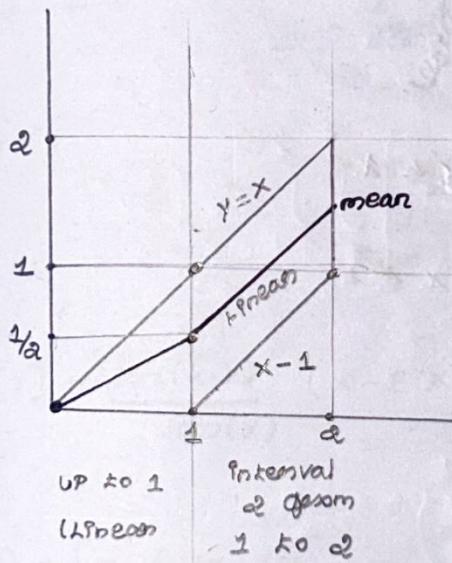
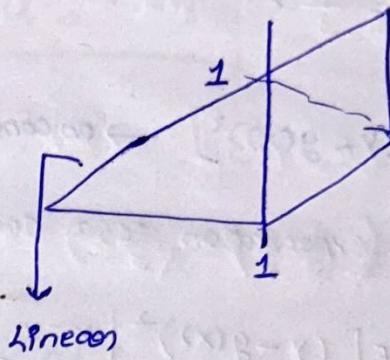


1) find $g(x)$ [LMS estimators]

$$g(x) = E[y|x]$$

up to $x=1$ (linear)



2 linear points.

* UP to 1

* from to 2) mean curve
(Linear)

$$E[y|x] = \begin{cases} \frac{1}{2}x & 0 \leq x \leq 1 \\ x - \frac{1}{2} & 1 \leq x \leq 2 \end{cases}$$

$$b) E[(y-g(x))^2 | x=x] \rightarrow \text{Conditional } \begin{pmatrix} \text{mean variance} \\ \text{variance in estimate} \end{pmatrix}$$

↓ ↓
Actual Estimate

$$E[(y - E[y|x])^2 | x=x] = \text{Var}(y|x)$$

$$\text{Var}(y|x) = \frac{\text{width}}{12} = \frac{(b-a)^2}{12}$$

$$\left\{ \begin{array}{l} \left(\frac{x^2}{12} \right), \quad 0 \leq x \leq 1 \\ \frac{1}{12}, \quad 1 \leq x \leq 2 \end{array} \right.$$

$1 \leq x \leq 2$

rotated width: $\frac{(x-x-1)^2}{12} = \frac{1}{12}$

$[x|y] = (x) \frac{1}{12}$

c) $E[(y - g(x))^2] \rightarrow \text{unconditional} \quad [\text{total expectation theorem}]$

Iterative expectation of a conditional exp \rightarrow unconditional exp

$$= \int E[(y - g(x))^2 | x=x] f_x(x) dx$$

↓
over all x [Jpdf]

Fix x , estimate y
 ↓
 estimate x | continue

$$= \int E[(y - g(x))^2 | x=x] f_x(x) dx$$

$$f_x(x) = \begin{cases} \int_0^x \frac{2}{3} dy & 0 \leq x \leq 1 \\ \int_{x-1}^x \frac{2}{3} dy & 1 \leq x \leq 2. \end{cases}$$

$$= \begin{cases} \frac{2}{3}x & 0 \leq x \leq 1 \\ \frac{2}{3} & 1 \leq x \leq 2. \end{cases}$$

$$E[(y - g(x))^2] = \int_0^1 \frac{x^2}{12} \cdot \frac{2}{3} x dx + \int_1^2 \frac{1}{12} \cdot \frac{2}{3} dx$$

$$= \int_0^1 \frac{x^3}{12} \cdot \frac{2}{3} dx + \int_1^2 \frac{1}{12} \cdot \frac{2}{3} dx$$

$$= \frac{5}{72}$$

$$= E[\text{Var}(Y|x)] \quad \text{same}$$

$$= E[E[(y - g(x))^2 | x]]$$

$$E\{ (Y - g(X))^2 \} = E\{ \text{var}(Y|X) \}$$

$$E\left[E\left[(Y - g(X))^2 | X=x\right]\right] = E\left[(Y - g(x))^2\right]$$

$$E\left[(Y - g(x))^2\right] = E\left[(Y - E[Y|X])^2\right]$$

$$= \text{var}(Y|X)$$

$$= E[\text{var}(Y|X)] = \text{var}(Y|X)$$

$$E[\text{var}(Y|X)] = E[(Y - g(X))^2]$$

LMS estimator \rightarrow Forces estimator to have 0 mean m.s.

d)

$$\lambda(x) = ?$$

One case

Two linear pieces (kink)

$$\lambda(x) = E[Y] + \frac{\text{cov}(x, Y)}{\text{var}(x)} (x - E(x))$$

$$E[X] = \int x \cdot f_X(x) dx$$

$$= \int_0^1 x \cdot \left(\frac{2}{3}x\right) dx + \int_1^2 x \cdot \left(\frac{2}{3}\right) dx$$

$$= \frac{2}{3} \left(\frac{x^3}{3}\right)_0^1 + \left(\frac{2}{3}\right) \left(\frac{x^2}{2}\right)_1^2$$

$$= \frac{2}{9} (1) + \frac{4 - 1}{3}$$

$$= \frac{2}{9} + 1$$

$$= \frac{11}{9}$$

$$E[X^2] = \int x^2 f_X(x) dx$$

$$= \frac{2}{3} \left(\frac{x^4}{4}\right)_0^1 + \left(\frac{2}{3}\right) \left(\frac{x^3}{3}\right)_1^2$$

$$= 31/18$$

$$\text{Var}(x) = E[x^2] - (E[x])^2 = \frac{31}{18} - \left(\frac{11}{9}\right)^2$$

$$= \frac{37}{162}$$

$$\begin{aligned} E[y] &= E[E[y|x]] = \int E[y|x] f_x(x) \\ &= \int_0^1 \left(\frac{1}{2}x \right) \cdot \left(\frac{2}{3}x \right) dx + \int_1^2 \left(x - \frac{1}{2} \right) \left(\frac{2}{3} \right) dx \\ &= \frac{7}{9} \end{aligned}$$

$$\text{Cov}(x, y) = E[xy] - E[x]E[y]$$

$$\begin{aligned} E[xy] &= \int_x \int_y xy f_{x,y}(x,y) dy dx \\ &= \left[\int_0^1 \int_0^x xy \left(\frac{2}{3} \right) dy dx \right] + \left[\int_1^2 \int_{x-1}^x xy \left(\frac{2}{3} \right) dy dx \right] \\ &= \frac{41}{36} \end{aligned}$$

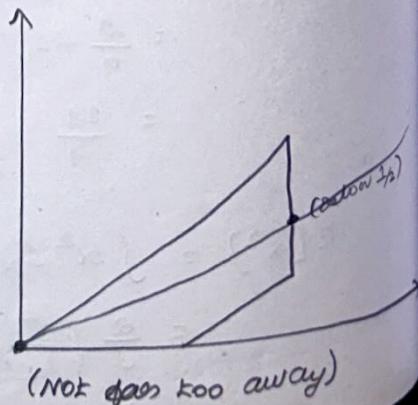
$$\text{Cov}(x, y) = \frac{41}{36} - \left(\frac{11}{9}\right)\left(\frac{7}{9}\right)$$

$$\boxed{\text{Cov}(x, y) = \frac{61}{324}}$$

$$L(x) = E[y] + \frac{\text{Cov}(x, y)}{\text{Var}(x)} (x - E[x])$$

$$= \frac{7}{9} + \frac{\frac{61}{324}}{\left(\frac{37}{162}\right)} \left(x - \frac{11}{9}\right)$$

$$L(x) = \frac{7}{9} + \frac{61}{74} \left(x - \frac{11}{9}\right) \rightarrow \text{plot}$$



LMS estimators: $g(x)$

Linear LMS estimator: $L(x)$

which one is better

$$E[(y - L(x))^2] \geq E[(y - g(x))^2]$$

↓
designed to minimize the error.

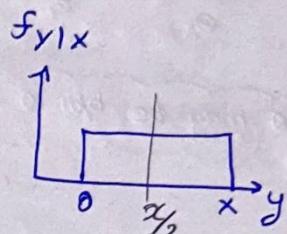
we have made approximation.

MAP estimation - pick out highest value
of y in distribution.

Our case: Uniform.

• No sensible way of choosing only one?

↓
cond. expectation is midpoint



• Every single point gives same distribution (conditional))

i) $\hat{\theta}$, conditional mean bivariate error

$$E[(\hat{\theta} - \theta)^2 | x = x] = \int_x^1 (\hat{\theta} - \theta)^2 f_{\theta|x}(\theta|x) d\theta$$

Another problem.

Recitation 2

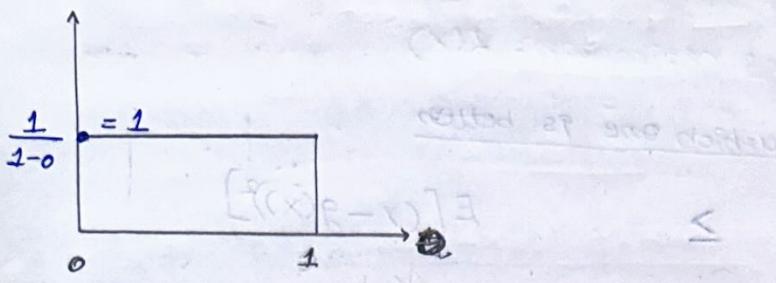
Updating a parameter of uniform

Romeo & Juliet start dating, but Juliet will be late by
date by a R.V x , uniformly dist $[0, 1]$. The parameter
 $\theta \rightarrow$ unknown \rightarrow modelled as the value of a R.V Θ , uniformly
distributed b/w zero & 1 hour.

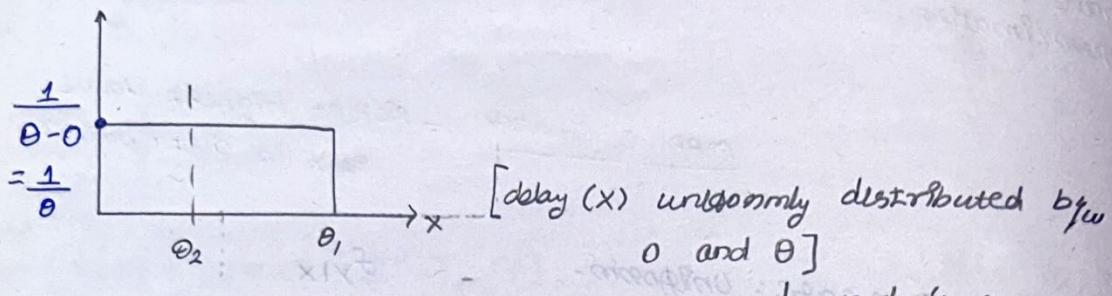
Ans:

a) Assume: Juliet was late by an amount x on their
first date, how should Romeo use this info to update the
distribution of θ ?

Ans:



distribution of θ



[delay (x) uniformly distributed b/w 0 and θ]

↳ we don't know θ .

$\therefore \theta$ may be b/w 0 and 1

⇒ players of uncertainty

$$f_{\theta}(\theta) = \begin{cases} 1, & 0 \leq \theta \leq 1 \\ 0, & \text{anywhere else} \end{cases} \rightarrow \text{prior}$$

$$f_{x|\theta}(x|\theta) = \begin{cases} \frac{1}{\theta}, & 0 \leq x \leq \theta \\ 0, & \text{anywhere else.} \end{cases} \rightarrow \text{Data.}$$

$$f_{\theta|x}(\theta | x) = \frac{f_{\theta}(\theta) f_{x|\theta}(x|\theta)}{f_x(x)} = \begin{cases} \frac{1 \cdot \frac{1}{\theta}}{f_x(x)}, & 0 \leq \theta \leq 1 \\ 0, & 0 \leq x \leq \theta \end{cases}$$

$\therefore (x \text{ can go to anywhere})$

$$f_x(x) = \int_{\theta}^1 f_{\theta}(\theta) f_{x|\theta}(x|\theta) d\theta$$

$$= \int_x^1 \left(\frac{1}{\theta} \right) \cdot 1 \cdot d\theta = \left[\log \theta \right]_x^1 = \log 1 - \log x = -\log x = \log x$$

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta} & 0 \leq \theta \leq 1 \\ |\log(\frac{x}{\theta})| & \theta < 0 \end{cases}$$

elsewhere.

$$\therefore f_x(x) \rightarrow +\infty$$

$$= \begin{cases} \frac{1}{\theta \log x} & 0 \leq \theta \leq 1 \\ 0 & \text{elsewhere.} \end{cases}$$

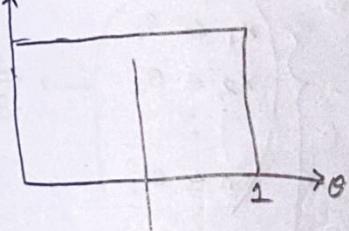
Started with posogn $\rightarrow \theta \sim [0, 1]$ uniform.
update we got

Assume 30 minutes late

$$0 \leq \theta \leq 1$$

$$\text{when } x = \frac{1}{2} \rightarrow$$

overs case



$$x = \frac{1}{2}$$

$$\theta > x > 0$$

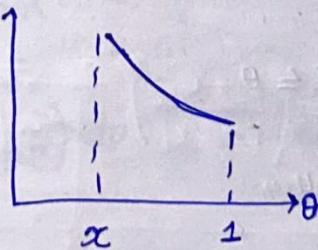
$$1 > \theta > 0$$

$$\frac{1}{\theta \log x} = (\theta/x) \cdot \theta^{-1}$$

So over delay must be $x = \frac{1}{2}$, so θ must start with $\frac{1}{2}$.

$$x \leq \theta \leq 1$$

$$f_{\theta|x}(\theta|x)$$



θ less than $x \rightarrow$ no way she arrives late by x time,

'More prob of being x' \rightarrow high prob curve.

b) more data

n samples, $x = (x_1, \dots, x_n)$ (more days)

'Conditionally independent x_1, \dots, x_n' \rightarrow Assumption.

'prior hasn't changed'

$$f_\theta(\theta) = \begin{cases} 1 & , 0 \leq \theta \leq 1 \\ 0 & , \text{ anywhere else} \end{cases}$$

$$f_{x|\theta}(x|\theta) = f_{x_1, \dots, x_n|\theta}(x_1, x_2, \dots, x_n|\theta) \quad \begin{matrix} \text{conditionally} \\ \rightarrow \text{Independent} \end{matrix}$$
$$= f_{x_1|\theta}(x_1|\theta) \cdots f_{x_n|\theta}(x_n|\theta)$$

\therefore Each one is uniformly dist, b/w 0 and θ

$$\boxed{f_{x|\theta}(x|\theta) = \frac{1}{\theta^n}} \quad \begin{matrix} 0 \leq x \leq \theta \\ 0 \leq \theta \leq 1 \end{matrix}$$

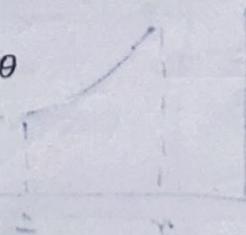
$$\left\{ \begin{array}{l} x_1 \leq \theta_1 \\ x_2 \leq \theta_2 \\ \vdots \\ x_n \leq \theta_n \end{array} \right\} \rightarrow \frac{1}{\theta^n}$$

$0 \rightarrow 0.w.$

Let find:

$$\bar{x} = \max\{x_1, \dots, x_n\}$$

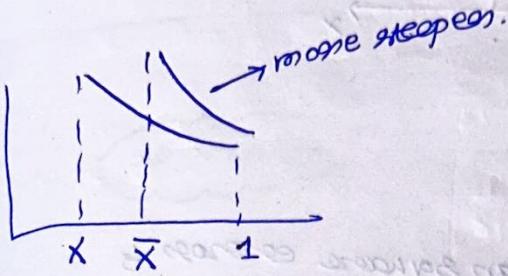
$$f_{x|\theta} = \begin{cases} \frac{1}{\theta^n}, & \bar{x} \leq \theta \\ 0 & \text{otherwise} \end{cases}$$



$$f_{\theta|\bar{x}}(\theta|\bar{x}) = \frac{f_\theta(\theta) f_{x|\theta}(\bar{x}|\theta)}{f_x(\bar{x})} = \frac{\frac{1}{\theta^n}}{f_x(\bar{x})}, \bar{x} \leq \theta$$

$$f_X(x) = \int_x^1 \frac{1}{\theta^n} d\theta$$

$$= \begin{cases} \frac{1}{n} \frac{1}{\theta^{n-1}} & , \bar{x} \leq \theta \\ \frac{1}{n} \frac{1}{\bar{x}^{n-1}} & , \text{others} \end{cases}$$



one single point

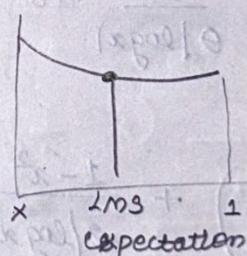
Posterior distribution: MAP

o) $\hat{\theta}_{\text{map}} = x$ (x has highest value)

↓

30 minutes late
[worst case]

on take expectation



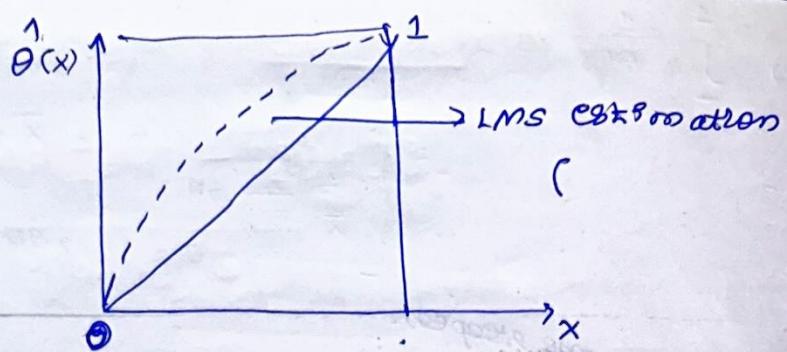
a) LMS = $E[\theta | x = x]$

$$= \int_x^1 \theta f_{\theta|x}(\theta|x) d\theta$$

$$= \int_x^1 \theta \left(\frac{1}{\log x} \right) d\theta \quad \theta / \left(\frac{1}{\log x} \cdot \theta \right) d\theta$$

$$= \int_x^1 \frac{1}{\log x} d\theta$$

$$= \frac{1-x}{\log x} \quad \rightarrow \text{LMS estimation}$$



Point-2

(i) $\hat{\theta} \rightarrow$ Initial Conditional mean square error is

$$E[(\hat{\theta} - \theta)^2 | x = x] = \int_{-\infty}^{\infty} (\hat{\theta} - \theta)^2 f_{\theta|x}(\theta|x) d\theta$$

(Value assigned to x)

$$= \frac{1}{\int_x^{\infty} (\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2) \frac{1}{\theta|\log x|} d\theta}$$

$$= (\cancel{\hat{\theta}} - \cancel{\theta})^2$$

$$= \hat{\theta}^2 (1) - 2\hat{\theta} \int_x^{\infty} \theta \frac{1}{\theta|\log x|} d\theta + \int_x^{\infty} \theta^2 \frac{1}{\theta|\log x|} d\theta$$

$$= \hat{\theta}^2 - 2\hat{\theta} \cdot \left(\frac{1-x}{|\log x|} \right) + \frac{1-x^2}{2|\log x|}$$

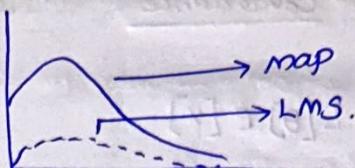
Now:

$$\hat{\theta} = x \quad \text{or} \quad \hat{\theta} = \frac{x}{|\log x|}$$

$$= x^2 - 2x \left(\frac{1-x}{|\log x|} \right) + \frac{1-x^2}{2|\log x|}$$

LMS: $\hat{\theta} = \frac{1-x}{|\log x|}$

$$\Rightarrow \frac{1-x^2}{2|\log x|} - \left(\frac{1-x}{|\log x|} \right)^2$$



LMS is better.

Linear LMS

Linear: $\theta x + b$

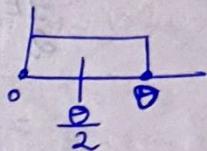
$$1 - \rho^2 \rightarrow \text{Not IP random}$$

$$\log x \rightarrow \text{Not IP random.}$$

$$\hat{\theta}_{LLMS} = E[\theta] + \frac{\text{Cov}(x, \theta)}{\text{Var}(x)} (x - E[x])$$

$$E[\theta] = \frac{1}{\alpha}$$

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



$$= E\left[\frac{\theta}{\alpha}\right]$$

$$= \frac{1}{\alpha} E[\theta]$$

$$= \frac{1}{\alpha} \left(\frac{1}{2}\right) = \frac{1}{4}.$$

$$\text{Var}(x) = E[\text{Var}(x|\theta)] \neq \text{Var}(E[x|\theta])$$

$$\text{variance} = \frac{(b-a)^2}{12}$$

$$= E\left[\frac{\theta^2}{\alpha}\right] + \text{Var}\left(\frac{\theta}{\alpha}\right)$$

$$= \frac{1}{\alpha} E[\theta^2] + \frac{1}{4} \text{Var}(\theta)$$

$$= \frac{1}{\alpha} \left[\text{Var}(\theta) + (E[\theta])^2 \right] + \frac{1}{4} \text{Var}(\theta)$$

$$= \frac{4}{12} \text{Var}(\theta) + \frac{1}{12} (E[\theta])^2$$

$$= \frac{1}{3} \text{Var}(\theta) + \frac{1}{12} (E[\theta])^2$$

$$= \frac{1}{3} \left(\frac{1^2}{12}\right) + \frac{1}{12} \left(\frac{1}{\alpha}\right)^2$$

$$= \frac{1}{3} \left(\frac{1}{12}\right) + \frac{1}{12} \cdot \frac{1}{4}$$

$$= \frac{1}{144}$$

Covariance

$$\text{Cov}(\theta, x) = E[\theta x] - E[\theta] E[x]$$

$$E[\theta x] = E[E[\theta x | \theta]] \quad \bullet$$

$$\text{Cov}(\theta, x) = E[E[\theta x | \theta]] - \frac{1}{4} \cdot \frac{1}{12}$$

$$= E[\theta E[x | \theta]] - \frac{1}{48}$$

$\therefore \theta$ is known
(constant)

$$= E\left[\theta \cdot \frac{\theta}{2}\right] - \frac{1}{48}$$

$$= \frac{1}{2} E[\theta^2] - \frac{1}{48}$$

$$E[\theta^2] = \frac{1}{12} + \frac{1}{4} = \frac{1}{3}$$

$$= \frac{1}{2} \left(\frac{1}{3} \right) - \frac{1}{48}$$

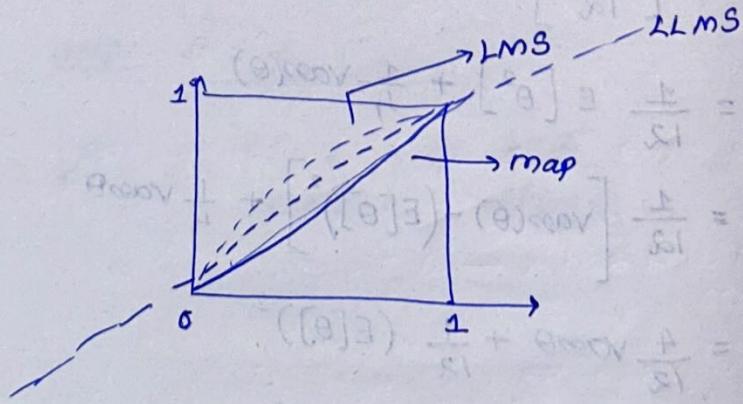
$$= \frac{1}{24}$$

$$\left[\frac{S}{S} \right]_2 =$$

LLMS

$$\hat{\theta}_{LLMS} = \frac{1}{2} + \frac{\frac{1}{24}}{\frac{7}{44}} \left(x - \frac{1}{4} \right)$$

$$\hat{\theta}_{LLMS} = \frac{6}{7}x + \frac{2}{7}$$

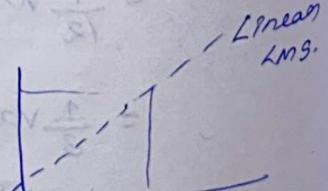


$x = 1$ (hours late)

LLMS $\rightarrow 8/7$ (greatest trans)

θ bounded by 1.

\hookrightarrow side effect of LLMS.

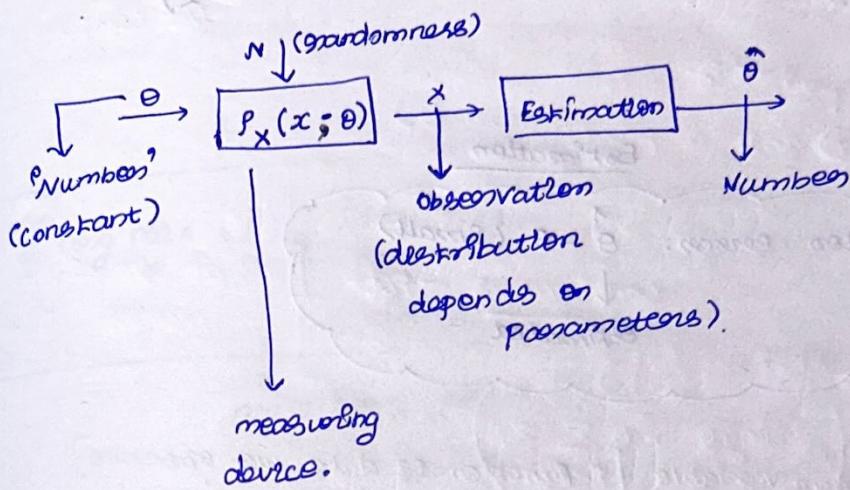


Advantage = simple

classical inference : 2 - lecture : 23

- ↓
★ don't want to assume prior distribution.
★ confidence interval.

(1) \rightarrow Classical statistics



problem: when we have two R.v (Hence only 1 R.v & a constant)

$$P_x(x; \theta)$$

↳ x has a certain distribution
(has a certain unknown parameter)

Estimation

Takes x & result $\hat{\theta}$ (unknown parameters)

$$P_x(x; \theta) \rightarrow 1\text{-dimensional quantity.}$$

→ unknown parameters in background.

$$P_{x_1 \dots x_n}(x_1, \dots, x_n; \theta_1, \dots, \theta_n)$$

* Not conditional probabilities. (θ is not random).

eg:

A normal with this mean or with that mean

(Alternative candidates of
prob model).

own job: choose which one is the correct model.

e.g.: coin - small no. of models

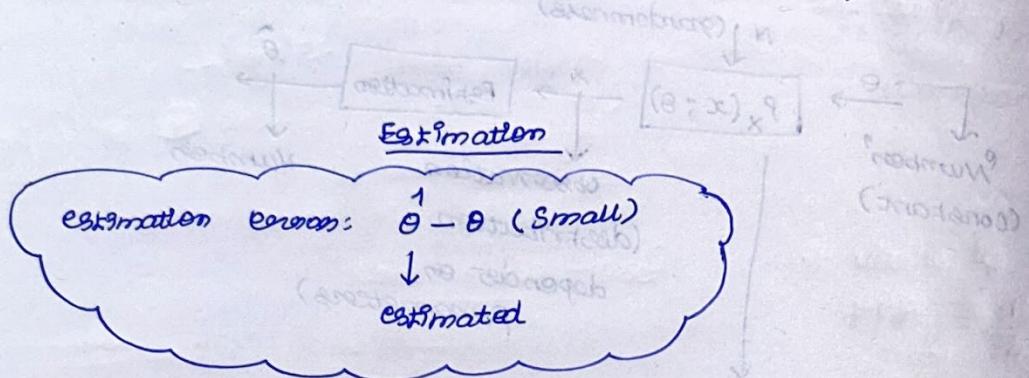
$$\text{Blas: } \frac{1}{10} \quad 0.01 \quad \frac{3}{4}$$

* decide: which is the correct model?

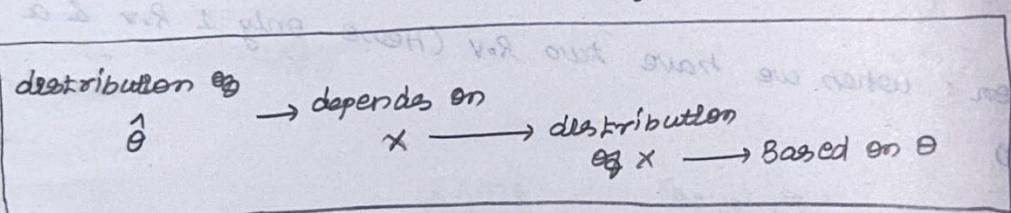
Hypotheses → choose one by estimation.

→ Blas → $(0.6, 0.7, \dots)$
→ unbiased → $(\frac{1}{2})$

lots of alternative models.



$\hat{\theta} \rightarrow$ Random variable [\because Function of data we observe]
even though constant θ .



Problem

Pick a θ , under which the data observed (x 's most likely to occur).

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

\therefore choose that $\theta \rightarrow$ which causes our observed data more likely.

[Under a particular $\theta \rightarrow x$ has this prob of occurring].

In Bayesian MAP estimation

maximum a posteriori distribution

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

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

Using ~~posterior~~ distribution θ , we will find the distribution (posterior) of θ using the data x .



In our case $\hat{\theta}_{ML} = \arg \max_{\theta} P_X(x; \theta)$

x for a part.

) opposite

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

$[\theta \rightarrow \text{Prior}]$

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

→ Bayesian.

Doesn't depend on θ .

most likely value of θ

∴ $\hat{\theta}_{MAP} \rightarrow \text{maximize}$

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

denominators don't have θ

If $P_{\theta}(\theta)$ (prior) is constant →

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

same as

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

Same form when prior is constant.

[similar]

In Bayesian: we assumed a prior
Hence: No prior

x_1, x_2, \dots, x_n , exponential (θ)

$$\max_{\theta} \prod_{i=1}^n \theta e^{-\theta x_i}$$

$\theta \rightarrow 1$ dimensional (several data)

↓
pmf
(discrete)
for a particular $\hat{\theta}_{ML} = \arg \max_{\theta} P(x; \theta)$

density.
(value of θ goes to a particular θ)

↓
Not conditional

continuous - pdf [what value of θ
makes x (observed)
more likely]

$$\max_{\theta} \left(n \log \theta - \theta \sum_{i=1}^n x_i \right)$$

maximize

(By taking log
then derivative)
w.r.t θ

$$\hat{\theta}_{ML} = \frac{n}{x_1 + \dots + x_n}$$

[ESTIMATOR] $\hat{\theta}_n = \frac{n}{x_1 + x_2 + \dots + x_n}$ → Sample mean
 \downarrow
 $(R.V) \leftarrow R.V$ (function of R.V)

In exponential distribution: $\theta = \frac{1}{\text{mean}}$

[mean = $\frac{1}{\lambda}$]
 $x = \theta$

Good or bad estimate?

$\hat{\theta}_n \rightarrow$ close to θ [less error]

Prop:

* Takes: x_1, \dots, x_n & produces $\hat{\theta}$ [random - Not exact]

* Ideally: we don't want errors [on average: correct value]

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

prob dest θ → affected by x_1 → In turn by θ .

$$E[\hat{\theta}_n] = \int f_{\hat{\theta}_n}(\hat{\theta}; \theta) \cdot \hat{\theta} d\hat{\theta}$$

$\hat{\theta} \rightarrow$ function of x .
density of x is affected by θ

\downarrow

$\therefore x_1$ affected
by θ

$\therefore \hat{\theta} \rightarrow$ goes a particular value of θ .

Under a particular θ

$$E[\hat{\theta}_n] = \theta \rightarrow \text{Irrespective of } \theta \text{ (look)}$$

No matter of $\theta \rightarrow$ we don't want to be biased.

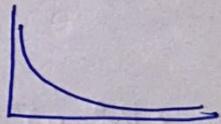
Irrespective of θ (look), $E[\hat{\theta}_n] = \theta$ [Real value]

Our estimate - Biased

Sample Size : 1

$$\theta = \frac{1}{x_1} \quad (x_1 \text{ has a less amount of density in } 0)$$

$$= \frac{1}{0} = \infty \neq \theta$$



\therefore (Biased upwards a lot)

'Under some conditions - we can make them asymptotically unbiased' - By having more data

'more data' - more correct data.

In poll problem

* More sample size - more accurate answers

'Desirable property of estimators'

more data: more accurate θ (unknown parameters)

↓
maximum likelihood estimation
(good approach in this case = more data)

$$\rightarrow \theta_n \rightarrow \theta$$

$$\hat{\theta}_n = \frac{n}{x_1 + \dots + x_n} \rightarrow \frac{1}{E[x]}$$

from weak law of convergence (large numbers)

Sample mean converges to the expectation.

$$\frac{x_1 + \dots + x_n}{n} \rightarrow E[x] = \frac{1}{\theta} \left[\frac{1}{\lambda}, \text{our case } x = \theta \right]$$

$$\frac{1}{E[x]} = \theta$$

$$\boxed{\hat{\theta}_n = \theta} \rightarrow \text{certainly correct one. [irrespective to } \theta \text{ [true data]]}$$

Small - Squared error

$$E[(\hat{\theta} - \theta)^2] = \text{var}(\hat{\theta}) + (E[\hat{\theta}] - \theta)^2$$

↓
Fix a little θ

$$= \text{var}(\hat{\theta}) + (\text{bias})_\theta^2$$

↳ depends on θ

$$\text{var}(x) = E[x^2] - (E[x])^2$$

∴ Average we fall above or below (more or less than actual) value

↓
Bias

meansquare error

→ Bias
→ $\text{var}(\hat{\theta})$

↳ estimation

) Contributions

Design: Both of the bias & $\text{var}(\hat{\theta})$ small [one's well]

Assume

$$x \sim N(\theta, 1)$$

$$\hat{\theta} = 100$$

$$E[(\hat{\theta} - \theta)^2] = 0 + (\text{Bias})^2 = (100 - \theta)^2$$

If $\theta = 0 \rightarrow \text{large error.}$

moar: ways, we can make variance small in cost of large bias
(need to take care of both)

Properties

Other method: Highschool method (say)

Estimate mean:

* x_1, \dots, x_n iid mean θ , variance σ^2

* $x_p = \theta + w_p$ (random) \rightarrow measurement

* $w_p \rightarrow$ iid, mean = 0, variance = σ^2

$$\hat{\theta}_n = M_n = \text{Sample mean} = \frac{x_1 + \dots + x_n}{n}$$

'Reasonable way'

↓
No need of variance
[No need of distribution of x]

Properties: good?

$$* E[\hat{\theta}_n] = \theta = E\left[\frac{x_1 + x_2 + \dots + x_n}{n}\right] = \frac{1}{n} n E[x] \\ = \theta \quad [E[x] = \theta]$$

* weak large numbers law: converges to true parameters
(consistency)

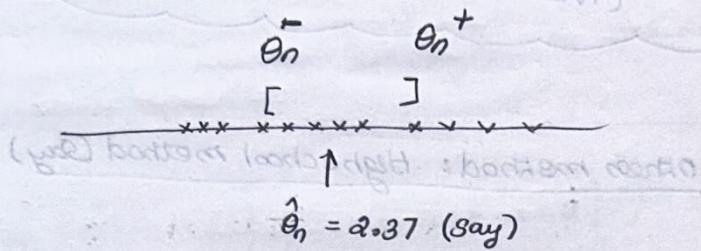
Mean square error

$$MSE = E[(\hat{\theta}_n - \theta)^2] \approx \frac{var(\hat{\theta}_n)}{n} + \left[var(\hat{\theta}_n) + E[(\hat{\theta}_n - \theta)^2] \right]$$

$\approx \frac{\sigma^2}{n}$

AS $n \rightarrow \infty$, $MSE = 0$.

Sample mean: simple, easy, nice



So we need to give some interval: how likely? → Report end points

Report $[\theta_n^-, \theta_n^+]$

depend on data (after seeing our data)
R.V

Confidence Interval

$x \rightarrow$ Estimation $\rightarrow \hat{\theta}$

$[\hat{\theta}^-, \hat{\theta}^+]$

/ Producing an interval too.

$[x] = n \cdot \frac{1}{n} = [x]$
we want it to highly likely to have our output result.

$$P(\hat{\theta}_n^- \leq \theta \leq \hat{\theta}_n^+) \geq 1 - \alpha$$

$\therefore \alpha \approx 0.05$ (large error)

↓
prob 0.95 this will happen!

with a 95% prob = b/w 1.98 to 2.56

numbers [not random]

$\therefore \theta \rightarrow$ not random.

'where is the randomness'



Interval \rightarrow with 95% prob it has real value.

\therefore Probability: 95% confidence of being the real value

[] \rightarrow unlucky

95% of days - we will be lucky.

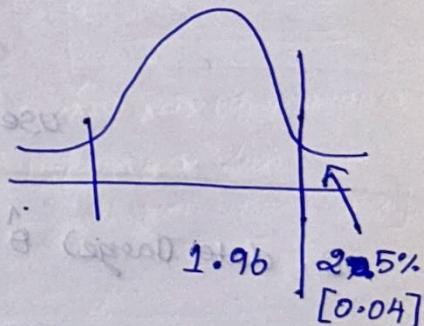
θ (real)

Confidence Interval

$$\hat{\theta}_n = (x_1 + \dots + x_n) / n \quad [\text{sample mean}]$$

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

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



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

Lower end

Upper end.

more generally: let z be S.O.F $\Phi(z) = 1 - \frac{\sigma}{2}$

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

* when n is large — more cons. abt — small interval [very little noise]

* vice versa.

$$\phi(z) = 1 - \frac{\sigma}{\sqrt{2}}$$

↓
Used Central Limit theorem (approximate) $\approx 95\%$. may not be accurate.

difficulty: I don't know σ .

case: unknown σ

- 1) Use upperbound on σ : pf X_i Bernoulli $\leq \frac{1}{2}$ [most biggest confidence level]

- 2) Estimate from data.

$$*\sigma = \sqrt{\theta(1-\theta)} \rightarrow \text{Bernoulli}$$

we don't know θ [But we have estimate]

$$= \sqrt{\hat{\theta}(1-\hat{\theta})}$$

use S.D. in the construction of confidence interval.

as data (large) $\hat{\theta} \approx \theta$ (good estimate)

$$\sigma = \sqrt{\hat{\theta}(1-\hat{\theta})} \text{ also good.}$$

3)

don't have any formulas like case ②

Genetic method: estimate variance \rightarrow mean θ $\left[(x_i - \theta)^2 \right]$

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \theta)^2 \rightarrow \hat{\theta}^2$$

(By weak law of large numbers)

obtain lots of measurement then take average

we don't know θ

[$\theta \rightarrow \text{mean}$]

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\theta}_n)^2 \rightarrow \sigma^2$$

(unbiased) $E[S_n^2] = \sigma^2$

AS large data.

$n-1$ [\therefore for being unbiased \rightarrow factors]

$$\therefore \frac{1}{100} \approx \frac{1}{99} \rightarrow \text{not a big deal.}$$

overall

1) Estimate Sample mean

2) S.D variance

3) use \rightarrow to find confidence interval

Confidence Interval \downarrow major \rightarrow use $\hat{\theta}$ to find

Pretend

S.D.

Sample mean

Normal distribution

(Central limit theorem)

Linear regression: → Today

* Likelihood \propto based on θ
 * Similar to Bayesian without prior.

Picking more likely θ ←
 pick θ evok syntactically but approx
 where x is observed more likely.

$\max_{\theta} P_X(x; \theta) \rightarrow$ choose more likely θ .

Sample mean estimate:

$$\hat{\theta}_n = \frac{x_1 + \dots + x_n}{n} \rightarrow \text{common mean [Take average]}$$

$1-\alpha$ confidential interval: $P(\hat{\theta}_n^- \leq \theta \leq \hat{\theta}_n^+) \geq 1-\alpha, \forall \theta$
 ↓
 Particularly

For Sample mean - CI

$$z \text{ be s.t } \phi(z) = 1 - \alpha/2$$

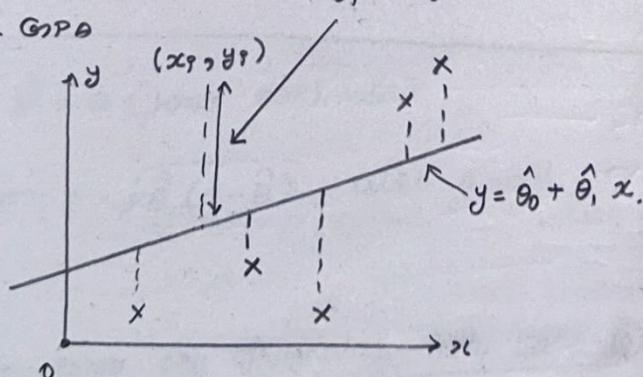
$$P(\hat{\theta}_n - \frac{z\sigma}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{z\sigma}{\sqrt{n}}) \approx 1 - \alpha$$

The interval has prob abt 95% eg having the real value in it.
 ∵ θ is not random [constant]

pretend: Sample mean \approx Normal distribution (good when n is large)

$x \rightarrow$ SAT \rightarrow high school

$y \rightarrow$ MIT - GPO



Data: $(x_1, y_1), \dots, (x_n, y_n)$

Model: $y \approx \theta_0 + \theta_1 x \rightarrow$ Linear relationship

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \rightarrow \text{Square error}$$

• we can't expect the line to pass through all points?

↓
we expect them to be close.
(randomness).

For a particular θ_0, θ_1 : calculate errors
choose θ_0, θ_1 with minimum errors.
↓
Square errors

one linear penetration:

$$* Y_g = \theta_0 + \theta_1 x_g + w_g, \quad w_g \sim N(0, \sigma^2), \text{ IID}$$

↓
GPA ↓
random noise.

SAT + randomness.

Random noise: Independent $\sim N(0, \sigma^2)$ 9.9-d

Likelihood:

likelihood: $f_{x, y | \theta}(x, y; \theta) : ?$

For a particular w : $C e^{-w^2 / 2\sigma^2}$

$$\text{for } Y_g = \theta_0 + \theta_1 x_g + w_g$$

Likelihood is

$$C \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{g=1}^n (y_g - \theta_0 - \theta_1 x_g)^2 \right\}$$

* Take logs, take derivatives

* Least squares \leftrightarrow penetrate w_g 9.9-d normal.

maximize this

$$\max \theta_0, \theta_1$$

After log:

$\max : (\text{since } -\text{ve sign} \rightarrow \text{smallest as possible})$

Linear regression

Conclusion: choose to do linear regression on maximum likelihood

$$\text{where } y_i = \theta_0 + \theta_1 x_i + w_i \quad (w_i \sim N(0, \sigma^2))$$

Assumes this

carrying out

$$\min_{\theta_1, \theta_0} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

then only C. exp $\left\{ -\frac{1}{2\sigma^2} \right\}$ will be maximum

differentiate $\rightarrow 0$

$\frac{d}{d\theta}$ (quadratic function) \rightarrow linear system

Set derivatives to zero:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}, \bar{y} = \frac{y_1 + \dots + y_n}{n}$$

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Solving quadratic

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

Interpretation

w - independent of x'

$$E[y] = \theta_0 + \theta_1 E[x] + 0 [w \text{ has zero mean}]$$

$$\theta_0 = E[y] - \theta_1 E[x]$$

$$\boxed{\theta_0 = \bar{y} - \theta_1 \bar{x}}$$

Estimate $\hat{\theta}_0$

$$\boxed{\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}}$$

Estimate $\hat{\theta}_1$

Assume x is simple. $E[x] = E[w] = 0$

$$y = \theta_0 + \theta_1 x + w$$

Mult by x

$$x y = \theta_0 x + \theta_1 x^2 + w x$$

$$\text{cov}(xy) = \theta_0 E[x] + \theta_1 E[x^2] + \text{cov}(wx)$$

$$\text{cov}(xy) = 0 + \theta_1 E[x^2] + 0$$

$$\text{cov}(xy) = \theta_1 E[(x - 0)^2]$$

$$\text{cov}(xy) = \theta_1 \text{var}(x)$$

$$\boxed{\theta_1 = \frac{\text{cov}(x,y)}{\text{var}(x)}}$$

when mean is not zero $E[x] \neq 0$

$$\theta_1 = \frac{\text{cov}(x,y)}{\text{var}(x)} = \frac{E[(x - E[x])(y - E[y])]}{E[(x - E[x])^2]}$$

Linear regression

model $\rightarrow y \approx \theta_0 + \theta_1 x$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

GPA may affected by several factors?

multiple linear regression

data: $(x_g, x_g', x_g'', y_g), g=1, \dots, n$

model: $y \approx \theta_0 + \theta_1 x + \theta_2 x' + \theta_3 x''$

formulation:

$$\min_{\theta_0, \theta_1, \theta_2} \sum_{g=1}^n (y_g - \theta_0 - \theta_1 x_g - \theta_2 x'_g - \theta_3 x''_g)^2$$

↓
Take derivation (set to zero)

↓
system of linear eqns.

choosing right variables: → (quadratic may fit) On some other,

model $y \approx \theta_0 + \theta_1 h(x)$

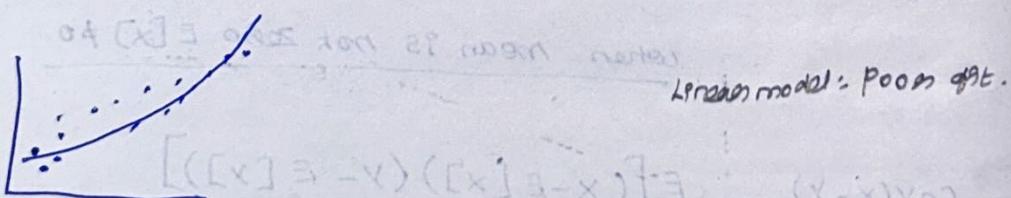
e.g.: $y \approx \theta_0 + \theta_1 x^2$

'work with data points $(y_1, h(x_1))$

formulation:

$$\min_{\theta_0, \theta_1} \sum_{g=1}^n (y_g - \theta_0 - \theta_1 h_1(x_g))^2$$

e.g.: GPA has something to do with SAP
(or) Success of SAP?



'A quadratic model may be good?'

choose: x^2 as planetary variable.

How to deal with $h(x)$

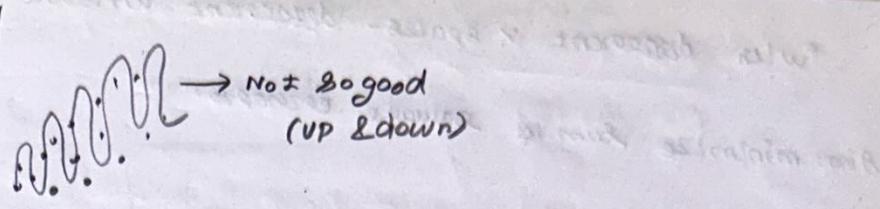
'multiple linear regression'

Linear

Any others (quadratic?)

may be 10m degree polynomial

(prediction errors extremely small)



* Avoid: too many parameters when little data is available'

Confidence Intervals

* Estimate σ [how much randomness?]

* $R^2 \rightarrow$ measure of explanatory power.

$$y \approx \theta_0 + \theta_1 x + w$$

Estimate σ^2 [std. errors]

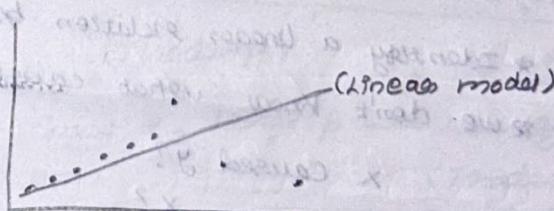
* How good we can predict

$\frac{Var(y|x)}{Var(y)}$ → Comparing: randomness of y knowing x
& randomness of y .

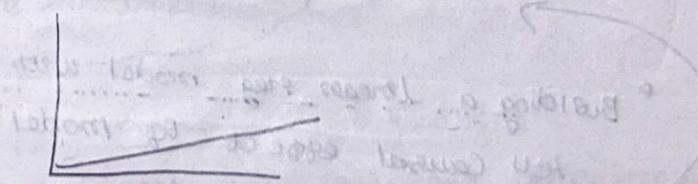
* How much randomness of y is attached with x .

Concerns

1) heteroskedasticity



'w has small variance'

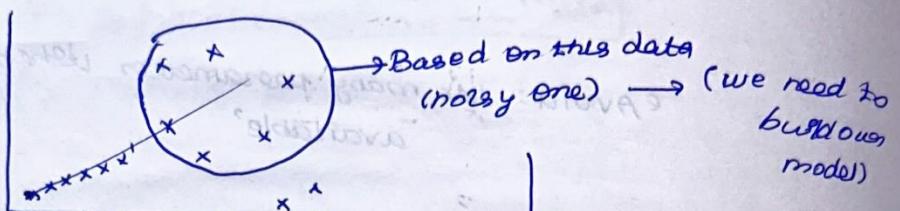


'w has large variances'

* w is not identically distributed - variance of w depends on x

With different x spaces - different amount of noise'

Aim: minimize sum of square errors.



Not the right thing.
(deal in many ways)

multiple explanatory variable

$$y_{GPA} = \theta_0 + \theta_1 SAT_{1st \text{ attempt}} + \theta_2 SAT_{2nd \text{ attempt}}$$

very likely to be close.

'Ignore one score = go with others'

'sensitive: data = θ_1 big, θ_2 small'

↓
our answer drastically changes

'Careful' - make special test.

Run linear regression, Run 1/Pearson model:

* Identify a linear relation b/w x & y

* we don't know what caused which?

x caused y?

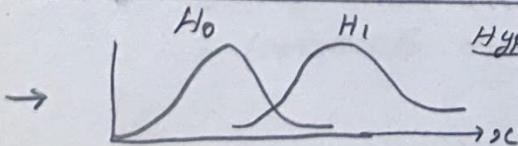
y " x?

x, y caused by any?

Building a linear reg model with small errors doesn't tell causal effect of model'

↓
Tells

close association b/w the variables.



Hypothesis testing

which distribution is correct?
 H_0 or H_1