = Y}_{Y=X^2} \quad \underbrace{E[X^4]}_{E[Y^2]} - \underbrace{(E[X^2])^2}_{(E[Y])^2}$

Kurtosis of  $N(0,1)$   $E[X^4] = 3$  (check) derived & known ✓

$$\underline{\underline{Y_N}} = \frac{Y_N - N \cdot 1}{\sqrt{2N}} \approx \underline{\underline{N(0,1)}} \text{ by CLT.}$$

as  $n \rightarrow \infty$ .

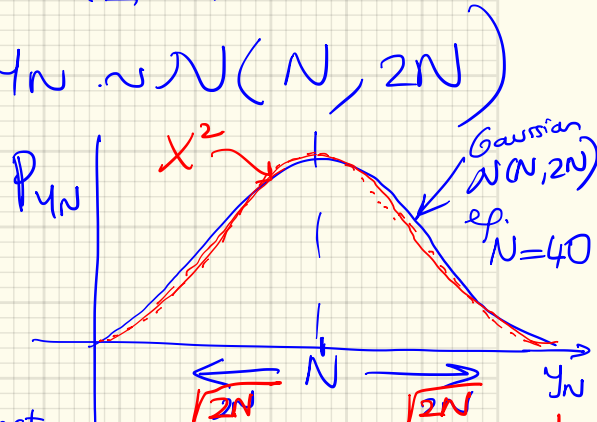
$$\rightarrow Y_N = \sqrt{2N} \tilde{Y}_N + N \rightarrow Y_N \sim N(N, 2N)$$

Gaussian? Yes Gaussian  $\rightarrow$  w/ mean = 0, std = 1

$$E[Y_N] = \sqrt{2N} E[\tilde{Y}_N] + N \quad \checkmark$$

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

$Y_N$  is  $\chi^2$ : chi-squared distrib.  $Y_N \sim \chi^2$  exact.



as  $N \uparrow$  approx. of chi-squared by a Gaussian becomes better.

# Ex: Pollster's problem using CLT

- $p$ : fraction of the population that prefers "something"
- $X_i$ :  $i$ th randomly selected person,  $X_i \stackrel{\text{Bernoulli}}{=} \begin{cases} 1, & \text{if yes} \\ 0, & \text{if no.} \end{cases}$

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

- We define 2 specifications for the poll  $\equiv$  2 parameters

$P(|M_n - p| \geq \underbrace{0.01}_{\text{accuracy}}) \leq 0.05 = 1 - \underbrace{\text{confidence}}_{95\%}$

~~Want~~  $\equiv$  Want probability 95% that our estimate  $M_n$  is within 1% of the true  $p$  value

event of interest

rewrite

$$|M_n - p| \geq 0.01$$

standardizing:

$$= \left| \frac{X_1 + \dots + X_n - np}{n} \right| \geq 0.01$$

$$= \left| \frac{X_1 + \dots + X_n - np}{\sqrt{n} \sigma} \right| \geq \frac{0.01 \sqrt{n}}{\sigma}$$

Standardized r.v.  $\equiv Z_n$ .

$$\sigma_{M_n}^2 = \frac{\sigma^2}{n}$$

$$\sigma_{M_n} = \frac{\sigma}{\sqrt{n}}$$

divide by this.

$$\rightarrow P(|M_n - p| \geq 0.01) = P(|Z_n| \geq \frac{0.01\sqrt{n}}{\sigma})$$

$$\approx P(|Z| \geq \frac{0.01\sqrt{n}}{\sigma})$$

result of the ball  $\downarrow$   
 the value  $\downarrow$

Using CLT, we calculate this prob

where  $Z$  is a standard normal r.v.:  $Z \sim N(0,1)$

$\rightarrow$  Note one difficulty: we don't know  $\sigma$ ? But we know an upper bound

For Bernoulli:  $\text{Var}(X) = p(1-p)$ . Recall last time:  $\sigma^2 \leq \frac{1}{4}$



$\rightarrow \sigma \leq \frac{1}{2}$   $\rightarrow$  we'll use this upper bound.

$$= P(|Z| \geq \frac{0.01\sqrt{n}}{\sigma}) \leq P(|Z| \geq 0.02\sqrt{n})$$

$$\left( \frac{1}{\sigma} \geq 2 \right)$$

lower bound.

(i) Given  $n = 10,000$

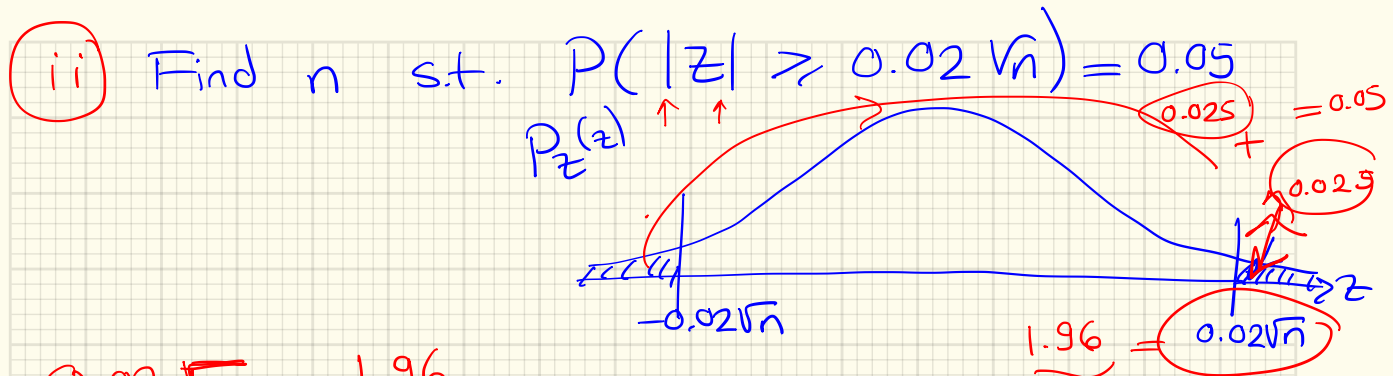
$$P(|Z| \geq 0.02\sqrt{10^4}) = P(|Z| \geq 2) = 2P(Z \geq 2)$$

$$= 2(1 - \underbrace{P(Z < 2)}_{\Phi(2)}) = 2(1 - 0.9772)$$

$$= 0.0456 \stackrel{\text{prob. of error}}{<} 5\%$$

Use Normal Table  $\approx \%4.6$





$$0.02\sqrt{n} = 1.96$$

$$\rightarrow n = 9604$$

w/  $n = 9604$  people in the poll  
our prob. of error is 0.05

CLT: in polling

- i) start w/  $n$  & calculate probabilities
- ii) start w/ specs (probs) & calculate  $n$ .

$$\Phi(c) = 1 - 0.025$$

$$\Phi(c) = 0.975$$

$\rightarrow c = 1.96$  from the table

we use  
CLT in different ways

- ① polling
- ② approximating distribution by a standard Gaussian.

Ex: CLT: Apply to Binomial approx.

$X_i$ : Bernoulli ( $p$ ), i.i.d.  $0 < p < 1$ .

$$S_n = X_1 + \dots + X_n \equiv \underline{\text{Binomial}(n, p)}$$

Binomial r.v.:  $\frac{\text{mean: } np}{\text{Variance: } np(1-p)} \checkmark$

CDF of  $\frac{S_n - np}{\sqrt{np(1-p)}}$   $\xrightarrow{\text{CLT}}$  Std Normal Distrib.  
Standardized

Check whether this approx is good

Let  $n=36$ ,  $p=0.5$

Find.

$$P(S_n \leq 21) ?$$

$$\begin{aligned} \text{mean} &= np = 18 \\ \text{Var} &= np(1-p) = 9 \rightarrow \sigma = 3 \end{aligned}$$

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

$$Z_n \leq 1$$

$\approx Z \leq 1$   $\xrightarrow{\text{CLT}}$   
std normal.

from the table.

$$\rightarrow \Phi(1) = P(Z \leq 1) = \boxed{0.843} \quad \text{approx. answer?}$$

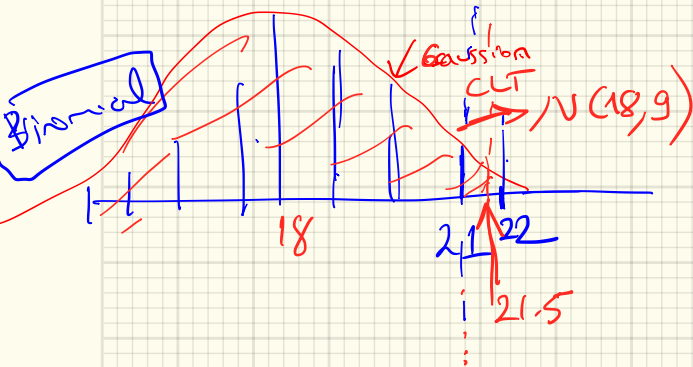
Exact answer:  $\sum_{k=0}^{21} \binom{36}{k} \left(\frac{1}{2}\right)^{36} = \boxed{0.8785}$

$P(S_n = k)$

b/c  $S_n$  is a discrete r.v.

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

$\therefore$  use  $P(S_n \leq \underline{21.5})$   $\rightarrow$  called  $\frac{1}{2}$  correction for Binomial approx.



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

Table:  $P(Z \leq 1.17)$

$$= \boxed{0.879}$$

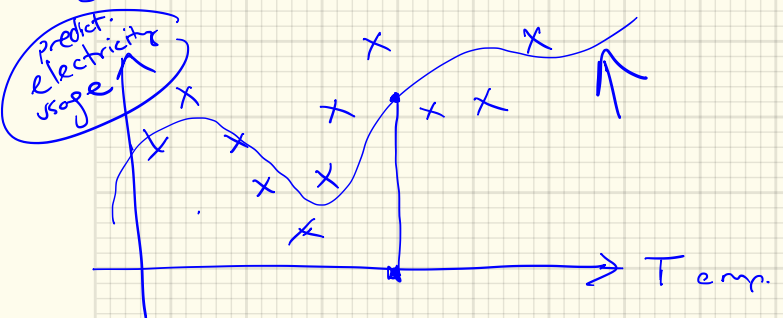
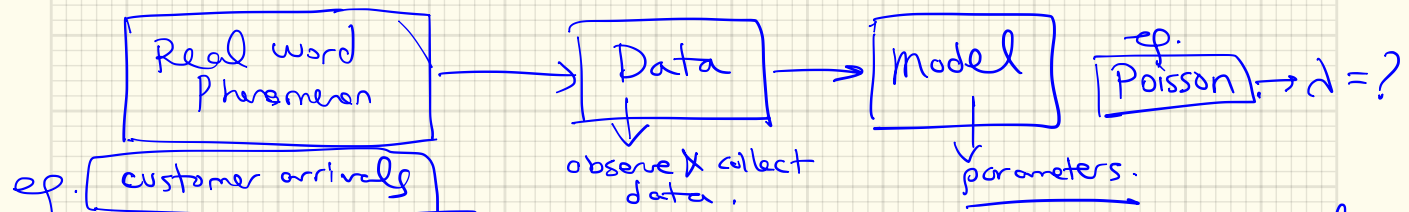
a much better approx.

w/ the

$\frac{1}{2}$  correction.

OK but not good!

# STATISTICAL INFERENCE ~ Applied Probability



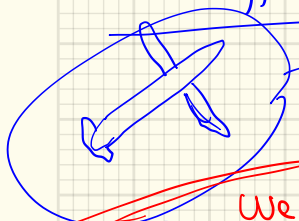
- come up w/ a model
- calculate its parameters
- make predictions about the real world.

— Polling Problems :  
estimation of preferences of populations

— Finance ..

- 1) Bayesian Statistical Inference
- 2) Classical Statistical Inference

Two types of problems (i) Hypothesis Testing: discrete quantities.



Radar measurement

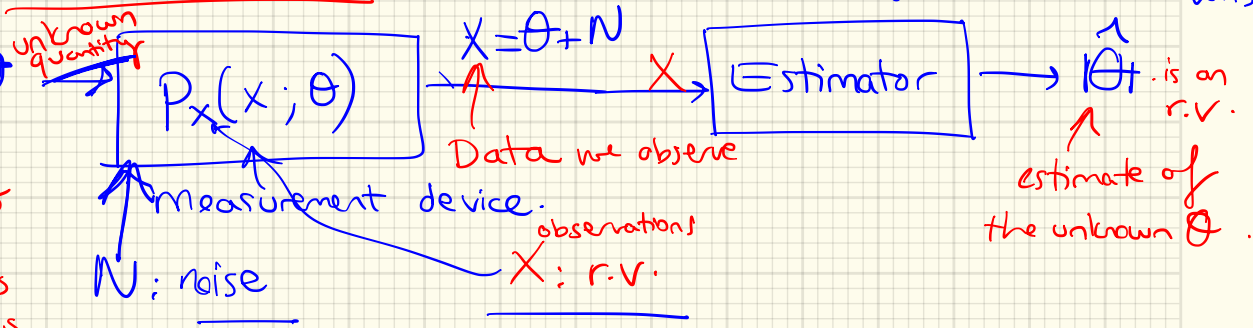
→ Detect a signal or not

interested in prob. of error.

We start w/ problem

(ii) Estimation: Want to measure a voltage  $\theta$ , in (V, V) volts

1st approach to Estimation  
we try to estimate  $\theta$   
In Estimation Unknown is a continuous quantity



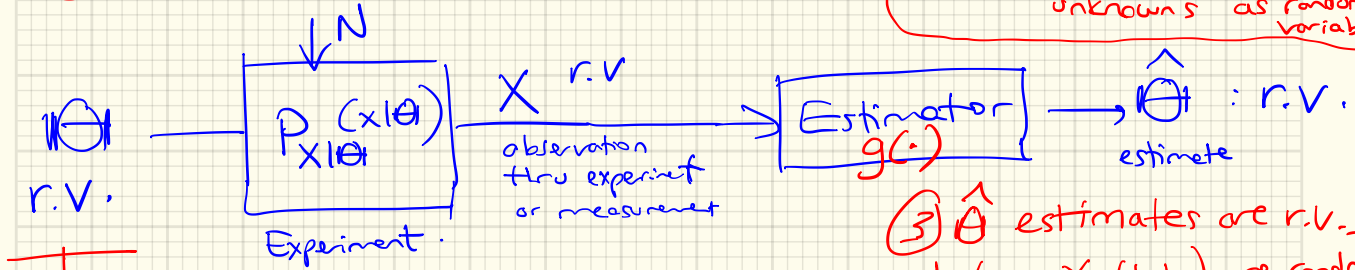
1 Classical Statistics approach

$\theta$ : unknown number; not a r.v.!



## ② Bayesian Statistics Approach:

- Note ① Classical stat. treats unknowns as unknown numbers
- ② Bayesian Stats treats unknowns as random variables



- ③  $\hat{\theta}$  estimators are r.v.s. b/c.  $X$  (data) are random variables

$P_{\theta}(\theta)$  : prior distribution:  
~~\*~~ our initial belief about  $\theta$  before the experiment

$\hat{\theta} = g(X)$   
 Recall: functions of r.v.s are r.v.s.

We rely on Bayes rule  $\rightarrow$  come up w/ a posterior distrib.

$$P_{\theta|X}(\theta|X)$$

: our revised beliefs once we obtain data  $X$ , i.e. after the experiments

~~\*~~ prior  $\rightarrow$  posterior ~~\*~~

# Bayesian Inference

conditional model of the experimental process = likelihood of observing the data

$$P(\theta | x) = \frac{P_{x|\theta}(x|\theta) \cdot P_{\theta}(\theta)}{P_x(x)}$$

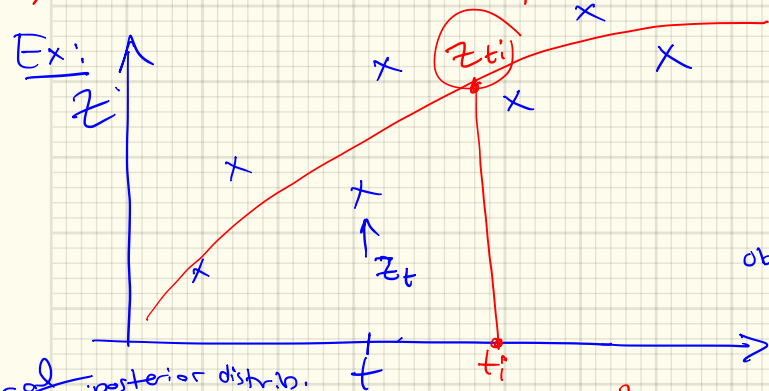
Prior: initial belief on  $\theta$  before the experiment

Valid for  $P$ : pmfs or pdfs ✓  
both discrete & continuous r.v.s.

evidence: can be calculated based on likelihood

$$\int_{\Theta} P_{x|\theta}(x|\theta) P_{\theta}(\theta) d\theta$$

Want to calculate posterior distribution of  $\theta$  given  $X$ .



position of the object.

$$z_t = \theta_0 + \theta_1 t + \frac{1}{2} \theta_2 t^2$$

$$X_t = Z_t + W_t$$

↑ observations      ↑ noise.

$$\underline{\theta} = [\theta_0, \theta_1, \theta_2]^T$$

(measurements we make are noisy).

Goal: Estimate  $\theta_0, \theta_1, \theta_2$

Goal: calculate posterior distrib.

$$P(\theta_0, \theta_1, \theta_2 | x_1, \dots, x_n)$$

Bayes ↓

$$= \frac{P(x_1, \dots, x_n | \underline{\theta}) P_{\theta}(\underline{\theta})}{P(x_1, \dots, x_n)}$$

likelihood

prior is given continuous (pdfs.) r.v.

Ex: Coin w/ unknown parameter  $\theta$  ( $\stackrel{\text{exp.}}{=} p$ ).  
 Estimate  $\theta$ .

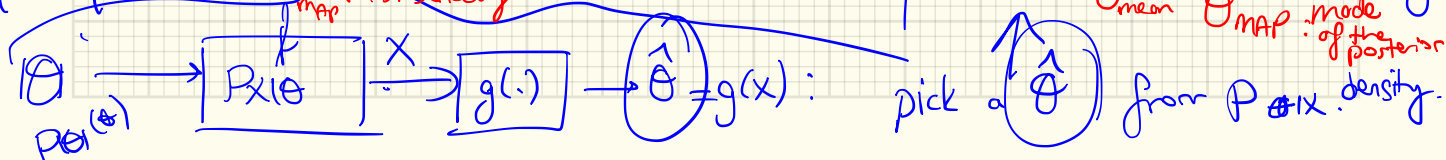
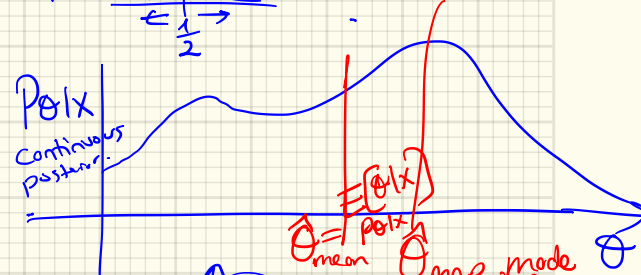
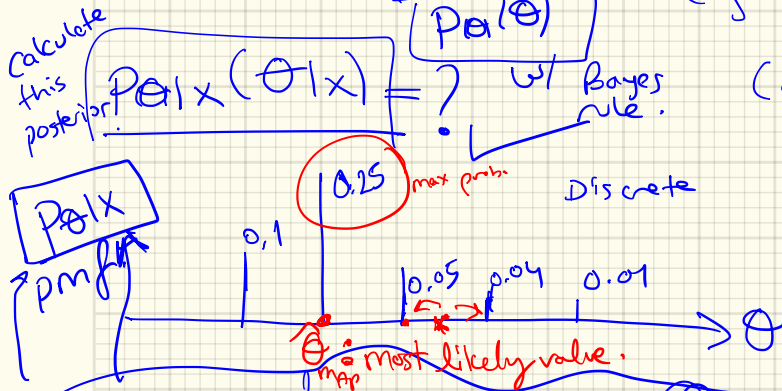
→ Classical statistics approach: Flip the coin many times

$$\hat{\theta}_m = \frac{S_n}{n}$$
 where  $S_n$  is the sum of  $X_i$ 's (heads in  $n$  trials) and  $n$  is the number of trials.

→ Bayesian approach: Assume a prior on  $\theta$ .

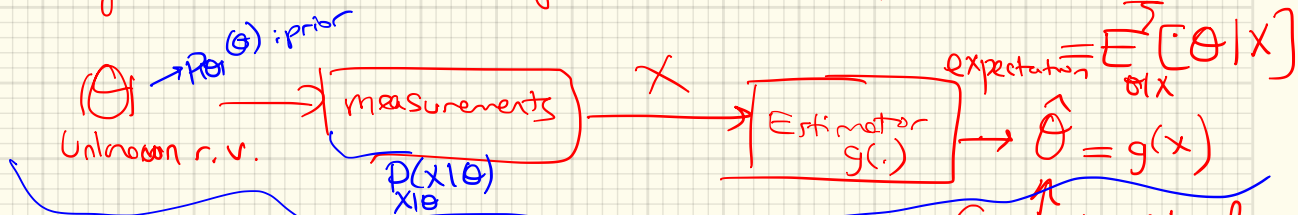
(eg. uniform prior if you don't know anything about  $\theta$ ).

(eg. another prior  $N(\mu, \sigma)$ )



→ Output of Bayesian inference :  $P_{\theta|X}(\theta|X)$  : whole distribution, <sup>posterior</sup>

→ If interested in a single output :  $\hat{\theta}_{MAP}$  or  $\hat{\theta}_{MLE}$ .



Once we estimate the posterior distrib  $P_{\theta|X}(\theta|X)$  : complete answer to the Bayesian inference problem. <sup>can be picked from the posterior.</sup>

point estimates of  $\theta$

(1) MAP: (Maximum a Posteriori) Estimate

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

(2) Conditional Expectation (Least Mean Squared <sup>Error</sup> LMS estimate)

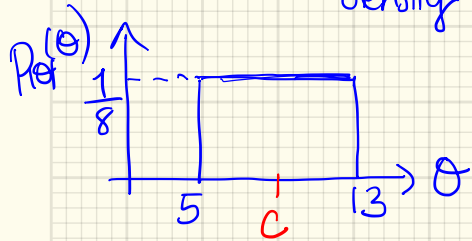
$$E[\theta|X=x] = \int \theta P_{\theta|X}(\theta|X) d\theta$$

: average of the posterior distrib.

# Least Mean-Squared Estimation (LMS):

Given a prior on  $\theta$  :  
density

Goal : come up w/ a point estimate  
using LMS.



(w/ no observed  $x$ ) we minimize LMS  
criterion:

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

This is called the  
Least Mean-Squared Estimation.

What is "optimal estimate  $c$ ?"  
in LMS sense

$$\min \frac{\partial}{\partial c} E[\theta^2] - 2E[\theta] \cdot c + c^2 = 0$$

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

$$\rightarrow c = E[\theta] = 9 \text{ for this specific example.}$$

"optimal" mean-squared error: To judge how good is our estimate?

$$E[(\theta - E[\theta])^2] = \text{Var}(\theta)$$

★ When we're minimizing the least mean-squared error,  
Expectation is the "best" estimate →

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

LMS estimation of  $\theta$  based on  $X$

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

is minimized by

$$c = E[\theta | X=x]$$

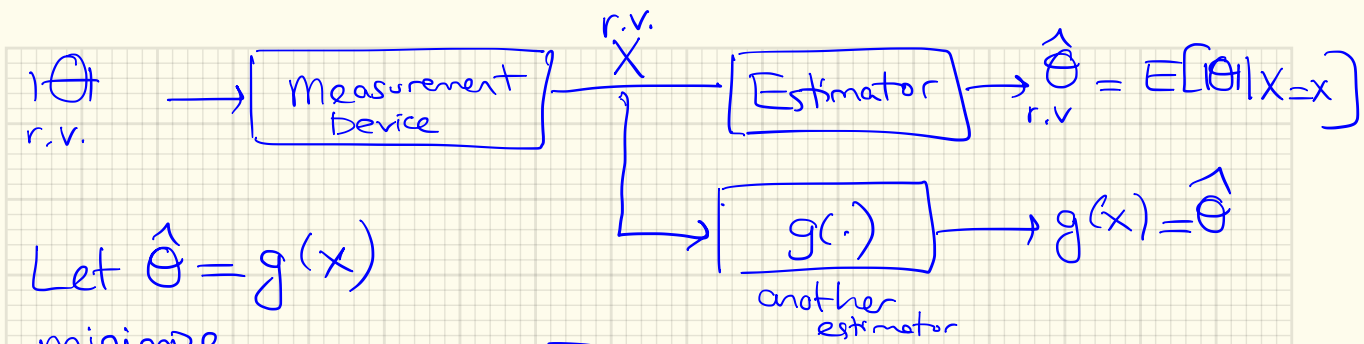
now we use  
conditional expectation  
(exp. of the posterior density).

: Given  $X$  r.v.

$\theta, X$  are r.v.s,

Given  $X=x$ ,

Evaluate in a  
conditional  
universe.



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

minimize

$$\min_{X, \Theta} E[(\Theta - g(X))^2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\Theta - g(x))^2 P_{X, \Theta}(x, \Theta) dx d\Theta$$

$$= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} \underbrace{(\Theta - g(x))^2}_{c} P_{\Theta|X}(\Theta|x) d\Theta \right) \underbrace{P_X(x) dx}_{P_X(x) \geq 0} \checkmark$$

We can minimize

$$c = g(x) \quad \frac{\partial}{\partial c} \int_{-\infty}^{\infty} (\Theta - c)^2 P_{\Theta|X}(\Theta|x) d\Theta = f(c) \quad \uparrow \text{unknown.}$$

$$\frac{\partial f}{\partial c} = 0 \quad \frac{\partial}{\partial c} \int_{-\infty}^{\infty} (\Theta - c) P_{\Theta|X}(\Theta|x) d\Theta = 0$$

$$\rightarrow \int_{-\infty}^{\infty} \Theta P_{\Theta|X}(\Theta|x) d\Theta = c \cdot \int_{-\infty}^{\infty} P_{\Theta|X}(\Theta|x) d\Theta \Rightarrow$$

$$c = E[\theta | X] = \int \theta P_{\theta|X}(\theta|x) d\theta$$

$g(x)$

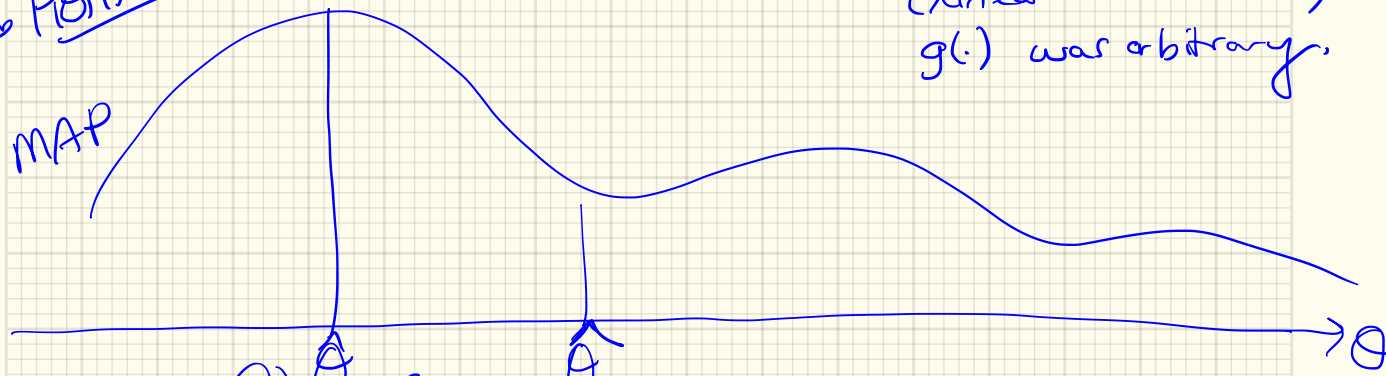
$$\hat{\theta}(x) = E[\theta | X]$$

is the "optimal" estimator  
(in sense of LMSE)

among all estimators  
(linear or nonlinear)  
 $g(\cdot)$  was arbitrary.

posterior dist.  $P_{\theta|X}(\theta|x)$ : complete answer to a Bayesian inference.

① MAP



①  $\theta_{MAP}$

②  $\theta_{LMSE} = E[\theta | X]$ : conditional mean of the posterior density



# Issues

- Unknown  $\theta$  : vector of  $\theta$ 's.
- Observations  $X = (X_1, \dots, X_n)$   $\leftarrow$  several measurements

Calculate w/ Bayes rule

$$P_{\theta | X_1, \dots, X_n}(\theta | X_1, \dots, X_n) \rightarrow \text{LMSE} \quad E[\theta | X_1, \dots, X_n]$$

$$\rightarrow P_{\theta | X_1, X_2, \dots, X_n} = \frac{P_{X_1, \dots, X_n | \theta} P_{\theta}}{P_{X_1, \dots, X_n}}$$

$$P_{X_1, \dots, X_n} = \iiint P_{X_1, \dots, X_n | \theta} P_{\theta}(\theta)$$

- 1) Calculations may become intractable  $\leftarrow$  multi-dimensional integrals  $\leftarrow$   $P_{\theta | X}$  & its expectation
- 2) Come up w/ a plausible prior

Ex 8.11 from Bertsekas book.

## Linear LMS estimator:

Now, consider a simpler estimator of  $\theta$ :

Let  $g(x) = \hat{\theta}_L = ax + b$  : affine mapping w/  
parameters  $a$  &  $b$   
unknown.

$$\begin{aligned} \text{minimize } & E[(\theta - \hat{\theta}_L)^2] \\ & = E[(\theta - \underbrace{(ax+b)}_{g(x)})^2] \end{aligned}$$

generic  
estimator

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

exercise Derive:

$$\begin{aligned} & E[\theta^2 - (ax+b)^2 - 2\theta(ax+b)] \\ \text{minimize w.r.t } & \left\{ E[\theta^2] + a^2 E[X^2] + 2ab E[X] + b^2 - 2a E[\theta \cdot X] - 2b E[\theta] \right\} \\ & \frac{\partial}{\partial a} = 0, \quad \frac{\partial}{\partial b} = 0 \quad \checkmark \end{aligned}$$

$$a = \frac{\text{cov}(X, \theta)}{\text{Var}(X)}, \quad b = E[\theta] - \frac{\text{cov}(X, \theta)}{\text{Var}(X)} E[X]$$

"Best" choice of  $a$  &  $b$  w.r.t. LMS criterion:

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

Linear LMS estimator of  $\theta$ .