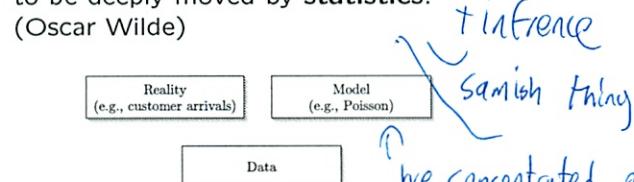


LECTURE 21

- Readings: Sections 8.1-8.2

"It is the mark of truly educated people to be deeply moved by **statistics**."
(Oscar Wilde)



- Design & interpretation of experiments

– polling, medical/pharmaceutical trials...

- Netflix competition

movies (try to predict rating)
of movie
have not seen



- Signal processing

– Tracking, detection, speaker identification,...

How use this table

- Sophisticated probabilistic model

Last unit: statistics

models only useful if

related to reality

delicate affair
manipulate data in sound manner
so data speaks for itself

Collect data to validate model
Extrapolate how stock's doing
"technical stock analyst"

Measure something, trying to figure out
what is happening
best fit to the data

Types of Inference models/approaches

- Model building versus inferring unknown variables. E.g., assume $X = aS + W$
 - Model building:
know "signal" S , observe X , infer a
 - Estimation in the presence of noise:
know a , observe X , estimate S .
- Hypothesis testing:** unknown takes one of few possible values; aim at small probability of incorrect decision - radar example
- Estimation:** aim at a small estimation error

Main methodologies

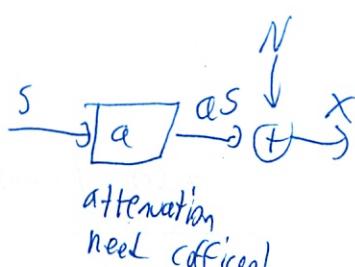
Type of interpretation

1. Have data, build model

2. Have model, get observations

of outcomes X

Mathematically the two problems are the same



2 fundamental methodologies

- Classical statistics:** Some # - guess of unknown value

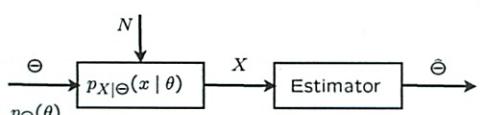
Mass of electron
 N noise
 $p_X(x; \theta)$ prob dist.
 X come up w/ estimate

prob dist. Not conditional prob

θ is a constant - but we don't know it

closer to what we were doing in class

- Bayesian:** Use priors & Bayes rule



X know theta is about here ~ RV
past accumulated beliefs of theta

Collect data
improve knowledge of theta

) will focus on this
for next few lectures

Bayesian inference: Use Bayes rule

- Hypothesis testing

- discrete data

$$p_{\Theta|X}(\theta | x) = \frac{p_{\Theta}(\theta) p_{X|\Theta}(x | \theta)}{p_X(x)}$$

4 Versions
depends if discrete/continuous
or
continuous data

posterior
of theta-
given everything

$$p_{\Theta|X}(\theta | x) = \frac{p_{\Theta}(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

- Estimation; continuous data

(continues)
discrete hypothesis/unknown quantity

$$f_{\Theta|X}(\theta | x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

$Z_t = \Theta_0 + t\Theta_1 + t^2\Theta_2$ find thetas that best match data

$$X_t = Z_t + W_t, \quad t = 1, 2, \dots, n$$

Bayes rule gives:

$$f_{\Theta_0, \Theta_1, \Theta_2 | X_1, \dots, X_n}(\theta_0, \theta_1, \theta_2 | x_1, \dots, x_n)$$

theta vector
relative likelihood
of diff thetas

Prior $\rightarrow P_{\Theta_0, \Theta_1, \Theta_2}(\theta_0, \theta_1, \theta_2)$

Then model of measurement

- dist of X_s
- which is dist of Z - given by formula

then get posterior dist

Estimation with discrete data

$$f_{\Theta|X}(\theta | x) = \frac{f_{\Theta}(\theta) p_{X|\Theta}(x | \theta)}{p_X(x)}$$

$$p_X(x) = \int f_{\Theta}(\theta) p_{X|\Theta}(x | \theta) d\theta$$

- Example:

- Coin with unknown parameter θ
- Observe X heads in n tosses

- What is the Bayesian approach?

- Want to find $f_{\Theta|X}(\theta | x)$
- Assume a prior on Θ (e.g., uniform)

Bayesian will use Bayes rule

- prior belief -

- if no data \rightarrow use uniform

$f_{X|\theta}$ = model of measurement \sim binomial here

find prob dist of θ , given X_s observed

Don't know bias of coin

f = prob of heads

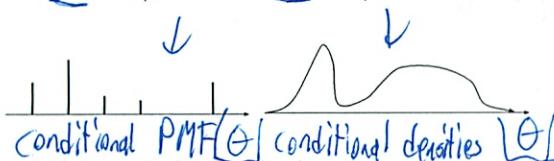
$\hat{\theta}$ = estimate

basic estimate $\hat{\theta} = \frac{X}{n}$ & classical stat

Output of Bayesian Inference

- Posterior distribution:

- pmf $p_{\Theta|X}(\cdot | x)$ or pdf $f_{\Theta|X}(\cdot | x)$



- If interested in a single answer:

- Maximum a posterior probability (MAP):

$$p_{\Theta|X}(\theta^* | x) = \max_{\theta} p_{\Theta|X}(\theta | x)$$

minimizes probability of error;
often used in hypothesis testing

$$f_{\Theta|X}(\theta^* | x) = \max_{\theta} f_{\Theta|X}(\theta | x)$$

- Conditional expectation:

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

- Single answers can be misleading!

want prob of making mistake is as small as possible

highest bar \Rightarrow call it θ^*

The sum of the height of other bars

for continuous - how receivers in communication systems work

maximum density point
mode



or look at median
or look at expected value / mean

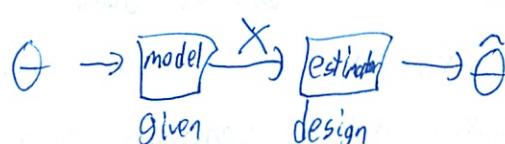
3 possible choices!

all reasonable

Plausible estimates

- need to compare different ways of producing estimates

most popular = mean squares



Function of data to produce an estimate $\hat{\theta} = g(x)$

$$\theta - \hat{\theta} = \text{error}$$

look at square, b/c don't like errors

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

Everything here is random!

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

Finding

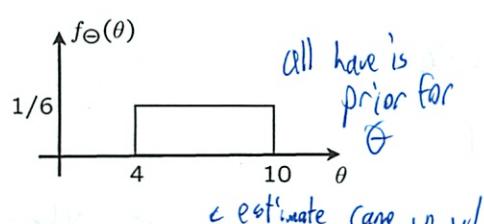
$$E[(\theta - c)^2] = E[\theta^2] - 2E[\theta]c + c^2$$

$$\downarrow \frac{\partial}{\partial c} = 0 : -E[\theta] + c = 0$$

$$E[\theta] = c$$

Least Mean Squares Estimation

- Estimation in the absence of information



- find estimate c , to:

$$\text{minimize } E[(\theta - c)^2]$$

- Optimal estimate: $c = E[\Theta]$

- Optimal mean squared error:

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

w/o any data!

- (this is what we did in 6.01)

LMS Estimation of Θ based on X

- Two r.v.'s Θ, X
- we observe that $X = x$
 - new universe: condition on $X = x$
- $E[(\Theta - c)^2 | X = x]$ is minimized by

$$c = E[\Theta | X=x] \text{ because we had an observation}$$

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

$$\leq E[(\Theta - g(x))^2 | X = x] \quad \begin{matrix} \text{echo 4} \\ \downarrow \text{when } X=x \\ \text{thinking of conditionals as RVs} \end{matrix}$$

- $E[(\Theta - E[\Theta | X])^2 | X] \leq E[(\Theta - g(X))^2 | X]$

$$E[\quad] \quad E[\quad] \quad \begin{matrix} \downarrow \\ \text{can do iterated expectations} \end{matrix}$$

- $E[(\Theta - E[\Theta | X])^2] \leq E[(\Theta - g(X))^2]$
 $\downarrow g - g \text{ is estimator}$
is optimal

$E[\Theta | X]$ minimizes $E[(\Theta - g(X))^2]$
over all estimators $g(\cdot)$

Smallest possible over all estimators

lots of data

- can still make same argument

Sometimes hard to calculate

- lots of joint prob

- sometimes may have to settle for less than
this

- devil is in the details

before just prior dist $P_\theta(\theta)$

after observing $P_{\theta|x}(\theta | x)$

estimate = # $\hat{\theta} \in E[\theta | x]$
estimator = formula or algorithm $\hat{\theta}$
conditional expectation $E[\theta | X]$

LMS Estimation w. several measurements

- Unknown r.v. Θ
- Observe values of r.v.'s X_1, \dots, X_n
- Best estimator: $E[\Theta | X_1, \dots, X_n]$
- Can be hard to compute/implement
 - involves multi-dimensional integrals, etc.

LECTURE 22

- Readings: pp. 225-226; Sections 8.3-8.4

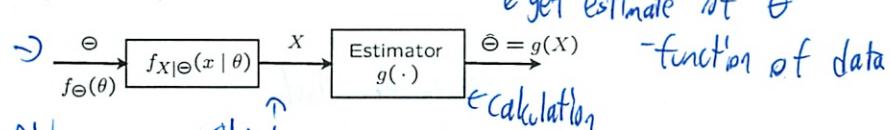
Topics

2nd lecture

- (Bayesian) Least means squares (LMS) estimation

- (Bayesian) Linear LMS estimation

Unknown quantity
treat as RV



- MAP estimate: $\hat{\theta}_{MAP}$ maximizes $f_{\theta|X}(\theta | x)$

- LMS estimation:

- $\hat{\theta} = E[\theta | X]$ minimizes $E[(\theta - g(X))^2]$
over all estimators $g(\cdot)$

- for any x , $\hat{\theta} = E[\theta | X = x]$
minimizes $E[(\theta - \hat{\theta})^2 | X = x]$
over all estimates $\hat{\theta}$

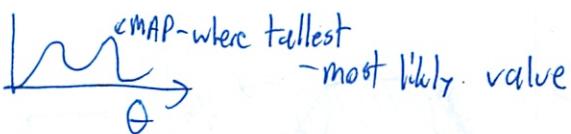
treat as #
if give specific
value of data

how far
estimate
is from
squared value
simple solution
 $E[\theta | X]$

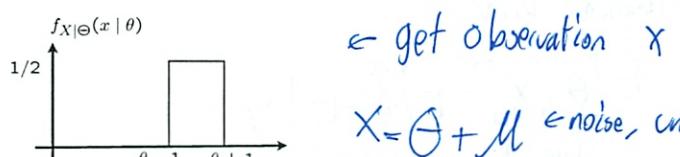
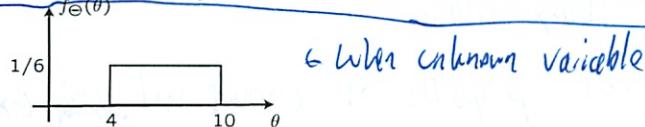
Sometimes good estimator

- What is good way of getting estimate

Obtain Posterior Dist



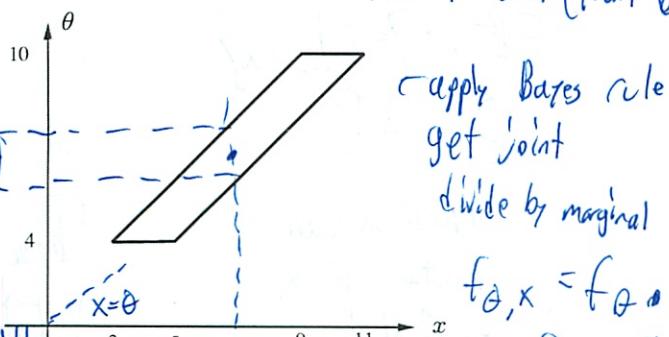
= recalculate in new universe → find mean + is best estimate



$X = \theta + \eta$ ← noise, uniform [-1, 1]

Observation error (from 6.01 I think)

~~for R2R2R2R2R2~~
~~= S2S2S2S2S2~~
~~for R2R2R2R2R2~~
~~= S2S2S2S2S2~~



$$f_{\theta|x} = \frac{f_{\theta,x}}{f_x}$$

$$f_{\theta,x} = f_{\theta} \cdot f_{x|\theta} = \frac{1}{12}$$

for possible values

Can see possible
values of θ
Uniform dist
Mean/midpoint
is in the middle
ends up being

?
When told
an x

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

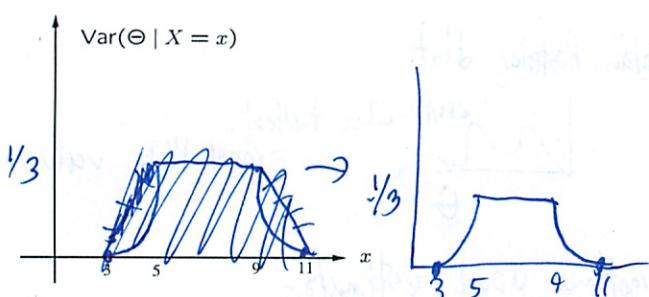
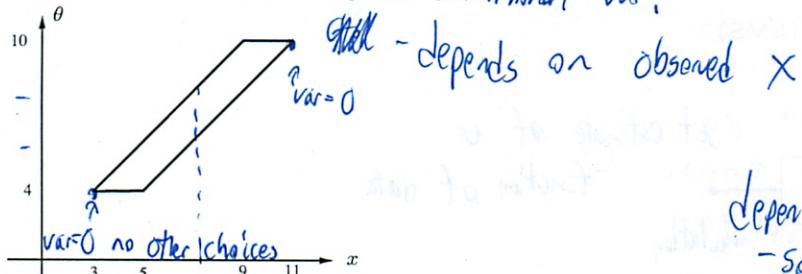
(nice - I like this!)

How good is your estimate?

Conditional mean squared error

- $E[(\Theta - E[\Theta | X])^2 | X = x]$ trying to minimize

- same as $\text{Var}(\Theta | X = x)$: variance of the conditional distribution of Θ just conditional var!



depends on data

- some data very useful
- since small var

(Var = estimation error)

Some properties of LMS estimation

"Heavy lifting"

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

- Estimation error: $\tilde{\Theta} = \hat{\Theta} - \Theta$

What is the avg value of estimation error

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

- $E[\tilde{\Theta}h(X)] = 0$, for any function h

- $\text{cov}(\tilde{\Theta}, \hat{\Theta}) = 0$

- Since $\Theta = \hat{\Theta} - \tilde{\Theta}$:

$$\text{var}(\Theta) = \text{var}(\hat{\Theta}) + \text{var}(\tilde{\Theta})$$

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

linearity of expectation

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

$$= \hat{\Theta} - \Theta$$

depends on value of X
function of X
so $\hat{\Theta}$ itself

$$= 0$$

$E[\tilde{\Theta} | X=x] = 0$ still 0 when have specific x

law of iterated expectations

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

basic principle
 $E[g(X)|X=x] = g(x)$

Linear LMS

Simplify a bit

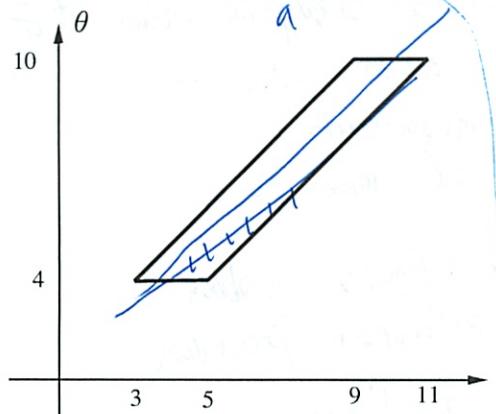
- Consider estimators of Θ , of the form $\hat{\Theta} = aX + b$

- if too difficult to calc. conditional expectation

- Minimize $E[(\Theta - aX - b)^2] = h(a, b)$ instead of getting optimal estimator

- Best choice of a, b ; best linear estimator:

$$\hat{\Theta}_L = E[\Theta] + \frac{\text{Cov}(X, \Theta)}{\text{var}(X)}(X - E[X])$$



- just get a single line
- easier to calculate estimate
- losing optimality

Find constants a, b , so mean squared line is least

Expand as function of a, b

$$\frac{\partial h}{\partial a} = 0 = \frac{\partial h}{\partial b}$$

- will be quadratic in a, b
- will get linear eq
- solve

Optimal form is this

Linear LMS properties

Single RV

if no data, a disappears or if they are ind.
if have data, add correction information new info

$$\hat{\Theta}_L = E[\Theta] + \frac{\text{Cov}(X, \Theta)}{\text{var}(X)}(X - E[X])$$

$$E[(\hat{\Theta}_L - \Theta)^2] = (1 - \rho^2)\sigma_\Theta^2$$

mean squared estimation error

Linear LMS with multiple data

multiple observations

- Consider estimators of the form:

$$\hat{\Theta} = a_1X_1 + \dots + a_nX_n + b$$

linear to keep it simple

- Find best choices of a_1, \dots, a_n, b

- Minimize:

$$E[(a_1X_1 + \dots + a_nX_n + b - \Theta)^2]$$

- Set derivatives to zero

linear system in b and the a_i

- Only means, variances, covariances matter

of Θ reduced by a certain factor
- if ρ close to 1, error becomes small
- Exploiting correlations w/ b/w data + unknown quant
- Can reduce error if data somewhat correlated
choice of coefficients to get smallest possible error

lots of calculations
do w/ computer

R - also cov of cross terms

- detailed stats don't matter
(plots)

- only things that matter is these

Can be tedious

The cleanest linear LMS example

$$X_i = \Theta + W_i, \quad \Theta, W_1, \dots, W_n \text{ independent}$$
$$\Theta \sim \mu, \sigma_0^2 \quad W_i \sim 0, \sigma_i^2$$

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

(weighted average of μ, X_1, \dots, X_n)

- If all normal, $\hat{\Theta}_L = E[\Theta | X_1, \dots, X_n]$

less tedious, get clean formula

Observe several noisy measurements

Usually avg measurements

- which is roughly optimal

- but want to include prior mean of Θ

- but not all measurements are =

- must weight measurements

- weights inv of noise vars

Choosing X_i in linear LMS

- $E[\Theta | X]$ is the same as $E[\Theta | X^3]$

- Linear LMS is different:

- $\hat{\Theta} = aX + b$ versus $\hat{\Theta} = aX^3 + b$

- Also consider $\hat{\Theta} = a_1X + a_2X^2 + a_3X^3 + b$

Same info, so same estimation problem

- won't change estimation procedure

- unless looking at linear estimators

- need to look at different estimators

- do what you think will fit the data better

Some orientation

Review

- Standard examples:

- X_i uniform on $[0, \theta]$; uniform prior on θ
- X_i Bernoulli(p); uniform (or Beta) prior on p
- X_i normal with mean θ , known variance σ^2 ; normal prior on θ ;
 $X_i = \Theta + W_i$

$$\begin{bmatrix} X \\ X^2 \\ X^3 \end{bmatrix} \text{ different form again}$$

- art in what you call the data in problem

- LMS does not matter

- matters in linear

- Estimation methods:

- MAP
- MSE
- Linear MSE

$$\begin{aligned}
 E[\tilde{\theta} h(x) | x] &= E[\hat{\theta} h(x) | x] - E[\theta h(x) | x] \\
 &= \hat{\theta} h(x) - h(x) E[\theta | x] \\
 &= \hat{\theta} h(x) - h(x) \hat{\theta} \\
 &= 0
 \end{aligned}$$

↑
 "if x known,
 - just the values
 ↗
 mae outside
 since quantity known
 - nothing random about it

$$E[\tilde{\theta}^2 h(x)]$$

↑ uncorrelated w/ the error

No function of data ~~gives~~ that gives us useful info about the error

(could use to ↓ error further
 but if optimal, can not ↓ error further)

Use this to look at cov

$$\begin{aligned}
 \text{cov}(\tilde{\theta}, h(x)) &= E[\tilde{\theta} h(x)] - E[\tilde{\theta}] E[h(x)] \\
 &= 0 - 0 E[h(x)] \\
 &= 0
 \end{aligned}$$

(2)

Apply to ~~special~~ special case where $h(x) = \hat{\theta}$

$\hat{\theta}$ is just function of x

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

Error and estimate are uncorrelated

~~large~~
informative

- Since uncorrelated, add vars

AV

- So having big var in estimate is good

$\hat{\theta}$ big $\rightarrow \hat{\theta}$ affected by the data
↳ When the data is informative

helps make $\hat{\theta}$ small \rightarrow so small error

- Can prove w/ law of total variance

$$= \text{Var}(\text{estimate}) + \text{Var}(\text{error})$$

5.4 Central Limit Theorem

- according to we L #, distribution of sample mean ^{increasingly concentrated} in the near vicinity of true mean μ
 - Var tends to 0
 - so $S_n = X_1 + \dots + X_n = nM_n \rightarrow \infty$
and distribution does not converge to anything meaningful
 - do this scaling $Z_n = \frac{S_n - n\mu}{\sigma\sqrt{n}} = \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$
- $$E[Z_n] = \frac{E[X_1 + \dots + X_n] - n\mu}{\sigma\sqrt{n}} = 0$$
- $$\text{Var}(Z_n) = \frac{\text{Var}(X_1 + \dots + X_n)}{\sigma^2 n} = \frac{\text{Var}(X_1) + \dots + \text{Var}(X_n)}{\sigma^2 n} = \frac{n\sigma^2}{n\sigma^2} = 1$$

(Central Limit Theorem)

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \quad \text{mean} \quad \text{Var} = \sigma^2$$

CDF F^{2n} converges to st normal

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx \rightarrow \lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z) \quad \text{for every } z$$

(2)

Very general - only requires iid, finite mean, var

Shows sum of large # of iid. RVs is \sim normal

applies to many situations (been proven empirically)

Also makes hard calculations simpler

- ↳ allows us to calculate probabilities related to Z_n
- ↳ like treating S_n as normal
 - w/ mean $n\mu$
 - var $n\sigma^2$

Normal Approx based on Central Limit Theorem

$S_n = X_1 + \dots + X_n \leftarrow X_i$ iid RVs mean μ var σ^2

If n is large $P(S_n \leq c) \sim$ normal by:

1. Calculate mean $n\mu$ var $n\sigma^2$ of S_n

2. Calculate the normalized value $Z = \frac{(c - n\mu)}{\sigma\sqrt{n}}$

3. Use approx $P(S_n \leq c) \approx \Phi(z)$

↳ from st. normal tables

(3)

example

have 100 packages that weigh b/w 5 → 50 lbs

$$P(\text{weight} < 3000 \text{ lbs})$$

caiser to do w/ central limit

$$S_{100} = X_1 + \dots + X_{100}$$

$$\hookrightarrow \mu = \frac{5+50}{2} = 27.5$$

↑ where get this? oh dist of ~~each~~ uniform RV

$$\sigma^2 = \frac{(50 - 5)^2}{12}$$

Then normalize:

$$z = \frac{3000 - 100 \cdot 27.5}{\sqrt{168.75 \cdot 100}} = \frac{250}{129.9} = 1.92$$

Plug into st. normal table

$$P(S_{100} \leq 3000) \approx \varphi(1.92) = 0.9726$$

$$\begin{aligned} \text{So } P(S_{100} > 3000) &= 1 - P(S_{100} \leq 3000) \\ &= 0.0274 \end{aligned}$$

(4)

If $\text{Var } X_i$ is unknown can use upper bound

- example S&L polling

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

$$M_n \rightarrow \text{Mean } p \quad \text{Var } \frac{p(1-p)}{n}$$

(are they sure)
not X_i) So M_n is approx normal

Interested in $P(|M_n - p| \geq \epsilon)$

- prob that polling error is larger than
desired accuracy ϵ

$$P(|M_n - p| \geq \epsilon) \approx 2 P(M_n - p \geq \epsilon)$$

M_n

Symmetry of normal PDF around mean

$\text{Var } \frac{p(1-p)}{n}$ of $M_n - p$ depends on p so unknown

so assume largest possible $\text{Var } \frac{1}{4n}$ (rest of this)
which corresponds to $p = \frac{1}{2}$

(5)

Now use this ~~standardized~~ ~~var~~ to find standardized value

$$Z = \frac{E}{\sqrt{2n}}$$

$$\begin{aligned} P(M_n - p \geq \epsilon) &\leq 1 - \phi(z) \\ &= 1 - \phi(2\epsilon\sqrt{n}) \end{aligned}$$

So if $n=100$ $\epsilon = .1$

$$\begin{aligned} P(M_{100} - p \geq .1) &\approx 2P(M_n - p \geq .1) \\ &\leq 2 - 2\phi(2 \cdot .1 \cdot \sqrt{100}) \\ &= 2 - 2\phi(2) \\ &= 2 - 2 \cdot .0977 \\ &= .046 \end{aligned}$$

Much smaller / more accurate than Chebshov

Normal approx is increasingly accurate as $n \rightarrow \infty$

But usually we have a finite n

No simple ans when n is big enough

(6)

De Moivre-Laplace Approx to Binomial

- (did we do this in lecture?)

- Binomial RV S_n w/ param n, p

- Sum of n ind. Bernoulli RV X_1, \dots, X_n w/ common param p

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

$$\mu = E[X_i] = p \quad \sigma^2 = \text{Var}(X_i) = \sqrt{p(1-p)}$$

Now want prob $k \leq S_n \leq l$

$$\frac{k-np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{l-np}{\sqrt{np(1-p)}}$$

↑ approx st. normal

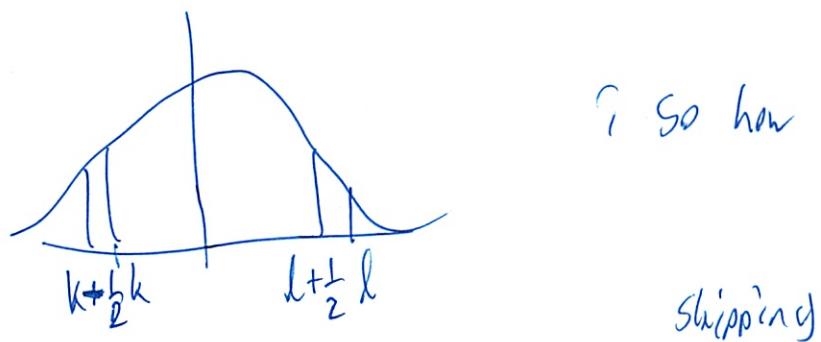
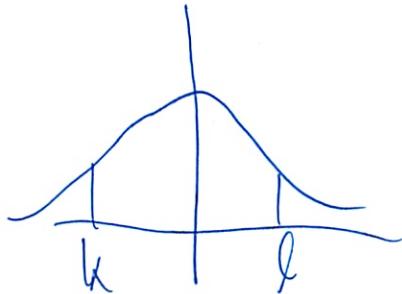
$$P(k \leq S_n \leq l) = P\left(\frac{k-np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{l-np}{\sqrt{np(1-p)}}\right)$$

$$\approx \Phi\left(\frac{l-np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k-np}{\sqrt{np(1-p)}}\right)$$

$$\approx \Phi\left(\frac{l+\frac{1}{2}-np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k-\frac{1}{2}-np}{\sqrt{np(1-p)}}\right)$$

6)

- got more accurate by replacing $k \rightarrow k - \frac{1}{2}$
 $l \rightarrow l + \frac{1}{2}$



? So how does that help?

skipping

5.5 Strong Law of Large

- similar to weak law ~~weak~~
- also about convergence of sample mean to true mean
- but it's a different type of convergence

X_1, X_2, \dots are iid RVs w/ mean M

$M_n = \underbrace{(X_1 + \dots + X_n)}_n$ converges to M ~~weakly~~

w/ prob 1

$$P\left(\lim_{n \rightarrow \infty} \frac{X_1 + \dots + X_n}{n} = M\right) = 1$$

⑧

States that all of prob concentrated on a particular subset of sample set

I kinda confused on what it says

Convergence w/ Prob 1

Let Y_1, Y_2, \dots be seq of RVs (not necessarily ind)

C = real #

Say Y_n converges to C w/ prob 1 (surely) if
 $P(\lim_{n \rightarrow \infty} Y_n = C) = 1$

All prob concentrated on seq that converge to C
- other seq unlikely

Convergence w/ prob 1 implies conv. in prob
~~✓~~

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
(Fall 2010)

Recitation 22
November 30, 2010

Examples 8.2, 8.7, 8.12, and 8.15 in the textbook

Romeo and Juliet start dating, but Juliet will be late on any date by a random amount X , uniformly distributed over the interval $[0, \theta]$. The parameter θ is unknown and is modeled as the value of a random variable Θ , uniformly distributed between zero and one hour.

- (a) Assuming that Juliet was late by an amount x on their first date, how should Romeo use this information to update the distribution of Θ ?
- (b) How should Romeo update the distribution of Θ if he observes that Juliet is late by x_1, \dots, x_n on the first n dates? Assume that Juliet is late by a random amount X_1, \dots, X_n on the first n dates where, given θ , X_1, \dots, X_n are uniformly distributed between zero and θ and are conditionally independent.
- (c) Find the MAP estimate of Θ based on the observation $X = x$.
- (d) Find the LMS estimate of Θ based on the observation $X = x$.
- (e) Calculate the conditional mean squared error for the MAP and the LMS estimates. Compare your results.
- (f) Derive the linear LMS estimator of Θ based on X .
- (g) Calculate the conditional mean squared error for the linear LMS estimate. Compare your answer to the results of part (e).

Recitation 22

- graded p-set 8+9 on shelves outside TA office

- most important Bayes problem in recitation 2

- the disease one

- (also in 6.01)

- archetypal example

- update beliefs when have info

- just more complex in this chap

- parameter θ

- random

- known prob density = prior

- $f_{x|\theta}(x|\theta)$ = math description of how obs generated

- likelihood function

- obs model

- generative model

) aka

- $f_{\theta|x}(\theta|x)$ Θ = capital θ

- want to get via Bayes rule

- posterior

$$f_{\theta|x}(\theta|x) = \frac{f_{\theta}(\theta) f_{x|\theta}(x|\theta)}{\int f_{\theta}(\theta) f_{x|\theta}(x|\theta) d\theta}$$

②

Is everything we know

- sometimes too much

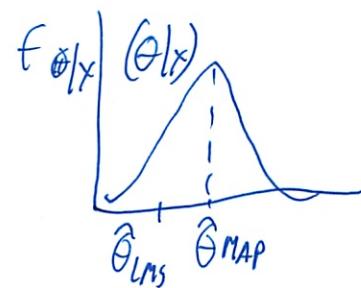
Want to boil down to 1 #

- maximum point on posterior = MAP

- Mean of posterior = LMS

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

$$- \text{Minimizes } E[(\hat{\theta} - \theta)^2 | X=x]$$



Romeo

1. Romeo and Juliet

- don't know θ

- as go on dates, she know late she is, increase knowledge of θ

$$f_{\theta}(\theta) = \begin{cases} 1 & 0 \leq \theta \leq p \\ 0 & \text{otherwise} \end{cases}$$

initial dist/prior

Lobs model

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

a) $f_{\theta|x}(\theta|x)$ ← saw one little x since went on 1 date

$$= \begin{cases} \frac{1}{\int_x^1 \frac{1}{\theta} d\theta} & \text{normalize } 0 \leq x \leq \theta \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

evaluate

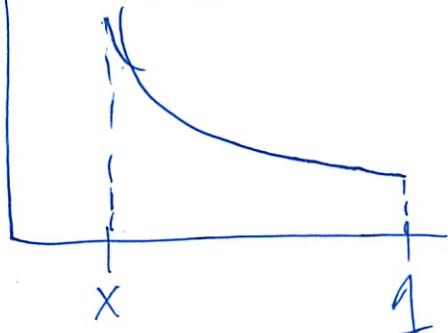
$$= \begin{cases} \frac{1}{\theta \log x} & 0 \leq x \leq \theta \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} \frac{1}{\theta \log(\frac{1}{x})} & \dots \\ 0 & \dots \end{cases}$$

another way to write

③

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



height

 $\frac{1}{\theta}$ shape

Pmt will late on 1st date
- after just normalization
So S to 1

b) Now go on n dates

$$f_{\theta|x_1, x_2, \dots, x_n}(\theta | x_1, x_2, \dots, x_n)$$

$$= \frac{f_{\theta}(\theta) f_{x_1, x_2, \dots, x_n | \theta}(x_1, x_2, \dots, x_n | \theta)}{\int \dots d\theta}$$

Write as product of marginal conditional dist.

$$= \frac{f_{\theta}(\theta) f_{x_1 | \theta}(x_1 | \theta) \dots f_{x_n | \theta}(x_n | \theta)}{\int \dots d\theta}$$

$$= \begin{cases} \frac{\dots}{\int \dots d\theta} & \dots \\ 0 & \text{otherwise} \end{cases}$$

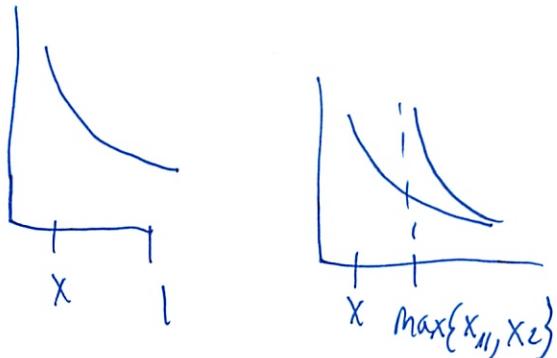
When is the first factor
between x_1 and 1 ?

confused

(9)

$$= \begin{cases} \frac{\frac{1}{\theta^n}}{\int_{\max\{x_1, x_2\}}^1 \frac{1}{\theta^n} d\theta} & \max\{x_1, x_2\} \leq \theta \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

For larger n \hookrightarrow becomes less flat/steeper



c) want a single # estimate

- here from 1 observation

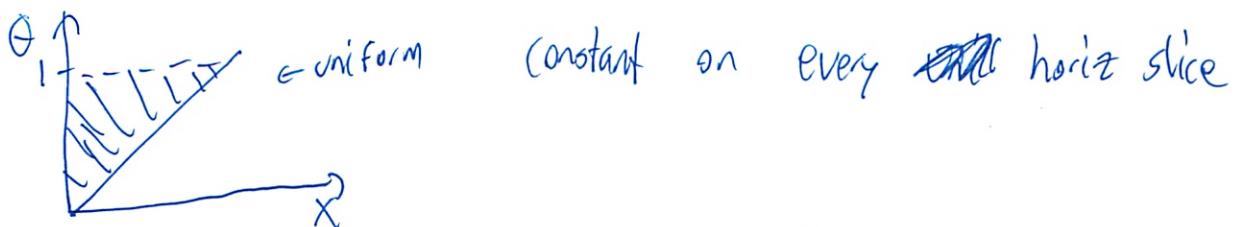
$\hat{\theta}_{MAP} = \underset{\text{value of } \theta}{\text{maximum of posterior}}$

notice that is decreasing function of θ

so is x

$$\hat{\theta}_{MAP} = x$$

extra) sketch joint dist

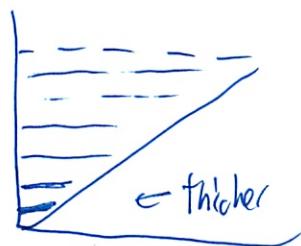


(5)

So need taller slices for tall values of θ
— for slices to ~~go~~ $\rightarrow 1$

Observe X taking slice in $|$ direction

— taller for small values of θ



MAP just takes left end of distribution

But picking | end might not be best estimate

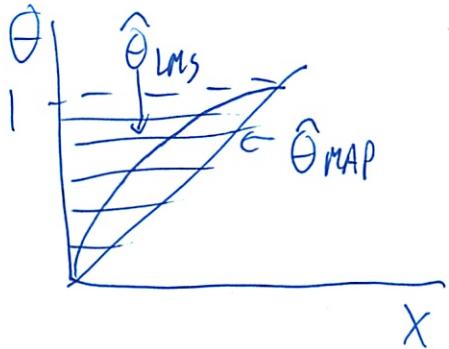
Minimize mean of square of error = LMS

Will be to left of midpoint

$$\begin{aligned} d) \hat{\theta}_{\text{LMS}} &= E[\hat{\theta} | X=x] \\ &= \int \theta f_{\hat{\theta}|X}(\theta|x) d\theta \\ &\quad \text{Posterior density} \\ &= \int_x^1 \theta \frac{1}{\theta \log(\frac{x}{\theta})} d\theta \\ &= \frac{1-x}{\log(x)} \end{aligned}$$

⑥

Then on our plot it would be
- at each point



e) Most tedious

Find conditional mean squared error for MAP, LMS

$$\mathbb{E}[(\theta - \hat{\theta})^2 | X=x] =$$

* Can make computations that apply in both cases

$$= \int (\theta - \hat{\theta})^2 f_{\theta|X}(\theta|x) d\theta$$

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

tedious, not hard calculation

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

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

for any deterministic function $\hat{\theta}$ - can put in LMS, MAP

7)

So if put in $\hat{\theta}_{MAP}$ \rightarrow Conditional Mean Squared Error (MSE)

$$x^2 + \frac{3x^2 - 4x + 1}{2 \log(\frac{1}{x})}$$

Comparison w/ $\hat{\theta}_{LMS}$

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

No insights to draw

Compare the two

LMS must be smaller than MAP by definition

for all possible value of X

Plots in textbook

f) Want to find the LLMS estimator

- minimizes LMS among all estimators that are functions of x

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

- defined in lecture
- "plausibly good"
- will not reprove

⑧

How to compute the γ constants?

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

$E[X] =$ have avoided calculating dist of X

can avoid calculating it

(conditioning on θ makes X simple)

so do iterated expectations

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

Conditioning on RV subtle constant

$$= E[E[X | \theta \neq \theta]]$$

Calculate for little θ

~~$E[\theta]$~~ $E[\theta] = \frac{\theta}{2}$

$$= E[\theta/2] \text{ then back to capital letters}$$

$$= \frac{1}{4}$$

~~Var~~ $\text{var}(X) =$ use law of total var

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

$$= E\left[\frac{1}{12}\theta^2\right] + \text{var}\left(\frac{\theta}{2}\right)$$

$$\text{var}[E[J]^2] = \frac{1}{12} \left(\frac{1}{12} + \left(\frac{1}{2}\right)^2 \right) + \frac{1}{4} \cdot \frac{1}{12}$$

trying to do w/ properties
of uniform dist

⑨

$$= \frac{1}{12} \cdot \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{1}{12} = \frac{1}{144}$$

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

$$= ? - \frac{1}{4} \cdot \frac{1}{2}$$

?

$$\overbrace{E[X\theta]} = E[E[X\theta|\theta]]$$

? condition on θ - easier

- find for little θ , then replace w/ θ

$$\overbrace{E[X\theta | \theta = \theta]}$$

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

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

$$= \theta \cdot \frac{\theta}{2}$$

$$= E\left[\frac{1}{2}\theta^2\right]$$

$$= \cancel{\frac{1}{2}} \cdot \cancel{\frac{1}{3}} \cdot \frac{1}{4} \text{ not did}$$

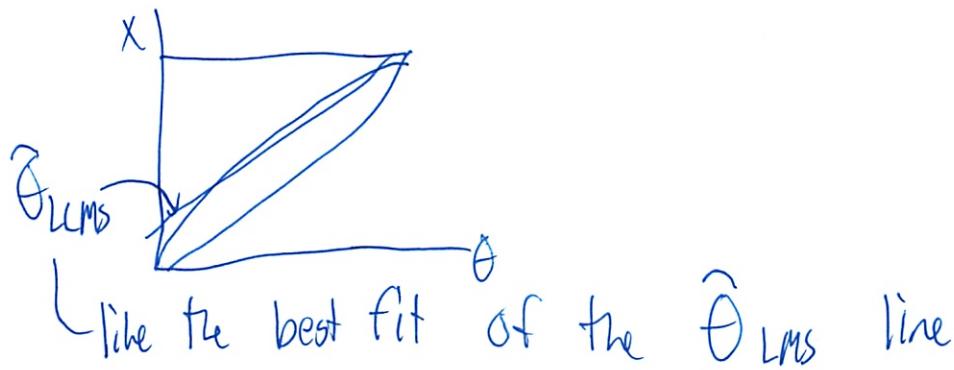
$$= \frac{1}{6}$$

$$\text{Cov}(X, \theta) = \frac{1}{24}$$

(10)

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

$$= \frac{1}{2} + \frac{6}{7} (x - \bar{x})$$



(Very cute how can draw this diagram - I want to become good at this)

Bayesian Statistical Inference Chap 8 Reading

11/2018
36

- process of extracting info about an unknown variable
of an unknown model from available data
- 2 main approaches; Bayesian + Classical
- Main categories of inference problems; parameter estimation, hypothesis testing, significance testing
- important methodologies; MAP, LMS, ~~max likelihood~~, regression, likelihood ratio tests

Probability vs Stats

- probability is a self-contained math subject
 - assume a fully specified prob model
 - use math to quantify consequences of model & ans qu
 - every qv has an answer
 - model taken for granted

Statistics more of an art form

- may be several reasonable methods \rightarrow several answers
- no 1 way to select right answer
- ~~also~~ require certain desirable principals
 - performance
 - past experience
 - common sense
 - community

(2) Bayesian vs Classical Stats

Bayesian - unknowns treated as RVs

(classical/frequentist) - deterministic quantity, ~~not~~ only unknown

Bayesian ~~can~~ tries to move stats back to probability

* dealing with multiple possible candidate probabilistic models *

(which are these again? EMAP, etc?)

Each one thinks their methods is better
(interesting discussion and comparison)

Model vs Variable Inference

Model inference - real phenomenon or process

want to construct/validate a model
have available data

to predict about the future or infer some underlying cause

Variable inference - estimate value of 1 or more unknown values
w/ noisy info

- Some overlap b/w the models
- May do both

(3)

A Rough Classification of Stat. Inference Problems

Estimation - model fully specified

- except for unknown, possibly multidimensional parameter θ
- want to estimate close to its exact value
- - RV (Bayesian)
- - unknown constant (classical)
- noisy transmission \rightarrow want to estimate a (model of channel)
- using polling data, estimate fraction of people that prefer one thing to the other
- estimate mean, var of movement of a stock

Binary hypothesis testing - start w/ 2 hyp. decide which is True

- use knowledge of a, X_i to decide if $s_i = 0$ or 1

$\begin{matrix} \uparrow & \uparrow \\ a & X_i \end{matrix}$
 noise received
 model bit

\uparrow
Sent bit

- with a noisy pic \rightarrow is there a person in there or not?
- decide which of 2 medical treatments is best

Mary hypothesis testing

- finite # m of hypotheses

Performance of a model judged by prob that it makes an erroneous decision

(4)

This chap focuses on estimation

Ignoring more advanced topics like non parametric $Y = g(X) + W$

\uparrow
 Unknown Function
 \uparrow
 Noise
 \Rightarrow Can't be described
 by Fixed # parameters

Principal Bayesian Inference Methods

Maximum a posteriori prob (MAP)

- the max on the posterior dist

Least Mean Squares (LMS) estimation

Select function/estimator that minimizes mean squared error b/w parameter and its estimate

Linear least Mean Squares (LLMs)

- linear function of data, like above
- simpler, but a little less accurate

(5)

8.1 Bayesian Inference + Posterior Distribution

Θ = capital Θ = unknown quantity of interest

↳ RV

↳ or finite # of RVs

↳ can represent physical quantities

$X = (X_1, X_2, \dots, X_n) = \text{Observations}$

Observation measurements
||
vector

know

joint dist Θ, X

OC

prior dist P_Θ / f_Θ

Conditional dist

$P_{X|\Theta} / f_{X|\Theta}$

Once X has been observed, get posterior dist

$P_{\Theta|X}(\Theta|X) / f_{\Theta|X}(\Theta|X)$

? Get via Bayes rule

prior P_Θ

Conditional
 $P_{X|\Theta}$

Obs. Process

\rightarrow

Posterior Calc

$P_{\Theta|X}(\cdot|X=x)$

Optional
Point estimate
Error analysis

①

Remember are 4 versions of Bayes rule for various combos
Continuous + discrete

Example - did today in recitation
(why math so complex - what going on again?)

- Other examples complex

Recursive inference - adding X_{n+1} observations does not
change posterior dist

- Can also have multiparameter
 - won't have closed form problems solution however
 - Sometimes will be just focused on 1 of the components

8.2 MAP

(how can there be a whole chap on this - so easy)

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{Arg\,max}} P_{\theta|x}(\theta|x)$$

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

discrete
continuous

Just the max



(7)

- maximizes the prob of a correct decision ~~given -~~
for any given value given x
- since tallest = most likely
- So only need to calculate numerator of posterior dist
(shortcut)

Point Estimation

~~Given the observed value x of X , posterior can~~

~~certain quantities that summarize posterior~~
~~single # value of θ~~

$$\hat{\theta} = \underline{\text{estimate}}$$

$$\hat{\theta} \doteq g(x) = \underline{\text{estimator}}$$

(realized value $= g(x)$ when X takes value x ,
depends on the observation)

MAP is a
estimator

Conditional Expectation estimator $\hat{\theta} = E[\theta | X=x]$

↳ Aka LMS estimator

(why is all this here - seems out of place)

(8)

If Posterior distribution θ is symmetric around its conditional mean and is unimodal (single max), Max is at mean
 $\text{MAP} = \text{LMS}$
(dvh)

if θ contains θ_{MAP} can be found analytically
- set to derivative to 0
- otherwise # search

Both MAP and LMS have no guarantees on accuracy

Hypothesis Testing

θ takes one of m values $\theta_1, \theta_2, \dots, \theta_m$
- $m=2$ when binary

Event $\{\theta = \theta_i\}$ = i^{th} hypothesis = H_i

After obs x , calculate posterior probabilities

$$P(\theta = \theta_i | X=x) = P_{\theta|x}(\theta_i | x) \text{ for each } i$$

Select θ_{mhyp} whose prob is highest (MAP)

(can calculate prob of a corresponding decision

$\theta_{\text{MAP}}(x)$ = observed hyp

$$P(\theta = \theta_{\text{MAP}}(x) | X=x)$$

⑨ Overall prob of correct decisions is

$$P(\hat{\theta} = \text{g}_{\text{MAP}}(X)) = \sum_i P(\hat{\theta} = \theta_i, X \in S_i)$$

(corresponding prob of error is

$$\sum_i P(\hat{\theta} \neq \theta_i, X \in S_i)$$

(skipping examples for now)

8.3 Bayesian Least Mean Squares Estimation (LMS)

aka the conditional expectations estimator

Simpler problem

Estimate θ w/ constant $\hat{\theta}$ w/o observation X

estimation error $= \hat{\theta} - \theta = RV$

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

- depends on $\hat{\theta}$

- can be minimized over $\hat{\theta}$

Best possible estimate $\rightarrow \hat{\theta} = E[\hat{\theta}]$

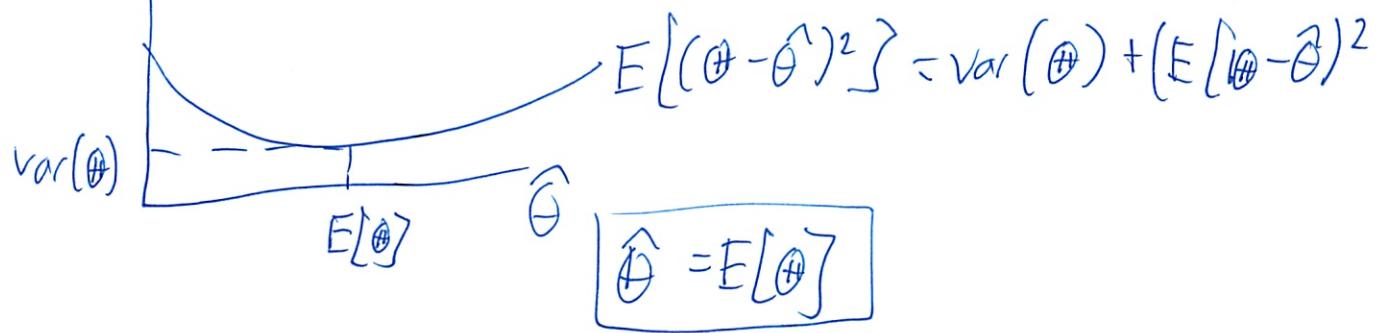
So for any estimate $\hat{\theta}$ have

$$\begin{aligned} E[(\hat{\theta} - \hat{\theta})^2] &= \text{Var}(\hat{\theta} - \hat{\theta}) + (E[\hat{\theta} - \hat{\theta}])^2 \\ &= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \hat{\theta})^2 \end{aligned}$$

(10)

- Uses $E[z^2] = \text{Var}(z) + (E[z])^2$

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



Now use observation X to estimate θ , so minimize MSE

Once know $X=x$ - now in new universe conditioned on $X=x$

So $E[\theta | X=x]$ minimizes $E[(\theta - \hat{\theta})^2 | X=x]$
conditional MSE over all constants $\hat{\theta}$

Unconditional w/ estimator

$$E[(\theta - g(x))^2]$$

$E[\theta | X]$ is estimator so minimized when $g(x) = E[\theta | X]$

Example from lecture

(11)

Some properties of estimation error

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

\uparrow \uparrow
 estimate
 estimation error

$\hat{\theta}$ is unbiased $\Leftrightarrow E[\hat{\theta}] = \theta$

$$E[\hat{\theta} | X=x] = \theta \quad \text{for all } x$$

$\hat{\theta}$ is uncorrelated w/ $\tilde{\theta}$

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

$\text{Var}(\hat{\theta})$ decomposes $\text{Var}(\theta) = \text{Var}(\hat{\theta}) + \text{Var}(\tilde{\theta})$

\uparrow if θ then info X is
Uninformative

Case of Multiple Observations and Multiple Parameters

if X is not just single RV, but vector of QVs

MSE estimation error minimized when use $E[\theta | X_1, \dots, X_n]$ as estimator

$$E[(\theta - E[\theta | X_1, \dots, X_n])^2] \leq E[(\theta - g(X_1, \dots, X_n))^2]$$

for all estimators $g(X_1, \dots, X_n)$

(12)

but hard to implement

- need a joint PDF $f_{\theta, X_1, X_2, \dots, X_n}$

- expected function hard to calculate even if have joint PDF

If want to estimate multiple parameters $\theta_1, \dots, \theta_m$

$$E[(\theta_1 - \hat{\theta}_1)^2] + \dots + E[(\theta_m - \hat{\theta}_m)^2]$$

and minimize over all estimators $\hat{\theta}_1, \dots, \hat{\theta}_m$

So decouple estimation problem

8.4 Bayesian Linear Least Mean Squares Estimation (LLMS)

- restricted class of estimators to those that are linear functions of observations
- might have higher mean squared error
- but much easier to calculate

$$\hat{\theta} = a_1 X_1 + \dots + a_n X_n + b$$

↑ scalars ↑ ↑

Corresponding MSE

$$E\left[\left(\theta - \sum_{i=1}^n a_i X_i - b\right)^2\right]$$

↑ ↑ ↑
choose these to minimize

but first where $n=1$

(3)

Single Observation

$E[(\theta - aX - b)^2]$ is associated w/ a linear estimator $aX + b$ of θ

- if a already chosen \rightarrow how chose b

$$b = E[\theta - aX] = E[\theta] - aE[X]$$

- with this b , now minimize w/ respect to a

$$E[(\theta - aX - E[\theta] + aE[X])^2]$$

- ~~partial~~

- rewrite as ~~partial~~

$$\text{Var}(\theta - aX) = \sigma_\theta^2 + a^2 \sigma_X^2 + 2\text{cov}(\theta, -aX)$$

$$= \sigma_\theta^2 + a^2 \sigma_X^2 - 2a \text{cov}(\theta, X)$$

- to min $\text{var}(\theta - aX)$, take deriv, set = 0

$$a = \frac{\text{cov}(\theta, X)}{\sigma_X^2} = \frac{\rho \sigma_\theta \sigma_X}{\sigma_X^2} = \rho \frac{\sigma_\theta}{\sigma_X}$$

- where $\rho = \frac{\text{cov}(\theta, X)}{\sigma_\theta \sigma_X}$ is Correlation Coefficient

Resulting MS estimation error of resulting linear estimator $\hat{\theta}$

$$\begin{aligned} \text{Var}(\theta - \hat{\theta}) &= \sigma_\theta^2 + a^2 \sigma_X^2 - 2a \text{cov}(\theta, X) \quad \text{given by} \\ &= \sigma_\theta^2 + \rho^2 \frac{\sigma_\theta^2}{\sigma_X^2} \sigma_X^2 - 2\rho \frac{\sigma_\theta}{\sigma_X} \rho \sigma_\theta \sigma_X \end{aligned}$$

(14)

$$= (1 - p^2) \sigma_{\theta}^2$$

Summary LLMS

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

$$= E[\theta] + p \frac{\partial \theta}{\partial x} (X - E[X])$$

$$p = \frac{\text{cov}(\theta, X)}{\sigma_{\theta} \sigma_X} \quad (\text{correlation coefficient})$$

$$\text{Resulting MSE estimation error} = (1 - p^2) \sigma_{\theta}^2$$

- easy to calculate since only mean, var, cov (θ, X)
- more intuitive

Multiple Observations + Multiple Parameters

analogous to LMS

(skipping section)

Linear Estimation + Normal Models

- different and inferior from LMS
- but if it happens to be linear - coincides w/ LLMS
(really!?! - sarcasm)

(15)

When estimation of normal RV θ based on

Observations $X_i = \theta + W_i$

\sim W_i ^{ind, 0 mean, noise terms}

Same as if original RVs were normal

Choice of Variables in Linear Estimation

- if transform $Y_i = h(X_i) \quad i=1, \dots, n$
- where Y is $1 \rightarrow 1$ then has same info
- so LMS estimator is same
- what you draw may not be good linear estimate
- so subtly transform
(confused)

end of chap

LECTURE 23

Replacement prof

- Readings:** Section 9.1
(not responsible for t -based confidence intervals, in pp. 471-473)

Outline

- Classical statistics
- Maximum likelihood (ML) estimation
- Estimating a sample mean
- Confidence intervals (CIs)
- CIs using an estimated variance

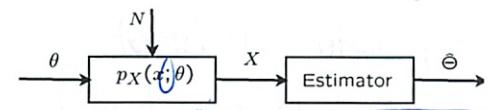
Starting chap 9

- big breakpoint

- 8+9 both about inference

- changing modeling \rightarrow unknown param unknown constants

- classical stats

Classical statisticsmodeling framework

θ = param of interest
still have probabilistic model

Use a semicolon

- also for vectors X and θ :
 $p_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta_1, \dots, \theta_m)$
- These are NOT conditional probabilities;
 θ is NOT random
 - mathematically: many models, one for each possible value of θ

From every θ , have a model for x
↳ observation model
估 estimate still a RV

Problem types:

- Hypothesis testing:
 $H_0: \theta = 1/2$ versus $H_1: \theta = 3/4$
- Composite hypotheses:
 $H_0: \theta = 1/2$ versus $H_1: \theta \neq 1/2$
- Estimation: design an estimator $\hat{\theta}$, to keep estimation error $\hat{\theta} - \theta$ small

Can address same issues as w/ Bayes stats

Maximum Likelihood Estimation

- Model, with unknown parameter(s):
 $X \sim p_X(x; \theta)$
- Pick θ that "makes data most likely"

$$\hat{\theta}_{ML} = \arg \max_{\theta} p_X(x; \theta) \quad \text{Value for } \theta \text{ that makes } X \text{ we observed most likely}$$

- Compare to Bayesian MAP estimation:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \frac{p_{X|\Theta}(x|\theta)p_{\Theta}(\theta)}{p_X(x)} = \arg \max_{\theta} p_{\Theta|X}(\theta|x) \quad \text{Maximize posterior definition}$$

- Example:** X_1, \dots, X_n : i.i.d., exponential(θ)

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

↑ drawn from many know
joint probability

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

want to maximize
over the choices
for θ

or instead $\hat{\theta}_{ML} = n/(x_1 + \dots + x_n)$

max log of the function
use cap. letters
 $\hat{\theta}_n = \frac{n}{X_1 + \dots + X_n}$) RV reciprocal of n RVs

take deriv, set = 0, max is

$$\frac{d}{d\theta} \rightarrow \frac{n}{\theta} - \sum_{i=1}^n x_i$$

$$\begin{cases} \theta e^{-\theta x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Very important estimator

- Use freq

- but only 1 possible estimator

Desirable properties of estimators

(should hold FOR ALL θ !!!)

- Unbiased:** $E[\hat{\theta}_n] = \theta$ expected value of estimator should = parameter of interest
 - exponential example, with $n = 1$: $E[1/X_1] = \infty \neq \theta$ (biased)

- Consistent:** $\hat{\theta}_n \rightarrow \theta$ (in probability)

- exponential example:

$$(X_1 + \dots + X_n)/n \rightarrow E[X] = 1/\theta$$

- can conclude that:

$$\hat{\theta}_n = n/(X_1 + \dots + X_n) \rightarrow 1/E[X] = \theta$$

- "Small" mean squared error (MSE)**

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= \text{var}(\hat{\theta} - \theta) + (E[\hat{\theta} - \theta])^2 \\ &= \text{var}(\hat{\theta}) + (\text{bias})^2 \end{aligned}$$

will be some function of θ

- non random

by weak law large # sample mean converges to the mean

General problem

Estimate a mean

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

$X_i = \theta + W_i$ *(Seq. RV fixed param of interest)*

W_i : i.i.d., mean, 0, variance σ^2

$$\hat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \dots + X_n}{n}$$

Properties:

- $E[\hat{\Theta}_n] = \theta$ (unbiased) *Elementary properties of expectation*
- WLLN: $\hat{\Theta}_n \rightarrow \theta$ (consistency)
- MSE: σ^2/n
? can compute $MSE = E[(\hat{\Theta}_n - \theta)^2] = E\left(\left(\frac{X_1 + \dots + X_n}{n} - \theta\right)^2\right) = \text{Var}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{1}{n^2}(\text{var}(X_1) + \dots + \text{var}(X_n)) = \frac{\sigma^2}{n}$
- Sample mean often turns out to also be the ML estimate.
E.g., if $X_i \sim N(\theta, \sigma^2)$, i.i.d.

$$\begin{aligned} E\left[\frac{X_1 + \dots + X_n}{n}\right] &= \frac{1}{n}[E[X_1] + \dots + E[X_n]] \\ &= \frac{1}{n} \cdot n\theta = \theta \\ E\left(\left(\frac{X_1 + \dots + X_n}{n} - \theta\right)^2\right) &= \text{Var}\left(\frac{X_1 + \dots + X_n}{n}\right) \\ &= \frac{1}{n^2}(\text{var}(X_1) + \dots + \text{var}(X_n)) \\ &= \frac{\sigma^2}{n} \end{aligned}$$

? many cases max likelihood estimates

Confidence intervals (CIs)

- An estimate $\hat{\Theta}_n$ may not be informative enough
- An $1 - \alpha$ confidence interval is a (random) interval $[\hat{\Theta}_n^-, \hat{\Theta}_n^+]$.
 - s.t. $P(\hat{\Theta}_n^- \leq \theta \leq \hat{\Theta}_n^+) \geq 1 - \alpha, \forall \theta$
 - often $\alpha = 0.05$, or 0.25 , or 0.01
 - interpretation is subtle
- CI in estimation of the mean
 $\hat{\Theta}_n = (X_1 + \dots + X_n)/n$
 - normal tables: $\Phi(1.96) = 1 - 0.05/2$

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

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

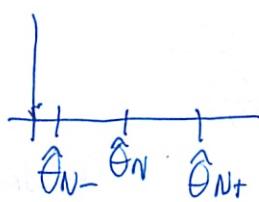
? So prob of this is ≈ 0.95

More generally: let z be s.t. $\Phi(z) = 1 - \alpha/2$

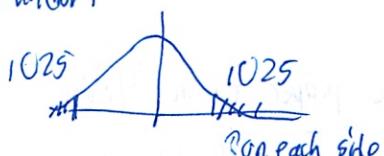
$$P\left(\hat{\Theta}_n - \frac{z\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{z\sigma}{\sqrt{n}}\right) \approx 1 - \alpha$$

How do we know if estimate is a good one?
 Very different from Bayes
 - Separate sheet

Want interval



True & not random



On each side

The case of unknown σ

Polls problem again

- Option 1: use upper bound on σ
 - if X_i Bernoulli: $\sigma \leq 1/2$
- Option 2: use ad hoc estimate of σ
 - if X_i Bernoulli(θ): $\hat{\sigma} = \sqrt{\hat{\theta}(1-\hat{\theta})}$
- Option 3: Use generic estimate of the variance

- Start from $\sigma^2 = E[(Y - \theta)^2]$

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

convergence prob by WLLN
(but do not know θ)

$$\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{\theta}_n)^2 \xrightarrow{\text{unbiased: } E[\hat{S}_n^2] = \sigma^2} \sigma^2$$

but need θ

Slight adjustment: multiply natural estimate

$$S_n^2 = \frac{n}{n-1} \cdot (\text{more obvious estimate})$$

Size of CI also depends on σ

Need to form estimate for $\text{var}(\cdot)$ to use in estimating CI

- Using option 3 \rightarrow form estimated var

An example of an exact CI

- X : exponential with parameter θ

Another example

- Analyze $[a/X, b/X]$ as a confidence interval for θ :

where have definite CI

\leftarrow reasonable estimate we considered previously

$$\begin{aligned} P\left(\frac{a}{X} \leq \theta \leq \frac{b}{X}\right) &= P\left(\frac{a}{\theta} \leq X \leq \frac{b}{\theta}\right) \leftarrow \text{elementary computation} \\ &= \int_{a/\theta}^{b/\theta} \theta e^{-\theta x} dx \\ &= -e^{-\theta x} \Big|_{x=a/\theta}^{x=b/\theta} = \frac{e^{-a} - e^{-b}}{\theta} \end{aligned}$$

No dependence on θ , so have a confidence interval

- Example: $\left[\frac{1}{4X}, \frac{4}{X}\right]$ is a 0.76 confidence interval ("76% CI")

- We know how to compute $P(X \in [\alpha, \beta])$

- Rearrange so prob in that form

confidence interval

no dependence on θ

if tried $[a\sqrt{X}, b\sqrt{X}]$ - would not have gotten prob ind. of θ

Much chose proper form "trick" to find CI not dependent on θ - could use instead CLT and bound on var

or θ - could use instead CLT and bound on var

bounds
on Var
- usually
- for
Bernoulli

Exponential example

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

$$\hat{\theta}_1 = \frac{1}{x_1}$$

$$E[\hat{\theta}_1] = E\left[\frac{1}{x_1}\right]$$

$$= \int_0^\infty \frac{1}{x} \theta e^{-\theta x} dx$$

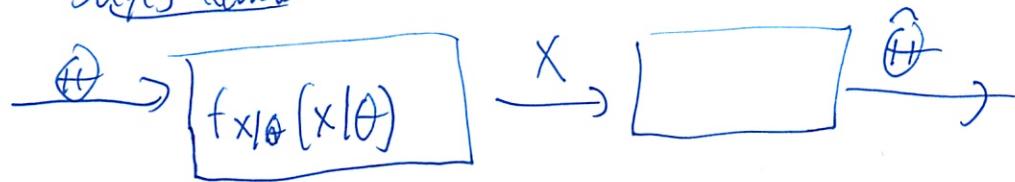
~~$\stackrel{x}{\rightarrow} \theta$~~ No

The S diverges for all θ values of θ

$= \infty$

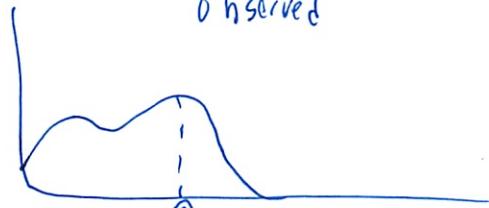
So this one is biased

(2)

Bayes Review

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

? little x
fixed size
observed



Boil down dist to 1 estimate, like the $\hat{\theta}_{MAP}$ estimate

(an ask qv abt quality of estimate

$E[\downarrow(\theta, \hat{\theta}_{MAP})]$ can compute error

But classical setting different θ is not random

Can give a confidence interval

1.5
2.0
0.5

Problem Set 10

Due December 2, 2010 Thur

1. **A financial parable.** An investment bank is managing \$1 billion, which it invests in various financial instruments ("assets") related to the housing market (e.g., the infamous "mortgage backed securities"). Because the bank is investing with borrowed money, its actual assets are only \$50 million (5%). Accordingly, if the bank loses more than 5%, it becomes insolvent. (Which means that it will have to be bailed out, and the bankers may need to forgo any huge bonuses for a few months.)
 - (a) The bank considers investing in a single asset, whose gain (over a 1-year period, and measured in percentage points) is modeled as a normal random variable R , with mean 7 and standard deviation 10. (That is, the asset is expected to yield a 7% profit.) What is the probability that the bank will become insolvent? Would you accept this level of risk?
 - (b) In order to safeguard its position, the bank decides to diversify its investments. It considers investing \$50 million in each of twenty different assets, with the i th one having a gain R_i , which is again normal with mean 7 and standard deviation 10; the bank's gain will be $(R_1 + \dots + R_{20})/20$. These twenty assets are chosen to reflect the housing sectors at different states or even countries, and the bank's rocket scientists choose to model the R_i as independent random variables. According to this model, what is the probability that the bank becomes insolvent?
 - (c) Based on the calculations in part (b), the bank goes ahead with the diversified investment strategy. It turns out that a global economic phenomenon can affect the housing sectors in different states and countries simultaneously, and therefore the gains R_i are in fact positively correlated. Suppose that for every i and j where $i \neq j$, the correlation coefficient $\rho(R_i, R_j)$ is equal to 1/2. What is the probability that the bank becomes insolvent? You can assume that $(R_1 + \dots + R_{20})/20$ is normal.
- (2) The adult population of Nowhereville consists of 300 males and 196 females. Each male (respectively, female) has a probability of 0.4 (respectively, 0.5) of casting a vote in the local elections, independently of everyone else. Find a good numerical approximation for the probability that more males than females cast a vote.
3. Let S_n be the number of successes in n independent Bernoulli trials, where the probability of success in each trial is $p = \frac{1}{2}$. Provide a numerical value for the limit as n tends to infinity for each of the following three expressions:
 - (a) $P\left(\frac{n}{2} - 10 \leq S_n \leq \frac{n}{2} + 10\right)$
 - (b) $P\left(\frac{n}{2} - \frac{n}{10} \leq S_n \leq \frac{n}{2} + \frac{n}{10}\right)$
 - (c) $P\left(\frac{n}{2} - \frac{\sqrt{n}}{2} \leq S_n \leq \frac{n}{2} + \frac{\sqrt{n}}{2}\right)$
4. Alice has two coins. The probability of heads for the first coin is 1/3; the probability of heads for the second coin is 2/3. Other than this difference in their bias, the coins are indistinguishable through any measurement known to man. Alice chooses one of the coins randomly and sends it to Bob. Let p be the probability that Alice chose the first coin. Bob tries to guess which of the two coins he received by flipping it 3 times in a row and observing the outcome. Assume that all coin flips are independent. Let Y be the number of heads Bob observed.

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
 Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
 (Fall 2010)

- (a) Given that Bob observed k heads, what is the probability that he received the first coin?
 - (b) Find values of k for which the probability that Alice sent the first coin increases after Bob observes k heads out of 3 tosses. In other words, for what values of k is the probability that Alice sent the first coin given that Bob observed k heads greater than p ? If we increase p , how does your answer change (goes up, goes down, or stays unchanged)?
 - (c) Help Bob develop the rule for deciding which coin he received based on the number of heads k he observed in 3 tosses if his goal is to minimize the probability of error.
 - (d) For this part, assume $p = 2/3$.
 - i. Find the probability that Bob will guess the coin correctly using the rule above.
 - ii. How does this compare to the probability of guessing correctly if Bob tried to guess which coin he received before flipping it?
 - (e) If we increase p , how does that affect the decision rule?
 - (f) Find the values of p for which Bob will never guess he received the first coin, regardless of the outcome of the tosses.
 - (g) Find the values of p for which Bob will always guess he received the first coin, regardless of the outcome of the tosses.
5. Consider a Bernoulli process X_1, X_2, X_3, \dots with unknown probability of success q . As usual, define the k th inter-arrival time T_k as

$$T_1 = Y_1, \quad T_k = Y_k - Y_{k-1}, \quad k = 2, 3, \dots$$

where Y_k is the time of the k th success. This problem explores estimation of q from observed inter-arrival times $\{t_1, t_2, t_3, \dots\}$.

You may find the following integral useful: For any non-negative integers k and m ,

$$\int_0^1 q^k (1-q)^m dq = \frac{k! m!}{(k+m+1)!}$$

Assume q is sampled from the random variable Q which is uniformly distributed over $[0, 1]$.

- (a) Compute the PMF of T_1 , $p_{T_1}(t_1)$
- (b) Compute the least squares estimate (LSE) of Q from the first recording $T_1 = t_1$.
- (c) Compute the maximum a posteriori (MAP) estimate of Q given the k recordings, $T_1 = t_1, \dots, T_k = t_k$.

For this part only assume q is sampled from the random variable Q which is now uniformly distributed over $[0.5, 1]$

- (d) Find the linear least squares estimate (LLSE) of the second inter-arrival time (T_2), from the observed first arrival time ($T_1 = t_1$).

6. The joint PDF of X and Y is defined as follows:

$$f_{X,Y}(x,y) = \begin{cases} cxy & \text{if } 0 < x \leq 1, 0 < y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
(Fall 2010)

- (a) Find the normalization constant c .
- (b) Compute the conditional expectation estimator of X based on the observed value $Y = y$.
- (c) Is this estimate different from what you would have guessed before you saw the value $Y = y$? Explain.
- (d) Repeat (b) and (c) for the MAP estimator.

OH! P-Set nothing new
→ just estimators
→ Bayes rule

I. A financial parable Investment bank

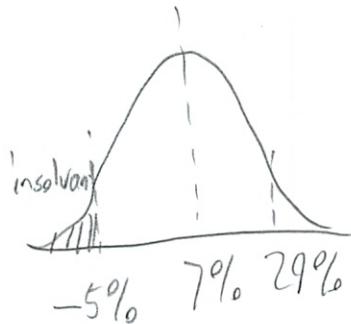
- managing \$1 billion in various housing assets
- only \$50 million / 5% actual assets
- so if it loses more than 5% it may be insolvent

a) Bank considers investing in a single asset

- gain over 1 year in % points is normal RV R
 $\mu = 7\%$ $\sigma = 10\%$

? Expected to yield a 7% profit

- what is prob that bank will become insolvent?



- ? Simple limit theorem / lecture 19 problem

- Well depends which method one uses

- Chebyshev

- Markov

- No! One RV, so can just use CDF

$$\Pr(R \leq -5\%) = \Phi\left(\frac{R - E[X]}{\sqrt{\text{Var}(X)}}\right)$$

$$\begin{aligned} \text{Var} &= \sigma^2 \\ \sqrt{\text{Var}} &= \sigma \end{aligned}$$

$$\begin{aligned}
 ② &= \Phi\left(\frac{R-7}{10}\right) \\
 &= \Phi\left(\frac{-5-7}{10}\right) \\
 &= \Phi(-1.2) \\
 &= 1 - \Phi(1.2) \\
 &= 1 - 0.8841 \\
 &= 0.1151
 \end{aligned}$$

11% risk of failure

- a bit more than 1 st dev which is ~15% - so makes sense

Would I bet the firm on this?

No - I am a more conservative investor.

b) In order to safeguard its position, bank diversifies

\$50 million in 20 different assets

R_i = ith one's gain

↳ each one mean 7 st dev 70

Bank's total gain rate = $\frac{R_1 + \dots + R_{20}}{20} = p$

bank's "rocket scientists" model each R_i as ind

Now what is prob bank is insolvent,

(I like this problem)

(3) Here limit theory of sums

Figuring how to do
grader skip pg

- Chebchier ?
- Markov ?
- Something more exact?
- Central Limit Theorem
 - Most exact
 - is it perfectly exact or just more exact?
- Is like Recitation 21 #1
- CLT is actual estimate, not exact though
- Chebisher + Markov are upper bounds
- How do actual again?

WP: Convolution of their densities

↳ alpha cross correlation

the shifting function

↳ common area

TB: Sum independent RVs

(4)

$$\begin{aligned}
 P(P \leq -5) &= P\left(\frac{P - E[R]}{\sqrt{\text{Var}(R)}} \leq \frac{-5 - E[R]}{\sqrt{\text{Var}(R)}}\right) \\
 &= P\left(\frac{P - E[R]}{\sqrt{\text{Var}(R)}} \leq \frac{-5 - 7}{10}\right) \\
 &= \Phi(-1.2) \\
 &= 1 - \Phi(1.2) \\
 &= 1 - .8849 \\
 &= .1151
 \end{aligned}$$

is same as before

Was not expecting, but makes sense since noise on both sides of mean

This is for ∞ n I believe

well $\lim_{n \rightarrow \infty}$

(5)

$$S_n = P \text{ here}$$

$$Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}}$$

$$= \frac{S_n - nE[X]}{\sqrt{n}\sigma} \quad \text{e.g. when sum of RVs}$$

$$P(P \leq -5) = P\left(\frac{P - nE[R_i]}{\sqrt{n}\sqrt{\text{Var}(R_i)}} \leq \frac{-5 - nE[R_i]}{\sqrt{n}\sigma_{R_i}}\right)$$

$$= P\left(\frac{P - nE[R_i]}{\sqrt{\text{Var}(R_i)}} \leq \frac{-5 - 20 \cdot 7}{\sqrt{20} \cdot 10}\right)$$

$$= \Phi(-3, 24)$$

Well wait $M_n = R$

and $S_n = nM_n$

$$= 1 - \Phi(3, 24)$$

$$= 1 - 0.9994$$

$$= 0.0006$$

Bank will almost certainly not go bankrupt
 Much more believable outcome, though I would not think
 that small

Remember Chebachev better than Markov!

c) Based on b) bank will diversify. But turns out QVs
are not independent \rightarrow positively correlate.

Suppose that for every i and j where $i \neq j$
the correlation coefficient $\rho(R_i, R_j) = \frac{1}{2}$

What is prob bank will go insolvent?

? Did we ever study this?

- CLT only when ind
- Markov + Cheb dep or ind
- but Cheb allows -, so can use

- p 220 chap 4

- must be careful to compute sum

- sum of var + sum of cov

both ways must count

- can derive, not in book

$$\text{Var}(X_1 + \dots + X_n) = E[(X_1 + \dots + X_n)^2] - [E(X_1 + \dots + X_n)]^2$$

Should simplify to expression w/ vars + cov

- last then from correlation coefficient to cov
sentence gives us a critical assumption

↳ that is normal (~~so no CLT~~) \rightarrow just normal CDF

$$\textcircled{7} \quad P(X, Y) = \frac{1}{2} \rightarrow \text{by definition}$$

$$\frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \quad \text{Can back at cov}$$

? mnemonic $\frac{\sigma_{xy}}{\sigma_x \sigma_y}$

$$\text{Var}(X) = \sigma_x^2$$

? st dev

$$\sigma_{xy} \in \text{Cov}(X, Y)$$

$$\text{Since } \text{Var}(X) = \text{Cov}(X, X)$$

$$E\left\{ \left(\sum_{i=1}^n X_i^2 + \dots \right) \right\}$$

? like terms

all shown on p 220

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) + \sum_{(i,j) | i \neq j} \text{Cov}(X_i, X_j)$$

? double summation

$$\sum_{i=1}^{20} \sum_{j=1}^{20}$$

$i \neq j$

so $\text{Cov}(X_1, X_2) + \text{Cov}(X_1, X_3) + \dots + \text{Cov}(X_2, X_1) + \text{Cov}(X_2, X_3)$

$\text{Cov}(X_1, X_1)$

? double count

(8)

Back out correlation coefficient

$$\frac{1}{2} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

$$\text{Var}(R) = \sigma^2 = 10^2$$

$$\sqrt{\text{Var}} = \sigma = 10$$

$$\frac{1}{2} = \frac{\text{Cov}(x, y)}{\sqrt{100 \cdot 100}}$$

$$50 = \text{Cov}(x, y)$$

$$\begin{aligned}\text{Var}\left(\sum_{i=1}^{20} X_i\right) &= \sum_{i=1}^{20} 100 + \sum_{(i,j) | i \neq j} 50 \\ &= 20 \cdot 100 + 19 \cdot 19 \cdot 50 \\ &= 20,050\end{aligned}$$

Var for each

 $E[\cdot]$ unchanged?

⑨

$$\begin{aligned}
 P(P \leq -5) &= P\left(\frac{P - nE[R_i]}{\sqrt{nVar(R_i)}} \leq \frac{-5 - 20}{\sqrt{20} \sqrt{20050}}\right) \\
 &= \Phi(-.2289) \quad \text{ha-actually worked out fairly well} \\
 &= 1 - \Phi(.2289) \\
 &= 1 - .5871 \\
 &= .4129
 \end{aligned}$$

Makes sense

- much higher
- but is it too high?

thankfully can't do this comply on final

(10)

2. Adult population of Nowhere ville

$$= 300 \text{ males}, 4 \text{ prob to vote}$$

$$= 196 \text{ females}, 5 \text{ prob to vote}$$

Find a good \approx approx for prob more males vote than females

$$M_i = \begin{cases} 1 & \text{yes} \\ 0 & \text{no} \end{cases} \quad \begin{cases} .4 & \text{male} \\ .6 & \text{"} \end{cases} \quad \begin{cases} .5 & \text{female} \\ .5 & \text{"} \end{cases}$$

$$M = M_1 + M_2 + \dots + M_i \quad \leftarrow \text{sum of } \stackrel{\text{ind}}{\text{Bernoulli's}} \leftarrow \text{Binomial}$$

$$F = F_1 + F_2 + \dots + F_i$$

$$P(M > F)$$

- So comparing 2 RVs

Can do p certain amt will vote, say 21

$$\text{exact } P(S_n \leq .21) = \sum_{k=0}^{21} \binom{300}{k} (.4)^k (.6)^{300-k} \leftarrow \text{because binomial}$$

(LT)

$$\approx P\left(\frac{S_n - E[S_n]}{\sqrt{Var(S_n)}} \leq \frac{.21 - 300 \cdot .4}{\sqrt{300 \cdot .4 \cdot .6}}\right) \quad \begin{aligned} E[S_n] &= np \\ Var(S_n) &= np(1-p) \end{aligned} \quad \text{Binomial}$$

$$\approx \Phi(-1.35)$$

but then usually add $\frac{1}{2}$

(11)

But want $P(M > F) = 1 - P(M \leq F)$

$$1 - P\left(\frac{M - E[M]}{\text{Var}(M)} \leq \frac{F - E[F]}{\text{Var}(F)}\right)$$

note this is S_n
so its not $nE[X]$

$$E[M] = np = 300 \cdot .4 = 120$$

$$\text{Var}(M) = np(1-p) = 300 \cdot .4 \cdot (.1-.4) = 72$$

$$E[F] = np = 196 \cdot .5 = 98$$

$$\text{Var}(F) = np(1-p) = 196 \cdot .5 \cdot (.1-.5) = 49$$

$$1 - P\left(\frac{M - 120}{72} \leq \frac{F - 98}{49}\right)$$

so now need dist of this

well want a single #

but how translate to st normal table

did we ever do a problem like this?

or do something else

+0.5

- both are normal distributions

- want overlap

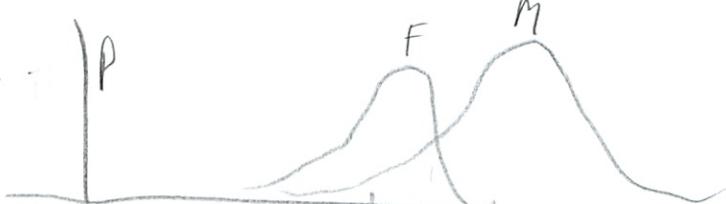
- No want prob of overlap

- well are prob dist

- but what graphically
is the ans?

- some weird convolution thing

- or add them into some new RV



Tile slanty piecing together

(12)

Sum of ind Bernoullis - almost normal

~~Could figure Mean, var M, F~~

So could do joint - had

$$\text{So do } V = M - F$$

And find distribution $P(V \geq 0)$

duh - should have thought of this
- was leaning towards

+0.5

- sum of ind normal is normal
↳ any linear comb

CLT - sum of any iid, standardised are approx normal

$E[V] = \text{expectation of sum of RV}$
linearity of expectation

$$= E[M] - E[F]$$

$$= 120 - 98$$

$$= 22$$

+0.5

$\text{Var}(V) =$ oh right was looking for this on Nov 20 P-set

$$= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

Find so 0

$$= 72 + 44$$

$$= 121$$

$$P(V \geq 0) = P\left(\frac{V - E[V]}{\text{Var}(V)} \geq \frac{0 - 22}{121}\right) \quad (-0.5)$$

(13)

$$= \Phi(-18)$$

$$= 1 - \Phi(18)$$

$$= 1 - 5714$$

$$\approx .4286$$

✓ Makes sense

1.5/2

(14)

3. Let S_n be # of successes in n ind. Bernoulli trials
where $p(\text{success}) = \frac{1}{2}$

Provide a numerical value as $n \rightarrow \infty$

This is DeMoivre - Laplace Approx to Binomial

- we did not go over much

$$\begin{aligned} M = E[X_i] &= p \\ &= \frac{1}{2} \end{aligned} \quad \begin{aligned} \sigma = \sqrt{\text{Var}(X_i)} &= \sqrt{p(1-p)} \\ &= \sqrt{\frac{1}{2}\left(1-\frac{1}{2}\right)} \\ &= \frac{1}{2} \end{aligned}$$

We will use approx from central limit theorem to provide approx for the prob $\{k \leq S_n \leq l\}$

Convert into a standardized RV using the equivalence

$$\frac{k-np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{l-np}{\sqrt{np(1-p)}}$$

$$P(k \leq S_n \leq l) = P(\text{ })$$

$$= \Phi\left(\frac{l-np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k-np}{\sqrt{np(1-p)}}\right) \quad \begin{array}{l} \text{then add } \frac{1}{2} \\ \text{for better estimate} \end{array}$$

CDF of one - the other

(15)

a) $P\left(\frac{n}{2} - 10 \leq S_n \leq \frac{n}{2} + 10\right)$

as $n \rightarrow \infty$
this gets larger

the n in there throws
a wrench in my plans

$$\Phi\left(\frac{\frac{n}{2} + 10 + \frac{1}{2} - n^{\frac{1}{2}}}{\sqrt{n^{\frac{1}{2}}(\frac{1}{2})}}\right) - \Phi\left(\frac{\frac{n}{2} - 10 - \frac{1}{2} - n^{\frac{1}{2}}}{\sqrt{n^{\frac{1}{2}}(\frac{1}{2})}}\right)$$

$$\Phi\left(\frac{5}{\sqrt{\frac{n}{4}}}\right) - \Phi\left(\frac{-5}{\sqrt{\frac{n}{4}}}\right) \quad \sqrt{\frac{n}{4}} = \frac{\sqrt{n}}{2}$$

As $n \rightarrow \infty$, denominator grows to ∞

so $\Phi(\) \rightarrow 0$, so value \rightarrow (15)

(-1)



(16)

$$b) P\left(\frac{1}{2} - \frac{n}{10} \leq S_n \leq \frac{1}{2} + \frac{n}{10}\right)$$

$$\Phi\left(\frac{\frac{n}{2} + \frac{n}{10} + \frac{1}{2} - n\frac{1}{2}}{\sqrt{n\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)}}\right) - \Phi\left(\frac{\frac{1}{2} - \frac{n}{10} - \frac{1}{2} - n\frac{1}{2}}{\sqrt{n\frac{1}{2}\left(\frac{1}{2}\right)}}\right)$$

$$\Phi\left(\frac{\frac{n}{10} + \frac{1}{2}}{\sqrt{n\frac{1}{4}}}\right) - \Phi\left(\frac{-\frac{n}{10} - \frac{1}{2}}{\sqrt{n\frac{1}{2}\left(\frac{1}{2}\right)}}\right)$$

So again denominator $\rightarrow \infty$) as $n \rightarrow \infty$

But now numerator also goes $\rightarrow \infty$

So $\Phi(\infty)$ which is indeterminate

But put whole thing into wolfram alpha $\Phi(\cdot) \rightarrow \infty$

Which leads value $\rightarrow 1$

Other side $\Phi(\cdot) \rightarrow -\infty$

So this is $1 - \Phi(\infty)$ which is $1 - 1 = 0$

$$\text{So } 1 - 0 = \textcircled{1}$$

+ |

(17)

$$P\left(\frac{n}{2} - \frac{\sqrt{n}}{2} \leq S_n \leq \frac{n}{2} + \frac{\sqrt{n}}{2}\right)$$

$$\Phi\left(\frac{\frac{n}{2} + \frac{\sqrt{n}}{2} + \frac{1}{2} - n\frac{1}{2}}{\sqrt{n\frac{1}{2}\left(\frac{1}{2}\right)}}\right) - \Phi\left(\frac{\frac{n}{2} - \frac{\sqrt{n}}{2} - \frac{1}{2} - n\frac{1}{2}}{\sqrt{n\frac{1}{2}\left(\frac{1}{2}\right)}}\right)$$

$$\Phi\left(\frac{\frac{\sqrt{n}}{2} + \frac{1}{2}}{\sqrt{\frac{n}{4}}}\right) - \Phi\left(\frac{-\frac{\sqrt{n}}{2} - \frac{1}{2}}{\sqrt{\frac{n}{4}}}\right)$$

Denom still $\rightarrow \infty$

Numerator $\rightarrow \infty$

Which is indeterminate

But whole thing together $\rightarrow 1$ (wolfram alpha)

$$\text{So } \Phi(1) \rightarrow ,8413$$

Other one together $\rightarrow -1$

$$\text{So } 1 - \Phi(1) = 1 - ,8413$$

$$\text{So } ,8413 - (1 - ,8413)$$

(16826)

+1

~~2/3~~

(18)

Q. Alice has 2 coins

$$P(\text{heads 1st coin}) = \frac{1}{3}$$

$$P(\text{heads 2nd coin}) = \frac{2}{3}$$

Alice chooses coin randomly + sends to Bob.

P = prob Alice picked first coin
 $(\frac{1}{2}, \text{right})$

Bob tries to guess by flipping his coin 3x

 γ = # of heads he sees.

Simple Bayesian problem

Use CLT
 then $\lim_{n \rightarrow \infty}$

 θ = which coin chosen

$$P_\theta = \begin{cases} \frac{1}{2} \text{ coin A} & -\text{well} \\ \frac{1}{2} \text{ coin B} & \end{cases} \quad \begin{cases} p & \text{coin A} \\ 1-p & \text{coin B} \end{cases}$$

But it said coins indistinguishable, so we can

assume $p = \frac{1}{2}$

$$P_\gamma | \theta = \begin{cases} \text{for } \theta = A \\ \frac{1}{3} \text{ heads} \\ \frac{2}{3} \text{ tails} \end{cases} \quad \begin{cases} \text{for } \theta = B \\ \frac{2}{3} \text{ heads} \\ \frac{1}{3} \text{ tails} \end{cases}$$

Just an application of Bayes Rule

- define RV or events

- Alice gives bob coin B

DOL

In general $P(A_i | Y=k) =$
 $\frac{\text{# heads you got}}{\text{got 1st coin}}$

$$= \frac{P(A_i) P(Y=k | A_i)}{\sum_{i=1}^n P(A_i) P(Y=k | A_i)} \leftarrow P(Y=k)$$

total prob theorem

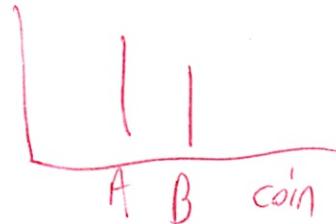
$$\Theta = \begin{cases} 1 & \text{1st coin} \\ 2 & \text{2nd coin} \end{cases} \leftarrow \text{just syntax}$$

Not implied that $p=\frac{1}{2}$

Which estimator MAP, LMS, LLMS

? not good for our purposes

$P_{\theta|Y}(\theta | k)$ ↪ just have 2 #



(oh so I was on right track)

So conditional prob of error = prob of other

So just do $\hat{\theta}_{MAP} = \arg \max P_{\theta|x}(\theta | x)$

Depends on p in C

$$\frac{P_\theta(\theta) P_{x|\theta}(x|\theta)}{\sum P_\theta(\theta) P_{x|\theta}(x|\theta)}$$

$\left. \begin{array}{l} \text{don't have} \\ \text{to worry} \\ \text{about denom} \end{array} \right\}$

②

$$P_0(1) P_{X|0}(x|1) \geq P_0(2) P_{X|0}(x|2)$$

↑ pick coin A

- in terms of k, p

- now if give φ - what it will be

- but can do which is bigger w/ fixed φ

b) right side in terms of k

both sides in terms of p

for which values of k is this true

(21)

? into chap 8 now

a) Given that he saw 3 heads, prob that received st
Coin

Y is binomial if $\begin{cases} \theta = A & p = \frac{1}{3} \\ \theta = B & p = \frac{2}{3} \end{cases}$ $n=3$

So must choose b/w those two models,

Now what is the simple math for this?

-remember 6.01

$$f_{\theta|x}(\theta|x) = \frac{f_{\theta}(\theta) f_{x|\theta}(x|\theta)}{\sum f_{\theta}(\theta) f_{x|\theta}(x|\theta)}$$

So first if $\theta = A$

$$\frac{1}{2} \cdot \sum_{k=0}^3 \binom{3}{k} \frac{1}{3}^k \left(\frac{2}{3}\right)^{3-k}$$

made bad
assumption, $p = \frac{1}{2}$

$$\frac{1}{2} \cdot \left(3 \cdot \frac{1}{3} \left(\frac{2}{3}\right)^2 + \frac{3!}{2!1!} \cdot \left(\frac{1}{3}\right)^2 \left(\frac{2}{3}\right) + 1 \cdot \left(\frac{2}{3}\right)^3 \right)$$

$$\frac{1}{2} \left(\frac{4}{9} + \frac{2}{9} + \frac{8}{27} \right)$$

1481

$$\theta = B \quad \frac{1}{2} \cdot \sum_{k=0}^3 \binom{3}{k} \frac{2}{3}^k \left(\frac{1}{3}\right)^{3-k}$$

$$\frac{1}{2} \left(3 \cdot \frac{2}{3} \left(\frac{1}{3}\right)^2 + 3 \cdot \left(\frac{2}{3}\right)^2 \left(\frac{1}{3}\right) + 1 \left(\frac{1}{3}\right)^3 \right)$$

(28)

$$\frac{1}{2} \left(\frac{2}{9} + \frac{4}{9} + \frac{1}{27} \right)$$

$$13518$$

$$\text{Sum } ,481 + ,3518 = ,83 = \frac{5}{6}$$

So $\Theta = \begin{cases} A & ,57 \\ B & ,42 \end{cases}$

There - I think that was the manually way

But what math shortcuts are there?

But wait the answer should be based on k heads

So what did I do?

Possibilities for Θ , for any/all values of k

So as function of k

$$\frac{\frac{1}{2} \cdot \binom{3}{k} \left(\frac{1}{3}\right)^k \left(\frac{2}{3}\right)^{3-k}}{\binom{3}{k} \left(\frac{1}{3}\right)^k \left(\frac{2}{3}\right)^{3-k} + \binom{3}{k} \left(\frac{2}{3}\right)^k \left(\frac{1}{3}\right)^{3-k}}$$

and this is only for 1st one

After OH

So 157 was right for map

$$\hat{\Theta}_{\text{MAP}} = A$$

(23) b) Find values for k for which the probability that Alice sent first coin increases after Bob observes k heads out of 3 tosses,

Aka for what values of k is the probability (Alice sent 1st coin | Bob observed k heads)

If $\neq p$, how does answer change?

So if $k=0$, then likely to have 1st coin "A"
 $k=3$, " " " " " 2nd coin "B"

But specifically what is prob

So for what values of k is

$$\frac{\frac{1}{2} \text{ or not } p \binom{3}{k} \left(\frac{1}{3}\right)^k \left(\frac{2}{3}\right)^{3-k}}{\binom{3}{k} \left(\frac{1}{3}\right)^k \left(\frac{2}{3}\right)^{3-k} + \binom{3}{k} \left(\frac{2}{3}\right)^k \left(\frac{1}{3}\right)^{3-k}} \rightarrow p$$

$$\text{So } k=0 \approx 0.8$$

$$1 \approx 0.6$$

$$2 \approx 0.3$$

$$3 \approx 0.1$$

$$w/p = 1$$

Since p is on the same side, changing it won't affect anything!

(24)

No - that does not make sense.

And tested it when $p = .25$ $p = .5$

$$\begin{array}{ll} k=0 & \approx .22 \\ 1 & \approx .16 \\ 2 & \approx .09 \\ 3 & \approx .01 \end{array}$$

$$\begin{array}{ll} k=0 & \approx .45 \\ 1 & \approx .32 \\ 2 & \approx .17 \\ 3 & \approx .07 \end{array}$$

? No this is wrong - answer is always larger than p

That is just numeric value for left hand side

? Supposed to be a PMF

Perhaps should do via an estimator

After OH

For which values of k is this true

? None of them - but is not right

Recur in wolfram alpha

$$p^9 \cdot 2^{-k-1} (3-k)! / k! > p$$

$k=0$	$\frac{27}{8} p > p$	✓	so here A is always
$k=1$	$\frac{9}{2} p > p$	✓	$>$ than p should
$k=2$	$\frac{9}{4} p > p$	✓	not be true
$k=3$	$\frac{27}{8} p > p$	✗	-asymptotic to 0

25

? Try out estimating function

So we have $Y = k$ $Y = F_1 + F_2 + F_3 = S_m$
 binomial bernulli

So we have sum $\Rightarrow CLT$?

$$Z_n = \frac{S_n - E[S_n]}{\sqrt{\text{Var}[S_n]}}$$

$E[S_n] = np$ but p varies w/ which
binomial one is chosen

$$\text{Var}(S_n) = np(1-p)$$

What is the problem even asking?

I remember how we did this in 6.01

$$Z_n = \frac{S_3 - 3p}{\sqrt{3p(1-p)}} \text{ but we are asked to find } p$$

Unless it's a combo of the two
-average

$$P_{\text{above}} = \begin{cases} \frac{1}{3} & \text{w/ prob } p \\ \frac{2}{3} & \text{w/ prob } (1-p) \end{cases}$$

$$P_{\text{above}} = \frac{1}{3}P + \frac{2}{3}(1-P)$$

(26)

$$Z_n = \frac{S_3 - 3\left(\frac{1}{3}p + \frac{2}{3}(1-p)\right)}{\sqrt{3\left(\frac{1}{3}p + \frac{2}{3}(1-p)\right)(1 - \left(\frac{1}{3}p + \frac{2}{3}(1-p)\right)}}}$$

Then can use LMS to find prob \rightarrow something

But should this not now be fn of k .

Well $S_3 = k$
 \approx actual value
 RV

$$\begin{aligned} P(S_3 \geq 0) &= \Phi\left(\frac{0 - \dots}{\dots}\right) \\ &= \Phi\left(\frac{-3}{p+1}\right) \end{aligned}$$

depends on p - can't use tables

? But can't find $\Phi\left(\frac{-3}{p+1}\right) \rightarrow P(p)$

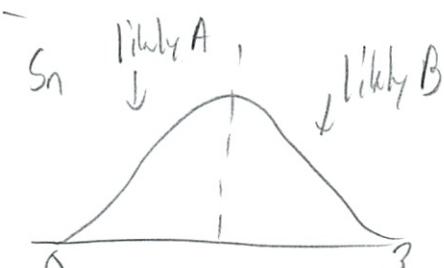
$$\begin{aligned} P(S_3 \geq 1) &= \Phi\left(1 - \dots\right) \\ &= \Phi\left(\text{(complicated)}\right) \end{aligned}$$



(27)

c) Help Bob develop the rule for deciding which coin he received based on # of heads^k he sees in 3 tosses

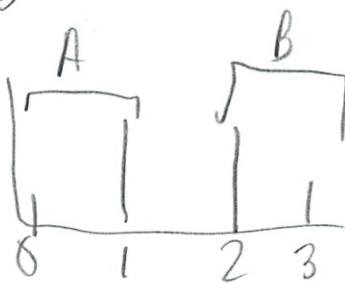
- To minimize prob of error = MAP
- well any estimator
- Well depends on p , but should be able to answer in terms of p



Perhaps not split

Since we don't know p

So



And this uses MAP \rightarrow tallest

(28)

I can't sort out all these variables in my head

d) Assume $p = \frac{2}{3}$

i) Find prob Bob will guess correctly

- So what guess is this again?

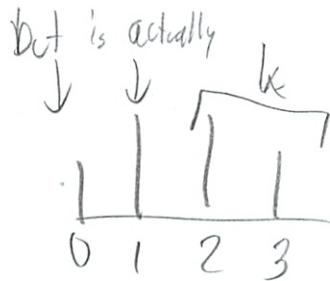
- After 3 observations

- Correct means after 3 data matches actual

- Wrong is error ??

- So MSE ??

- Error is if says A, and is B



$$\text{So } p P(k=2, k=3) + (1-p) P(k=0, k=1)$$

(29)

ii) How does this compare w/ prob of guessing correctly if Bob tried to guess which coin he recalled before flipping it.

Well before flipping it, he only knew he had p chance of getting A.

Now the certainty of knowing which is which is higher w/ more info.

(30)

e) If we $\cap p$, how does this affect the decision rule?

It will tilt the decision to one of the coins (the one more likely to be picked). This will affect the observed distribution which Bob sees.

(31)

f) Find the value of p for which Bob will never guess he received the lost coin, regardless of toss outcome.

If p is very small,

like if p is 0, he will likely see 2 heads and declare his coin to be B, but there is a chance he will see only 1 head

$$\binom{3}{1} \left(\frac{2}{3}\right)^1 \left(\frac{1}{3}\right)^2 \text{ if } B + \binom{3}{0} \left(\frac{2}{3}\right)^0 \left(\frac{1}{3}\right)^3$$

prob $\frac{1}{3}$ that will be only 1 head w/
coin B

So this is not possible, because with only 3 tosses, still a fairly significant chance will declare coin B to be A

(32)

9) Now always guess 1st coin

I am guessing this is the same, but reverse

If $p=1$ would always be A

But then chance to have 2,3 heads

$$\binom{3}{2} \left(\frac{1}{3}\right)^2 \left(\frac{2}{3}\right)^1 + \binom{3}{3} \left(\frac{1}{3}\right)^3 \quad (\text{if } p=1)$$

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

Same $\frac{1}{3}$ chance he will say its B and
be wrong

So no good value for p

(33)

5. Consider Bernoulli process $X_1, X_2, X_3 \dots$

w/ unknown prob of success q

k th interarrival time $T_k =$

$$T_1 = Y_1 \quad T_k = Y_k - Y_{k-1} \quad k=2,3,\dots$$

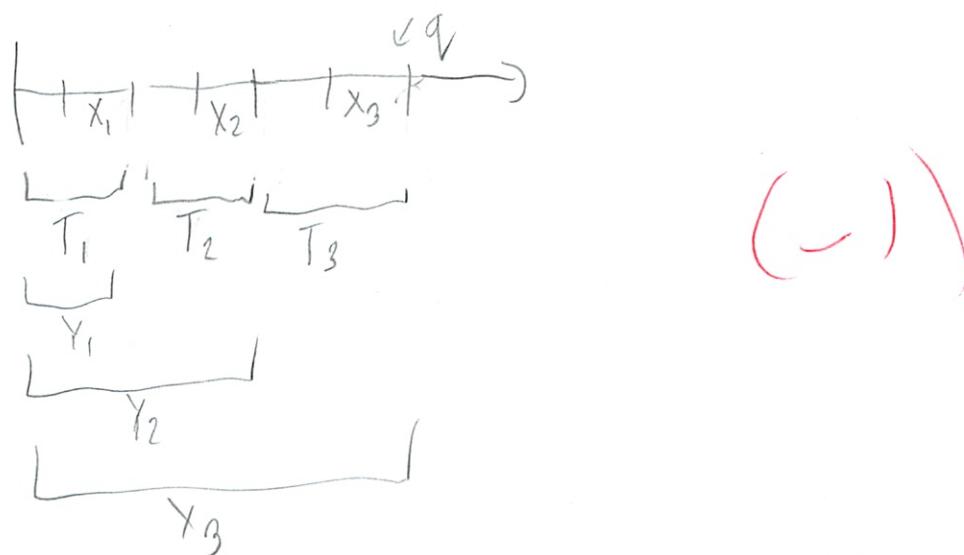
where Y_k is time of k th success

Explores estimation of q from observed inter-arrival time

Assume q is sampled from RV Q , which is uniformly distributed $[0,1]$

Hint $\int q^k (1-q)^m dq = \frac{k! m!}{(k+m+1)!}$

a) Compute PMF of T_1 , $p_{T_1}(t_1)$



(34)

b) Compute LSE of Q from 1st recording $T_1 = t_1$
 So now have one reading

? Draw joint



$$\begin{aligned}\hat{\theta}_{\text{LMS}} &= E[\theta | x=x] \\ &= E[Q | T_1=t_1] \\ &= \int_Q f_{Q|T_1}(q|x) dQ\end{aligned}$$

? posterior dist

? so what is the posterior dist?

want info on q

$$\int_0^1 q^k dq$$

(35)

So this probability of T_1 is time to first
occurrence of Bernoulli (q)

? But q is random $[0,1]$

- why type of problem is this?
- # trials to 1st success = geometric

$$(1-p)^{k-1} p \quad k=1$$

$$\begin{matrix} (1-p)^0 & p \\ & p \end{matrix}$$

and $p = q$.

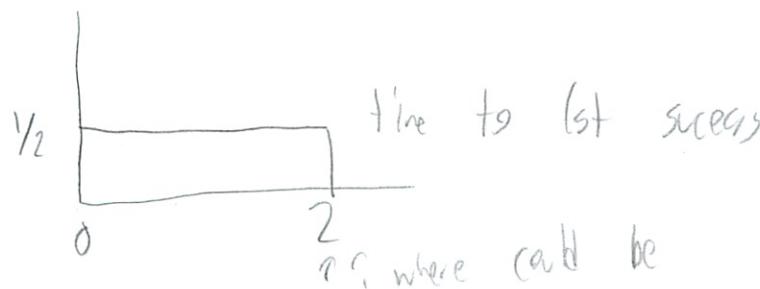
$$\frac{1}{b-a+1} = \frac{1}{1-0+1} = \left(\frac{1}{2}\right) \quad (-1)$$



Bernoulli



Geometric



(36)

i) Compute MAP estimate of Q given k recordings

$$t_1 = t_1, \dots, T_k = t_k$$

Now have k estimates

Conduct posterior dist, find max value in it

$$f_{Q|T}(q | t = t_1, \dots, t_k) = \frac{f_{T|Q}(t | q) f_Q(q)}{\sum f_{T|Q}(t | q) f_Q(q)}$$

$$f_{T|Q}(t | q)$$

? just the geometric

$$(1-p)^{k-1} \cdot p \quad k=1$$

$$p$$

$$p=q$$

So definitive value q

$$f_Q(q) = \text{Uniform} \quad \frac{1}{b-a+1} = \frac{1}{1-0+1} = \frac{1}{2}$$

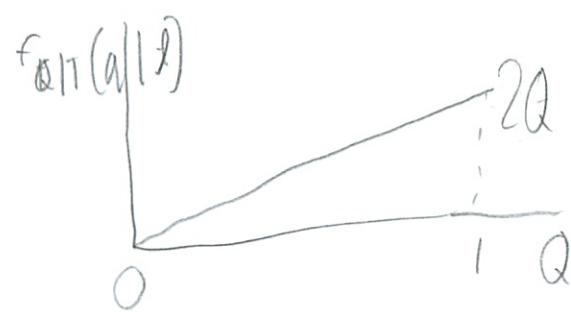
$$\int_0^1 \frac{1}{2} dq = \left. \frac{q^2}{4} \right|_0^1 = \frac{1}{4}$$

$$= \frac{\frac{1}{2}}{\frac{1}{4}} = 2q$$

if this right

(37)

Need a single value - the max



Max is $q=1$ w/ prob 2
can't be prob > 1

(-1)

(38)

$$f_{Q|T_1, \dots, T_n}(q | t_1, \dots, t_n)$$

$$P_Q(q) P_{T_1, T_2, \dots, T_n | Q}(t_1, \dots, t_n | q)$$

find q that maximizes this

* Given Q are ind. geometric

If not given Q , then not ind

Product in terms of $\underbrace{(1-q)}_{A(q)} - \underbrace{(q)}_{B(q)}$

Sum in the exponential

take deriv of ~~$\frac{d}{dq}$~~
w/ respect to q

$$\cancel{\frac{\partial A(q)}{\partial q}} + \cancel{\frac{\partial B(q)}{\partial q}}$$

$$\cancel{\frac{\partial A(q)}{\partial q}} \cancel{\frac{d}{dq} B(q)} + B(q) \cancel{\frac{d}{dq} A(q)}$$

Set = 0

- could be max, min

take 2nd deriv if min is 0

- well other way around

(39)

LSE/LMS

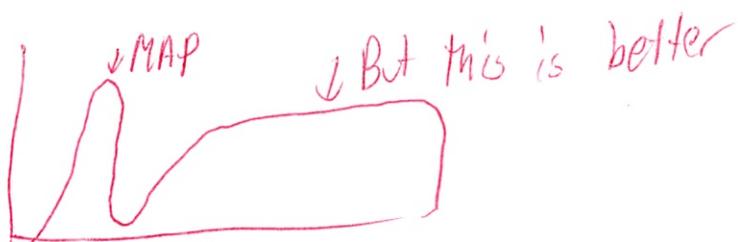
$$\hat{\theta}_{LMS}(x) = E[\theta | x=x]$$

Measured x
got this

Computing may not be easy

form of it is simple

When want to minimize the size of the error



Care about how big error is

- How far off are you

Not how often will you be wrong (MAP)

$$\hat{\theta}_{LMS} = aX + b$$

best you can do is

$$= E[\theta] + \frac{cov(x, \theta)}{var(x)} (x - E[x]) \quad \leftarrow \text{more intuitive}$$

$$= E[\theta] + \rho \frac{\sigma_\theta}{\sigma_x} \underset{\text{normalize}}{(x - E[x])}$$

- amt above / below $E[\theta]$

(40)

For this part only assume θ is sampled from RV Q
which is uniform $[1, 5]$

d) Find LLSE of 2nd interarrival time T_2

from Observed 1st arrival time $T_1 = t_1$

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

$$\hat{Q}_{LLSE} = E[Q] + \frac{\text{Cov}(Q, T)}{\text{Var}(T)} (T - E[T]) \quad +0.5$$

$$E[T] = \text{geometric} = \frac{1}{p} = \frac{1}{q}$$

$$\text{Var}(T) = \text{geometric} = \frac{1-p}{p^2} = \frac{1-q}{q^2}$$

$$E[Q] = \text{Uniform} = \frac{a+b}{2} = \frac{3}{4}$$

$$\text{Var}(Q) = \text{Uniform} = \frac{(b-a)(b-a+2)}{12} = \frac{(1-5)(1-5+2)}{12} = \frac{15 \cdot 2.5}{12} = \frac{5}{48}$$

$$\text{Cov}(Q, T) = E[QT] - E[Q] \cdot E[T]$$

(-0.5)

(-0.5)

$$? E[QT] = E\{E[QT | T=t]\}$$

$$? Q E[T | T=t]$$

$$Q \cdot \frac{1}{q}$$

$$Q \cdot \frac{1}{q} = 1$$

(41)

$$= E[1] = 1$$

$$\text{Cov}(Q, T) = 1 - \frac{3}{4} \cdot \frac{1}{q}$$

should it have a variable

$$\uparrow \frac{1}{3/4}$$

$= 0$ totally uncorrelated

(-0.5)

$$\text{LLMS} = \frac{3}{4} + \underline{0} \dots$$

(rest does not matter)

$$= \left(\frac{3}{4} \right)$$

- Where did we base it on T_1 ?

- T_1

- This is for any T_{1k}

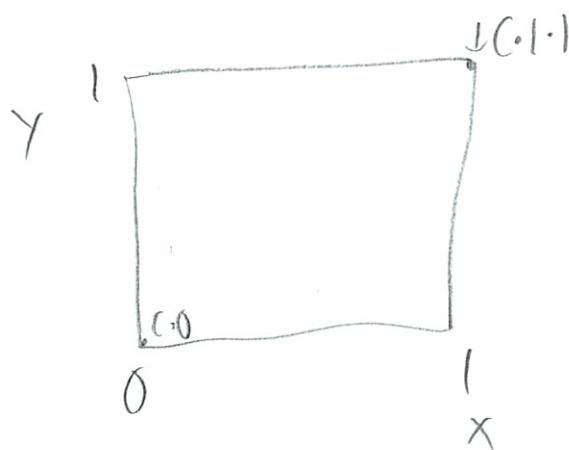
$0.5/5$

I have no clue if this is right

(42)

6. Joint PDF of X, Y defined

$$f_{X,Y}(x,y) = \begin{cases} Cxy & \text{if } 0 < x \leq 1, 0 < y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

a) Find normalization constant C 

$$\int_0^1 \int_0^1 Cxy \, dy \, dx$$

$$\left[\int_0^1 \frac{1}{2} Cx y^2 \, dx \right]_0^1$$

$$\int_0^1 \frac{Cx}{2} \, dx$$

$$\frac{Cx^2}{4} \Big|_0^1$$

$$\frac{C}{4} = \text{want} = 1$$

$$\frac{C}{4} = 1$$

$$C = 4$$

(43)

b) Compute the conditional expectation estimator of X
based on observed value $Y=y$

Conditional expectation estimator = LMS estimator

$$\begin{aligned}\hat{\theta}_{\text{LMS}} &= E[\theta | X=x] \quad \text{General variables} \\ &= E[X | Y=y] \quad \text{our variables} \\ &= \int x f_{X|Y}(x|y=y) dx \\ &\quad \text{Posterior def}\end{aligned}$$

$$f_{X|Y} = \frac{f_{X,Y}(x,y)}{f_Y(y)} \underbrace{4_{XY}}_{\text{observed so w/ certainty = 1}}$$

$$= \int x 4_{XY} dx$$

$$= \int_0^1 4_{XY} x^2 dx$$

$$\frac{4}{3} y^3 \Big|_0^1$$

$\frac{4}{3} y$ ← So not plug in observed y

(44)

Q Is this estimate different what you would have guessed before saw $y = \gamma$?

Yes, of course because you now have additional information which you did not have before

Before $4xy$

After $\frac{4}{3}y$

(45)

d) Repeat for MAP

↳ just max of posterior definition $\frac{q_{xy}}{I}$

What x will max q_{xy} ?

The largest possible $x \rightarrow x=1$

Yes this also provides slightly more information which we did not have before

Before q_{xy}

After q_y

It gives us an actual estimate, although it is a fairly inaccurate estimate. The new information is not very significant

Problem Set 10 Solutions

1. A financial parable.

- (a) The bank becomes insolvent if the asset's gain $R \leq -5$ (i.e., it loses more than 5%). This probability is the CDF of R evaluated at -5 . Since R is normally distributed, we can convert this CDF to be in terms of a standard normal random variable by subtracting away the mean and dividing by the standard deviation, and then look up the value in a standard normal CDF table.

$$\begin{aligned} E[R] &= 7, \\ \text{var}(R) &= 10^2 = 100, \\ P(R \leq -5) &= P\left(\frac{R-7}{10} \leq \frac{-5-7}{10}\right) = \Phi(-1.2) \approx 0.115. \end{aligned}$$

Thus, by investing in just this one asset, the bank has a 11.5% chance of becoming insolvent.

- (b) If we model the R_i 's as **independent** normal random variables, then their sum $R = (R_1 + \dots + R_{20})/20$ is also a normal random variable (see Example 4.11 on page 214 of the text). Thus, we can calculate the mean and variance of this new R and proceed as in part (a). Note that since the random variables are assumed to be independent, the variance of their sum is just the sum of their individual variances.

$$\begin{aligned} E[R] &= (E[R_1] + \dots + E[R_{20}])/20 = 7, \\ \text{var}(R) &= \frac{1}{20^2}(\text{var}(R_1) + \dots + \text{var}(R_{20})) = \frac{20 \cdot 100}{400} = 5, \\ P(R \leq -5) &= P\left(\frac{R-7}{\sqrt{5}} \leq \frac{-5-7}{\sqrt{5}}\right) = \Phi(-5.367) \approx 0.0000000439 = 4.39 \cdot 10^{-8}. \end{aligned}$$

Thus, by diversifying and assuming that the 20 assets have **independent** gains, the bank has seemingly decreased its probability of becoming insolvent to a palatable value.

- (c) Now, if the gains R_i are positively correlated, then we can no longer sum up the individual variances; we need to account for the covariance between pairs of random variables. The covariance is given by

$$\text{cov}(R_i, R_j) = \rho(R_i, R_j) \sqrt{\text{var}(R_i)\text{var}(R_j)} = \frac{1}{2} \sqrt{10^2 \cdot 10^2} = 50.$$

From page 220 in the text, we know that the variance in this case is

$$\begin{aligned} \text{var}(R) &= \text{var}\left(\frac{1}{20} \sum_{i=1}^{20} R_i\right) = \frac{1}{400} \left(\sum_{i=1}^{20} \text{var}(R_i) + \sum_{\{(i,j)|i \neq j\}} \text{cov}(R_i, R_j) \right) \\ &= \frac{1}{400} (20 \cdot 100 + 380 \cdot 50) = 52.5. \end{aligned}$$

Since we assume that $R = (R_1 + \dots + R_{20})/20$ is still normal, we can again apply the same steps as in parts (a) and (b):

$$\begin{aligned} E[R] &= (E[R_1] + \dots + E[R_{20}])/20 = 7, \\ \text{var}(R) &= 52.5, \\ P(R \leq -5) &= P\left(\frac{R-7}{\sqrt{52.5}} \leq \frac{-5-7}{\sqrt{52.5}}\right) = \Phi(-1.656) \approx 0.0488. \end{aligned}$$

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
 Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
 (Fall 2010)

Thus, by taking into account the positive correlation between the assets' gains, we are no longer as comfortable with the probability of insolvency as we thought we were in part (b).

2. Let M and N be the number of males and females, respectively, that cast a vote. We need to find $P(M > N)$, i.e., $P(M - N > 0)$. The central limit theorem does not apply directly to the random variable $M - N$. However, the central limit theorem implies that M and N are well approximated by normal random variables. So, $M - N$ is the difference of two independent approximately normal random variables. Since the difference of two normal random variables is itself normal, it follows that $M - N$ is approximately normal. The mean and variance of $M - N$ are found by

$$\begin{aligned} \mathbf{E}[M - N] &= 300 \cdot 0.4 + 196 \cdot 0.5 = 120 - 98 = 22, \\ \text{var}(M - N) &= \text{var}(M) + \text{var}(N) = 300 \cdot 0.4 \cdot 0.6 + 196 \cdot 0.5 \cdot 0.5 = 121. \end{aligned}$$

Thus, the standard deviation of $M - N$ is 11. Let Z be a standard normal random variable. Using the central limit theorem approximation, we obtain

$$\begin{aligned} \mathbf{P}(M - N > 0) &= \mathbf{P}\left(\frac{M - N - 22}{11} > -\frac{22}{11}\right) \\ &\approx \mathbf{P}(Z \geq -2) \\ &= 0.9772. \end{aligned}$$

A slightly more refined estimate is obtained by expressing the event of interest as $\mathbf{P}(M - N \geq 1/2)$. We then have

$$\begin{aligned} \mathbf{P}(M - N > 1/2) &= \mathbf{P}\left(\frac{M - N - 22}{11} \geq -\frac{21.5}{11}\right) \\ &\approx \mathbf{P}(Z \geq -1.95) \\ &= 0.974. \end{aligned}$$

3. (a) Using the Central Limit Theorem, we obtain $\mathbf{P}(\frac{n}{2} - 10 \leq S_n \leq \frac{n}{2} + 10) \approx \Phi(\frac{20}{\sqrt{n}}) - \Phi(-\frac{20}{\sqrt{n}}) \rightarrow 0$ as $n \rightarrow \infty$.
 - (b) The limit is 1, by the weak law of large numbers.
 - (c) Using the Central Limit Theorem, we obtain $\mathbf{P}(\frac{n}{2} - \frac{\sqrt{n}}{2} \leq S_n \leq \frac{n}{2} + \frac{\sqrt{n}}{2}) \rightarrow \Phi(1) - \Phi(-1) = 0.6826$.
 4. (a) Let C denote the coin that Bob received, so that $C = 1$ if Bob received the first coin, and $C = 2$ if Bob received the second coin. Then $\mathbf{P}(C = 1) = p$ and $\mathbf{P}(C = 2) = 1 - p$. Given C , the number of heads Y in 3 independent tosses is a binomial random variable.
- We can find the probability that Bob received the first coin given that he observed k heads using Bayes' rule.

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
 Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
 (Fall 2010)

$$\begin{aligned}
 P(C = 1 | Y = k) &= \frac{P(Y = k | C = 1) \cdot P(C = 1)}{P(Y = k | C = 1) \cdot P(C = 1) + P(Y = k | C = 2) \cdot P(C = 2)} \\
 &= \frac{\binom{3}{k} \cdot (1/3)^k (2/3)^{3-k} p}{\binom{3}{k} \cdot (1/3)^k (2/3)^{3-k} p + \binom{3}{k} \cdot (2/3)^k (1/3)^{3-k} \cdot (1-p)} \\
 &= \frac{2^{3-k} p}{2^{3-k} p + 2^k (1-p)} = \frac{1}{1 + \frac{1-p}{p} 2^{2k-3}}
 \end{aligned}$$

(b) We want to find k so that the following inequality holds.

$$\begin{aligned}
 P(C = 1 | Y = k) &> p \\
 \frac{2^{3-k} p}{2^{3-k} p + 2^k (1-p)} &> p
 \end{aligned}$$

Note that if $p = 0$ or $p = 1$, there is no value of k that satisfies the inequality. We now solve it for $0 < p < 1$:

$$\begin{aligned}
 \frac{2^{3-k}}{2^{3-k} p + 2^k (1-p)} &> 1 \\
 2^{3-k} &> 2^{3-k} p + 2^k (1-p) \\
 2^{3-k} (1-p) &> 2^k (1-p) \\
 2^{3-k} &> 2^k \\
 2k &< 3 \\
 k &< 3/2
 \end{aligned}$$

For $0 < p < 1$, $k = 0$ or $k = 1$ the probability that Alice sent the first coin increases. The inequality does not depend on p , and so does not change when p increases. Intuitively, this makes sense: lower values of k increase Bob's belief he got the coin with lower probability of heads.

(c) Given that Bob observes k heads, Bob must decide on whether the first or second coin was used. To minimize the error, he should decide it is the first coin when $P(C = 1 | Y = k) \geq P(C = 2 | Y = k)$. Thus, we have the decision rule given by

$$\begin{aligned}
 P(C = 1 | Y = k) &\geq P(C = 2 | Y = k) \\
 \frac{2^{3-k} p}{2^{3-k} p + 2^k (1-p)} &\geq \frac{2^k (1-p)}{2^{3-k} p + 2^k (1-p)} \\
 2^{3-k} p &\geq 2^k (1-p) \\
 2^{2k-3} &\leq \frac{p}{1-p} \\
 k &\leq \frac{3}{2} + \frac{1}{2} \log_2 \frac{p}{1-p}
 \end{aligned}$$

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
 Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
 (Fall 2010)

- (d) i. If $p = 2/3$, the threshold in the rule above is equal to $\frac{3+\log_2 2}{2} = 2$. Therefore, Bob will decide that he received the first coin when he observes 0, 1 or 2 heads, and will decide that he received the second coin when he observes 3 heads.

We find the probability of a correct decision using the total probability law:

$$\begin{aligned}\mathbf{P}(\text{Correct}) &= \mathbf{P}(\text{Correct} | C = 1) \cdot p + \mathbf{P}(\text{Correct} | C = 2) \cdot (1 - p) \\ &= \mathbf{P}(Y < 3 | C = 1) \cdot p + \mathbf{P}(Y = 3 | C = 2) \cdot (1 - p) \\ &= (1 - \mathbf{P}(Y = 3 | C = 1)) \cdot p + \mathbf{P}(Y = 3 | C = 2) \cdot (1 - p) \\ &= (1 - (1/3)^3)(2/3) + (2/3)^3(1/3) = 20/27 \approx .741\end{aligned}$$

- ii. In absence of any data, all Bob can do is decide he received the first coin with some probability q . Note that this rule includes the deterministic decisions that he received either the first coin ($q = 1$) or the second coin ($q = 0$).

In this case, the probability of correct decision is equal to

$$\begin{aligned}\mathbf{P}(\text{Correct}) &= \mathbf{P}(\text{Correct} | C = 1) \cdot p + \mathbf{P}(\text{Correct} | C = 2) \cdot (1 - p) \\ &= qp + (1 - q)(1 - p) = 1 - p + q(2p - 1) = \frac{1+q}{3}\end{aligned}$$

Clearly, the probability of the correct decision is maximized (or the probability of error is minimized) when $q = 1$, i.e., when Bob deterministically decides he received the first coin. In this case, $\mathbf{P}(\text{Correct}) = 2/3 \approx .667$. Observing 3 coin tosses increases the probability of the correct decision by $2/27 \approx .074$.

- (e) If p is increased, the threshold in the decision rule in part (c) goes up, i.e., the range of values of k for which Bob decides he received the first coin can only go up.
 (f) Bob will never decide he received the first coin if the threshold in the rule above is below zero:

$$\begin{aligned}\frac{3}{2} + \frac{1}{2} \log_2 \frac{p}{1-p} &< 0 \\ \log_2 \frac{p}{1-p} &< -3 \\ \frac{p}{1-p} &< \frac{1}{8} \\ p &< \frac{1}{9}\end{aligned}$$

If $p < 1/9$, the prior probability of receiving the first coin is so low that no amount of evidence from 3 tosses of the coin will make Bob decide he received the first coin.

- (g) Bob will always decide he received the first coin if the threshold in the rule above is equal to or above 3:

$$\begin{aligned}\frac{3}{2} + \frac{1}{2} \log_2 \frac{p}{1-p} &\geq 3 \\ \log_2 \frac{p}{1-p} &\geq 3 \\ \frac{p}{1-p} &\geq 8 \\ p &\geq \frac{8}{9}\end{aligned}$$

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
 Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
 (Fall 2010)

If $p \geq 8/9$, the prior probability of receiving the first coin is so high that no amount of evidence from 3 tosses of the coin will make Bob decide he received the second coin.

5. (a) Using the total probability theorem, we have

$$p_{T_1}(t) = \int_0^1 p_{T_1|Q}(t, q) f_Q(q) dq = \int_0^1 (1-q)^{t-1} q dq = \frac{1}{(t+1)t} \quad \text{for } t = 1, 2, \dots$$

- (b) The least squares estimate coincides with the conditional expectation of Q given T_1 , which is derived as

$$\begin{aligned} \mathbb{E}[Q | T_1 = t] &= \int_0^1 p_{Q|T_1}(q | t) q dq \\ &= \int_0^1 \frac{p_{T_1|Q}(t | q) f_Q(q)}{p_{T_1}(t)} q dq \\ &= \int_0^1 t(t+1)q(1-q)^{t-1} q dq \\ &= \int_0^1 t(t+1)q^2(1-q)^{t-1} dq \\ &= t(t+1) \frac{2(t-1)!}{(t+2)!} \\ &= \frac{2}{t+2} \end{aligned}$$

- (c) We write the posterior probability distribution of Q given $T_1 = t_1, \dots, T_k = t_k$

$$\begin{aligned} f_{Q|T_1, \dots, T_k}(q | t_1, \dots, t_k) &= \frac{f_Q(q) \prod_i^k P_{T_i}(T_i = t_i | Q = q)}{\int_0^1 f_Q(q) \prod_i^k P_{T_i}(T_i = t_i | Q = q) dq} \\ &= \frac{q^k (1-q)^{\sum_i^k t_i - k}}{c} \\ &= \frac{1}{c} q^k (1-q)^{\sum_i^k t_i - k}, \end{aligned}$$

where the denominator integrates out q so it could be viewed as a constant scalar c .

To maximize the above probability we set its derivative with respect to q to zero

$$kq^{k-1}(1-q)^{\sum_i^k t_i - k} - (\sum_i^k t_i - k)q^k(1-q)^{\sum_i^k t_i - k - 1} = 0,$$

or equivalently

$$k(1-q) - (\sum_i^k t_i - k)q = 0,$$

which yields the MAP estimate

$$\hat{q} = \frac{k}{\sum_{i=1}^k t_i}.$$

For this part only assume q is sampled from the random variable Q which is now uniformly distributed over $[0.5, 1]$

(d) The LLSE of T_1 given T_2 is

$$\hat{T}_2 = \mathbf{E}[T_2] + \frac{\text{cov}(T_1, T_2)}{\text{var}(T_1)}(T_1 - \mathbf{E}[T_1]),$$

where the coefficients are

$$\mathbf{E}[T_1] = \mathbf{E}[T_2] = \int_{0.5}^1 f_Q(q) \mathbf{E}[T|Q=q] dq = \int_{0.5}^1 2 * 1/q dq = 2 \ln 2,$$

and from the law of total variance

$$\begin{aligned} \text{var}(T_1) &= \text{var}(T_2) = \mathbf{E}[\text{var}(T_1 | Q)] + \text{var}[\mathbf{E}(T_1 | Q)] \\ &= \mathbf{E}\left[\frac{1-Q}{Q^2}\right] + \text{var}\left[\frac{1}{Q}\right] \\ &= \mathbf{E}[1/Q^2] - \mathbf{E}[1/Q] + \mathbf{E}[1/Q^2] - \mathbf{E}[1/Q]^2 \\ &= \int_{0.5}^2 f_Q(q) \frac{1}{q^2} dq - \int_{0.5}^2 f_Q(q) \frac{1}{q} dq + \int_{0.5}^2 f_Q(q) \frac{1}{q^2} dq - \left(\int_{0.5}^2 f_Q(q) \frac{1}{q} dq\right)^2 \\ &= 2 - 2 \ln 2 + 2 - (2 \ln 2)^2 \\ &= 4 - 2 \ln 2 - (2 \ln 2)^2, \end{aligned}$$

and their covariance

$$\begin{aligned} \text{cov}(T_1, T_2) &= \mathbf{E}[T_1 T_2] - \mathbf{E}[T_1] \mathbf{E}[T_2] \\ &= \mathbf{E}[\mathbf{E}[T_1 T_2 | Q]] - \mathbf{E}[T_1] \mathbf{E}[T_2] \\ &= \mathbf{E}[\mathbf{E}[T_1 | Q] \mathbf{E}[T_2 | Q]] - \mathbf{E}[T_1] \mathbf{E}[T_2] \\ &= \mathbf{E}[1/Q^2] - \mathbf{E}[T_1] \mathbf{E}[T_2] \\ &= 2 - 4(\ln 2)^2 \end{aligned}$$

Therefore we have derived the linear least squares estimator

$$\hat{T}_2 = 2 \ln 2 + \frac{2 - 4(\ln 2)^2}{4 - 2 \ln 2 - (2 \ln 2)^2} (T_1 - 2 \ln 2) \approx 1.543 + 0.113 T_1.$$

6. (a) To find the normalization constant c we integrate the joint PDF:

$$\int_0^1 \int_0^1 f_{X,Y}(x, y) dy dx = c \int_0^1 \int_0^1 xy dy dx = c \int_0^1 1/2x dx = c/4.$$

Therefore, $c = 4$.

(b) To construct the conditional expectation estimator, we need to find the conditional probability density.

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} = \frac{4xy}{\int_0^1 4xy dx} = \frac{4xy}{2y} = 2x, \quad x \in (0, 1]$$

Thus

$$\hat{x}_{\text{CE}}(y) = \mathbf{E}[X | Y = y] = \int_0^1 x \cdot 2x dx = 2/3.$$

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
Department of Electrical Engineering & Computer Science
6.041/6.431: Probabilistic Systems Analysis
(Fall 2010)

- (c) We first note that the conditional probability does not depend on y . Therefore, X and Y are independent, and whether or not we observe $Y = y$ does not affect the estimate in part (b). Another way to see this is to consider that if we do not observe y , we can compute the marginal $f_X(x) = \int_0^1 4xy dy = 2x$ which is equal to the conditional density, and will therefore produce the same estimate.
- (d) Since X and Y are independent, no estimator can make use of the observed value of Y to estimate X . The MAP estimator for X is equal to 1, regardless of what value y we observe, since the conditional (and the marginal) density is maximized at 1.

Recitation 23
December 2, 2010

1. Example 9.1, page 463 in textbook

Romeo and Juliet start dating, but Juliet will be late on any date by a random amount X , uniformly distributed over the interval $[0, \theta]$. The parameter θ is unknown. Assuming that Juliet was late by an amount x on their first date, find the ML estimate of θ based on the observation $X = x$.

2. Example 9.4, page 464 in textbook

Estimate the mean μ and variance v of a normal distribution using n independent observations X_1, \dots, X_n .

3. Example 9.8, page 474 of textbook

We would like to estimate the fraction of voters supporting a particular candidate for office. We collect n independent sample voter responses X_1, \dots, X_n , where X_i is viewed as a Bernoulli random variable, with $X_i = 1$ if the i th voter supports the candidate. We conducted a poll of 1200 people in North Carolina, and found that 684 were supporting the candidate. We would like to construct a 95% confidence interval for θ , the proportion of people who support the candidate. As we saw in lecture, using the central limit theorem, an (approximate) 95% confidence interval can be defined as

$$\hat{\Theta}^- = \hat{\Theta}_n - 1.96\sqrt{\frac{v}{n}}, \quad \hat{\Theta}^+ = \hat{\Theta}_n + 1.96\sqrt{\frac{v}{n}}$$

where $v = \text{Var}(X_i)$, and $\hat{\Theta}_n = (X_1 + \dots + X_n)/n$. Unfortunately, we don't know the value for v . Construct confidence intervals for θ using the following three ways of estimating or bounding the value for v (in each case simply assume that v is equal to the given estimate; note that this is a further approximation in cases (a) and (b)).

(a)

$$\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\Theta}_n)^2$$

(b)

$$\hat{\Theta}_n(1 - \hat{\Theta}_n)$$

(c) The most conservative upper bound for the variance.

Recitation 23

- More stats
 - highlights of stats
 - infer obs based on unknowns
 - Bayesian - estimate θ given X
 - $\hat{\theta}$ unknown param
 - x_{obs} obs
 - ~~Posterior~~ given $f_{\theta}(\theta)$ $f_{x|\theta}(x|\theta)$
 - prior
 - conditional
 - Observation
 - goal: ~~get~~ get posterior w/ Bayes rule

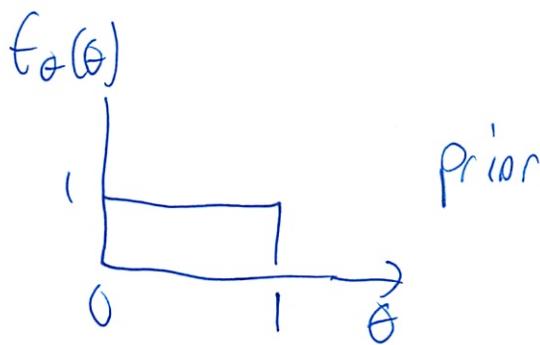
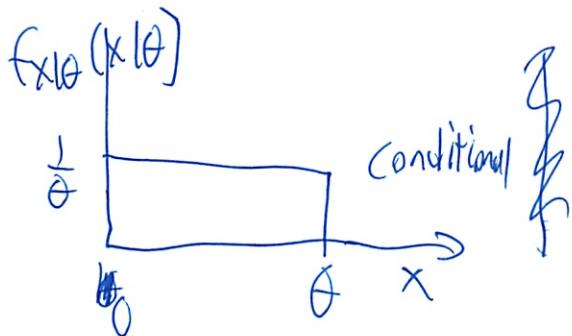
$$f_{\theta|x}(\theta|x) = \frac{f_{\theta}(\theta) f_{x|\theta}(x|\theta)}{f_x(x)}$$
 - lots of things to look at
 - $E[\theta|x]$
 - LSE
 - MAP
- 
- $\hat{\theta}_{\text{MAP}}$ does not matter
 - $\arg \max_{\theta} f_{\theta}(\theta) f_{x|\theta}(x|\theta)$
 - can solve in closed form
 - w/o calculation

(2)

1. Romeo + Juliet are back again

Juliet late by amt x

↳ uniform $[0, \theta]$



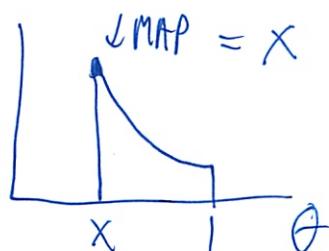
MAP rule

- maximize product of the two

$$\max_{\theta} \frac{1}{\theta} \quad \text{for}$$

- know $\theta > x$
 $\theta < 1$

So



- not so good
- check br calc var

③

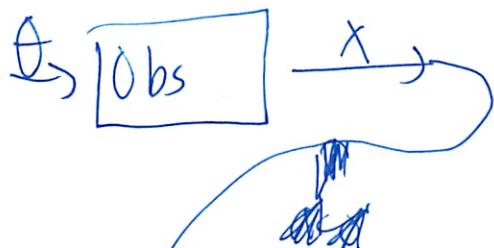
Classical Estimation

- given many prob models - one for each possible θ
- θ is unknown, not "random"

$f_x(x|\theta)$ ← one PDF for each value of θ

Observe x , try to find which θ it came from

One way \rightarrow Max Likelihood method (ML)



Consider all possible θ

for each θ , is some PDF for x
pick the one that is most likely

$$f_x(x|\theta_1)$$
$$f_x(x|\theta_2)$$
$$\dots$$
$$f_x(x|\theta_m)$$

max
 θ
in each
 θ

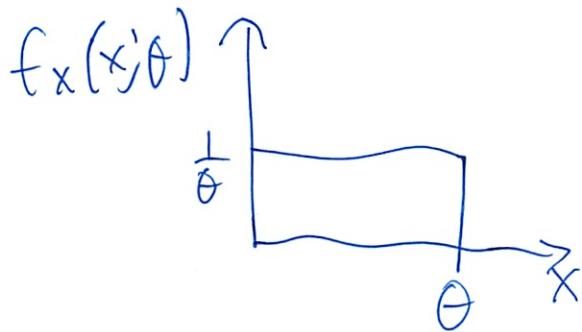
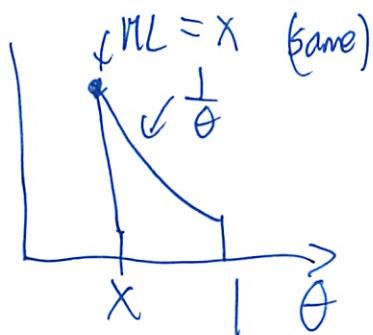
→ Max Likelihood $\theta(x)$

Then for each $x = x$

$$f_x(x|\theta)$$

find each value individually

(4)

R+J againGive likely hood for x for each value of θ θ is ~~any~~ unknown except $0 \leq \theta \leq 1$ max over $x \leq \theta \leq 1$ Plot likelihood function for each θ } ^{↑ from problem definition}Is same
- since it's a flat prior

- still 2 procedures are conceptually very different

(5)

2. multiple samples from normal - classical - ML

$X = X_1, \dots, X_n = \text{iid Normal RV}$

$$\mathbb{E}[X_i] = \mu$$

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

? we don't know σ^2
so need to estimate it

$$\theta = (\mu, \sigma^2)$$

? want, via likelihood function

$$\begin{aligned} \text{likelihood } f_n &\rightarrow f_n(x_1, \dots, x_n | \mu, \sigma^2) = \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-(x_i - \mu)^2 / 2\sigma^2} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot e^{-\sum_{i=1}^n (x_i - \mu)^2 / 2\sigma^2} \\ &\quad \text{fancy algebra} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \frac{1}{\sqrt{n}} \cdot e^{-n(S_n^2 + (\bar{x}_n - \mu)^2) / 2\sigma^2} \end{aligned}$$

where $\bar{x}_n = \text{sample mean} = \frac{\sum_{i=1}^n x_i}{n}$

$S_n^2 = \text{avg of sum of deviations} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$

- opt like sample var - but w/ sample mean

(6)

Can max log instead

- easier to max sums, instead of products

the "log likelihood"

$$\text{log } f_X(x_1, \dots, x_n | \mu, \nu) =$$

$$= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\nu) - \frac{n}{2\nu} S_n^2 - \frac{n(m_n - \mu)^2}{2\nu}$$

Max w/ respect to μ, ν

can do deriv

but shortcut μ only in last term
and is positive

So only set last part = 0

$$\boxed{\hat{\mu} = m_n}$$

So max w/ respect to ν of rest

- take deriv, = 0

$$\frac{\partial}{\partial \nu} = 0 \rightarrow -\frac{n}{2\nu} + \frac{n}{2} \frac{S_n^2}{\nu^2} = 0$$

~~$$\frac{S_n^2}{\nu} = 0$$~~

$$\boxed{\hat{\nu} = S_n^2}$$

⑦ What we did

To estimate mean + var of normal, given
n ind samples of it

We don't know if estimate is good

Check

$$E[\hat{\mu}] = \mu$$

Same as true mean
called "unbiased"

$$E[S_n^2] \neq \sigma^2$$

- but as $n \rightarrow \infty$

$$E[S_n^2] \rightarrow \sigma^2$$

called "asymptotical unbiasedness"

? ML has this in general

$\hat{\mu}_n \rightarrow \mu$ w/ prob 1

- special case WLLN

"Consistency"

↳ becomes exact

$\hat{V}_n \rightarrow V$ close to WLLN

also "consistency"

⑧

ML is one of the major methods for estimation

We had gotten estimates \bar{X}_n , S_n^2 - but how
confident are we?

- We introduce "confidence interval"

- again have multiple obs

- but don't make point estimate $\hat{\theta}_n$

- make lower + upper bounds of θ

↑ ("interval estimate")
 $[\theta_n^- \text{, } \theta_n^+]$

- For every value of θ , we know θ is within
the interval w/ high probability

$$\theta_n^- \leq \theta \leq \theta_n^+ \text{ w/ high prob}$$

↓ have as many intervals as have θ in classical
Only 1 interval in Bayesian

\sqcup \sqcup
RV - dist based on θ
- estimator

$$P(\theta_n^- \leq \theta \leq \theta_n^+) \text{ is a FF}$$

⑨

Say called $1-\alpha$ confidence interval

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

~~Want estimate~~

Things to remember

1. $\bar{\theta}_n, \theta_n^+$ are RVs (dist depends on θ)

2. Usually $\bar{\theta}_n = \hat{\theta}_n - \beta$

$$\theta_n^+ = \hat{\theta}_n + \beta \quad \text{some reasonable estimator}$$

Get end points w/ estimator

then go ~~up~~ left and right

adjust β till inequality holds

3. Often ~~use~~ approx - ^{using} the CLT

$\bar{\theta}_n, \theta_n^+$ calculated

⑩.2 again

$$X = X_1, \dots, X_n$$

mean μ ~~known~~ still unknown

var V known

⑩

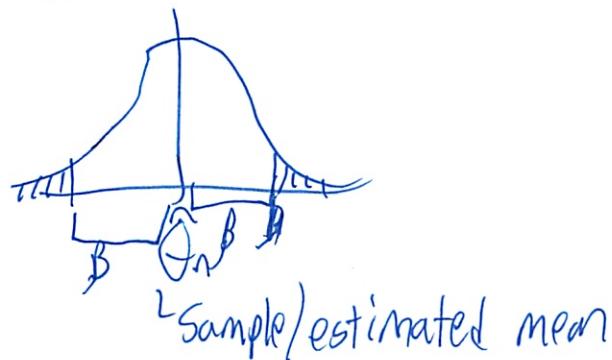
Adjust until threshold is met

- we know Var - so have idea of dist

$$\hat{\theta}_n = \frac{X_1 + \dots + X_n}{n} \quad \text{sample mean}$$

Try CI $\{\hat{\theta}_n - \beta, \hat{\theta}_n + \beta\}$

Find β



$$P(|\hat{\theta}_n - \theta| \leq \beta) \geq 1 - \alpha$$

$$P(|\hat{\theta}_n - \theta| \leq \beta) \geq .95$$

So $\beta \geq 1.96 \sqrt{\frac{V}{n}}$ ^{table}

find smallest β that satisfies

$$\left[\hat{\theta}_n - 1.96 \sqrt{\frac{V}{n}}, \hat{\theta}_n + 1.96 \sqrt{\frac{V}{n}} \right]$$

⑪

But what if X_i not normal, but are iid

- Or if var is unknown

- So look at large sample \rightarrow where sample mean \sim normal

- reasonably good approx

- but also need var - Use nature of samples to get
approx var

- 1. Pretend sample mean is normal
- 2. Using approx var
estimated

3. Pollster example

$X = X_1, \dots, X_n$ Bernoulli

$$E[X_i] = \theta$$

Construct CI $[\hat{\theta}_n - \beta, \hat{\theta}_n + \beta]$

? Sample mean = $\hat{\theta}_n$

- fraction of 1s in sequence

Again (CLT) $B = 1.96 \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$ estimate of var
- of Bernoulli
1. $\hat{\theta}(1-\hat{\theta})$
So $\hat{\theta}(1-\hat{\theta})$

(12)

- or
2. Sample var
- consistent estimator for normal or large # of samples
 - good estimator

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

or

3. Use upper bound of var
 t -or Barnalli = $\hat{t}_Y = \hat{V}_n$
- But will over estimate
 - conservative CI

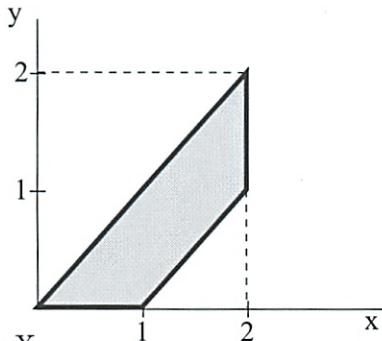
- Same type of calc in lecture 19
 except here estimate var

v) large # samples, not much different (> 100 samples)

Tutorial 11

1. Continuous random variables X and Y have a joint PDF given by

$$f_{X,Y}(x,y) = \begin{cases} 2/3 & \text{if } (x,y) \text{ belongs to the closed shaded region} \\ 0 & \text{otherwise} \end{cases}$$



We want to estimate Y based on X .

- (a) Find the LMS estimator $g(X)$ of Y .
 - (b) Calculate the conditional mean squared error $\mathbf{E}[(Y - g(X))^2 | X = x]$.
 - (c) Calculate the mean squared error $\mathbf{E}[(Y - g(X))^2]$. Is it the same as $\mathbf{E}[\text{var}(Y|X)]$?
 - (d) Derive $L(X)$, the linear LMS estimator of Y based on X .
 - (e) How do you expect the mean squared error of $L(X)$ to compare to that of $g(X)$?
 - (f) What problem do you expect to encounter, if any, if you try to find the MAP estimator for Y based on observations of X .
2. Consider a noisy channel over which you send messages consisting of 0s and 1s to your friend. It is known that the channel independently flips each bit sent with some fixed probability p ; however the value of p is unknown. You decide to conduct some experiments to estimate p and seek your friend's help. Your friend, cheeky as she is, insists that you send her messages consisting of three bits each (which you will both agree upon in advance); for each message, she will only tell you the total number of bits in that message that were flipped. Let X denote the number of bits flipped in a particular three-bit message.
- (a) Find the PMF of X .
 - (b) Derive the ML estimator for p based on X_1, \dots, X_n , the numbers of bits flipped in the first n three-bit messages.
 - (c) Is the ML estimator unbiased?
 - (d) Is the ML estimator consistent?
 - (e) You send $n = 100$ three-bit messages and find that the total number of bits flipped is 20. Construct a 95% confidence interval for p . If necessary, you may use a conservative bound on the variance of the number of bits flipped.
 - (f) What are some other ways to estimate the variance. How do you expect your confidence interval to change with different estimates of the variance.

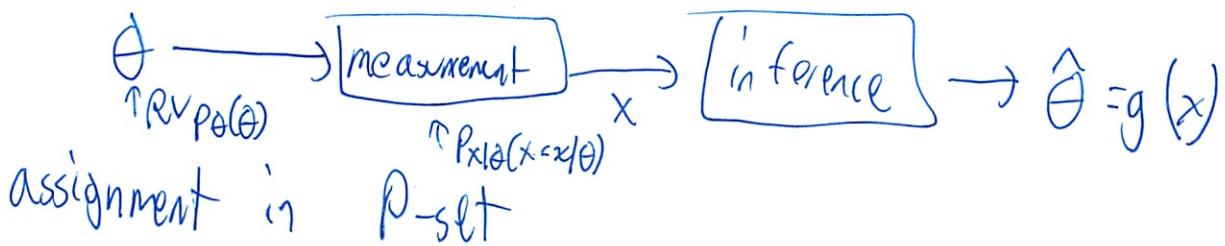
Tutorial 11

Bayesian Inference

have $\theta \rightarrow$ unknown, can't observe directly

$x = \text{observations}$

Want to make inference about θ



assignment in P -set

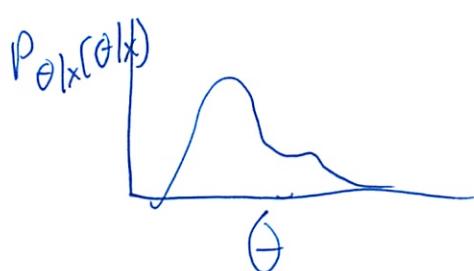
first step: get posterior prob

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

Bayes Rule

$$\stackrel{?}{=} \frac{P_{x|\theta} P_\theta(\theta)}{\sum_\theta P_{x|\theta} P_\theta(\theta)}$$

Now have



- a distribution

- want a single estimate

Error = result - actual



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



minimum probability
of error

(2)

Find by taking deriv, set = to 0, solve for θ

LMS/CEE

- minimizes mean square error

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

$$\hat{\theta}_{\text{LMS}} = E[\theta | X]$$

← often hard to find, complicated

LLMS

restrict estimator to linear to simplify

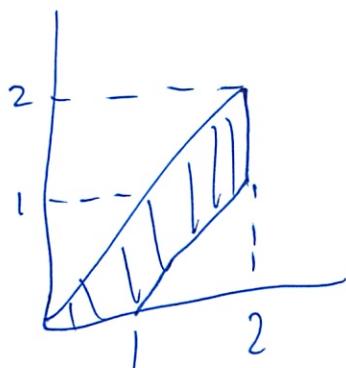
$$\hat{\theta} = E[\theta] + \frac{c(x, \theta)}{\text{Var}(\theta)} (x - E[x])$$

will have higher mean square error

\uparrow
or =

when posterior PDF is linear

#1 Given this joint PDF



Uniform in shaded region

(3)

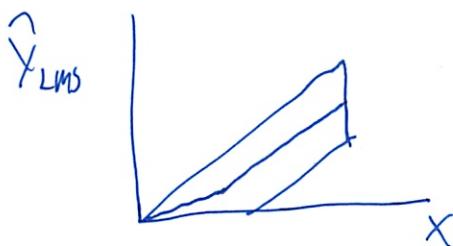
a) $\hat{Y}_{LMS} = E[Y|X]$

$Y|X$ is 

if uniform, $E[Y]$ is midpoint of line

$$= \begin{cases} \frac{x}{2} & \text{if } x \in [0, 1] \\ \frac{x+1}{2} = x - \frac{1}{2} & \text{if } x \in [1, 2] \end{cases}$$

Can plot estimator for each value of x



b) Now find the error that results from using this estimator

$$E[(Y - g(x))^2]_{X=x} = E[(Y - E[Y|X])^2]_{X=x} =$$

- is the conditional var in general

$$= \text{Var}(Y|X=x)$$

$$= \begin{cases} \frac{x^2}{12} & \text{if } x \in [0, 1] \\ \frac{1}{12} & \text{if } x \in [1, 2] \end{cases}$$

(4)

c) Now just mean squared error, no conditioning

$$E\{(Y - g(x))^2\} \stackrel{?}{=} E\left[\text{Var}(Y|X)\right]$$

? have

Use iterate expectation

Statement is true

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

Substitute from B

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

? in general this one

above was if have a specific $X=x$

$$= \int_0^2 \text{var}(Y | X=x) f_X(x) dy$$

need to find
by integrating over the ys
bounds are tricky

$$f_X(x) = \int_0^2 f_{X,Y}(x,y) dy$$

must break bounds down

$$= \begin{cases} \int_0^x \frac{2}{3} dy & x \in [0,1] \\ \int_{x+}^2 \frac{2}{3} dy & x \in [1,2] \end{cases}$$

$$(5) \quad = \begin{cases} \frac{2}{3}x & x \in [0, 1] \\ \frac{2}{3} & x \in [2, 3] \end{cases}$$

$$= \int_0^1 \frac{x^2}{12} \cdot \frac{2x}{3} dx + \int_1^2 \frac{1}{12} \cdot \frac{2}{3} dx \\ = \frac{5}{72}$$

d) LLMS

- just plug in formula

$$\hat{y} = L(x) = E[y] + \frac{\text{cov}(xy)}{\text{var}(y)} (x - E[x])$$

- p-set 4,5 style integrations and such
- skip

e) Compare LLMS to LMS

- less accurate
- LLMS error \leq LMS error
 ↑ part c
 can calculate
 w/ same process

as c

- also in book

$$\sigma_y^2(1-p^2)$$

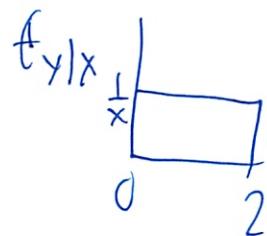
(6)

f) MAP for y based on X

$$\hat{y}_{MAP} = \underset{Y}{\operatorname{argmax}} f_{Y|X}(y|x)$$

? will be uniform for a particular X

But how find the max of something uniform?



Pick arbitrarily

No help

Classical Estimation

Still



? not an RV
no prior

$p(x; \theta)$
 $p_x(x)$ for each θ

Can't calc posterior

Can't do Bayes

So Max Likelihood $\hat{\theta}_n = \operatorname{argmax} P(x_1, x_2, \dots, x_n | \theta)$

7

Same ans as MAP in some cases

$$\text{MAP} \quad \underset{\theta}{\operatorname{argmax}} \quad P(\theta | x_1, \dots, x_n)$$

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{P(x_1, x_2, \dots, x_n | \theta)}_{\text{does not matter}} P_{\theta}(\theta)$$

\rightarrow When prior is uniform ~~MAP~~ $\hat{\theta}_{ML} = \hat{\theta}_{MAP}$

2. ~~Defn~~ PS 10 was hyp testing; Bayesian Inf

This is classical; and find p - not just test hyp

a) PMF of X

\sim binomial (p)

$$P_X(k; p) = \begin{cases} \binom{3}{k} p^k (1-p)^{3-k} & k = 0, 1, 2, 3 \\ 0 & \text{else} \end{cases}$$

b) $X_1, \dots, X_n = \# \text{ of bits flipped in 1st } n \text{ messages}$

$$\hat{p}_{ML} = \underset{p}{\operatorname{argmax}} \quad P(X_1, X_2, \dots, X_n | p)$$

- messages are ind

(8)

So can replace joint with a product when iid

$$= \arg \max_p \prod_{i=1}^n P(X_i; p)$$

$$= \underset{\text{? same as part a}}{\arg \max_p} \prod_{i=1}^n \binom{3}{k_i} p^{k_i} (1-p)^{3-k_i}$$

Take deriv, set = 0

- really hard to do w/ product
 - ↳ chain rule

- take log to get rid of product
 - how is this ok?

- if strictly monotonic increasing \rightarrow will be increasing

$$= \arg \max_p \log \left(\prod_{i=1}^n P(X_i; p) \right)$$

↳ 1 to 1 map

$$= \arg \max_p \sum_{i=1}^n \log \binom{3}{k_i} p^{k_i} (1-p)^{3-k_i}$$

log of product =
sum of products

$$= \arg \max_p \sum_{i=1}^n \log \binom{3}{k_i} + \sum_{i=1}^n k_i \log p + \sum_{i=1}^n (3-k_i) \log (1-p)$$

derivative

$$= 0 + \sum_{i=1}^n \frac{k_i}{p} - \sum_{i=1}^n \frac{(3-k_i)}{1-p}$$

Set = 0

$$\log a^x = x \log a$$

(9)

Simplify

$$\hat{P}_{n_{ML}} = \frac{1}{3n} \sum_{i=1}^n \cancel{k_i}$$

↑ "makes sense"

if 20 out of 300 flipped

$$p \approx \frac{20}{300}$$

$$= \frac{1}{3n} \sum_{i=1}^n x_i$$

↑ since valid for all k_i
(little x)

(c) Is estimator unbiased

- Bayesian \rightarrow calc MBE error, p error (to check if good)- Classical $\rightarrow E[\hat{P}_n] = p$ our guess for

$$= E\left[\frac{1}{3n} \sum_{i=1}^n x_i\right]$$

$$= \underbrace{\sum_{i=1}^n E[x_i]}_{3n}$$

$$= \frac{3pn}{3n} = p \quad \text{✓ yes}$$

(10)

↓) Consistent? $\hat{P}_n \xrightarrow{\text{in prob}} p$

~~WLLN~~ converges in prob $\lim_{n \rightarrow \infty} P(|\hat{P}_n - p| \geq \delta) = 0 \quad \forall \epsilon > 0$



← gets more + more centered
around p

$$\hat{P}_n = \frac{\sum_{i=1}^n X_i}{n} \rightarrow p$$

Sample mean
 $\frac{\sum X_i}{n} \rightarrow p$

~~WLLN~~ Due to WLLN

c) Confidence Interval

- report interval, not just estimate

$$E[\hat{P}_n^-, \hat{P}_n^+]$$

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

Given
Confidence interval = 95%

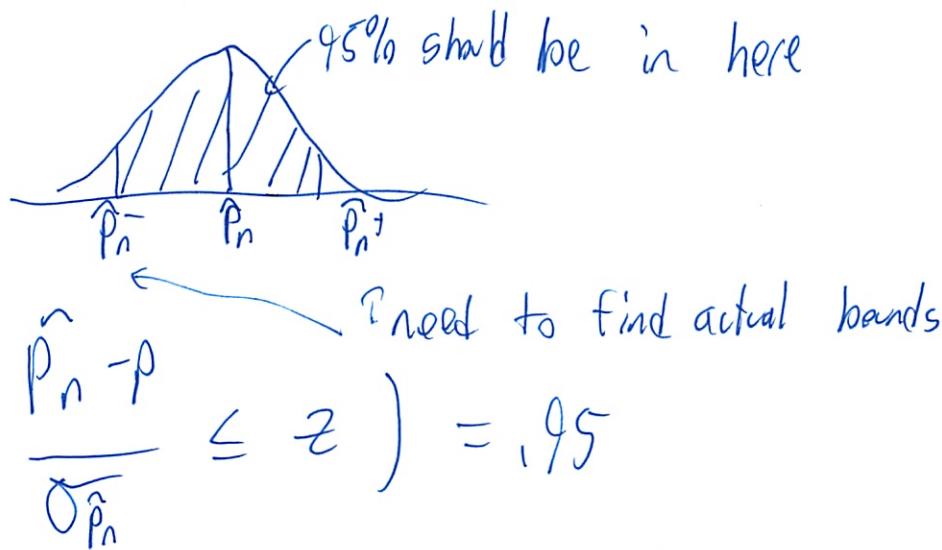
(11)

Start w) finding \hat{P}_n^-, \hat{P}_n^+

So start w/ \hat{P}_n

$$\hat{P}_n = \underbrace{\frac{1}{m} \sum_{n=1}^{3^m} Y_i}_{\begin{matrix} \text{\# total bits} \\ \text{\# bits} \end{matrix}} \quad \leftarrow \begin{matrix} \text{sum of 1 bits} \\ \in \text{sum up all bits,} \\ \text{not just all messages} \end{matrix}$$

$\hat{P}_n \approx$ normal for large n



$$P\left(-z \leq \frac{\hat{P}_n - P}{\sqrt{\hat{P}_n}} \leq z\right) = .95$$

↑ need to find actual bands

$$z = ?$$

Look at st normal so CDF is .025

$$z = 1.96 \quad \text{↑ going backwards}$$

~~$$\hat{P}_n^+ = \hat{P}_n + 1.96$$

$$\hat{P}_n^- = \hat{P}_n - 1.96$$~~

(2) but need to get it out of "normal form"
 $P(-1.96 \leq \frac{\hat{P}_n - P}{\sigma} \leq 1.96) \geq .995$ -denormalize

$$P(\hat{P}_n - 1.96 \sqrt{\frac{V}{n}} \leq P \leq \hat{P}_n + 1.96 \sqrt{\frac{V}{n}} \leftarrow \text{var})$$

\hat{P}_n total #

If don't ~~will~~ know var, guess!

Conservative to for bernoulli

Or use sample var $\hat{P}_n(1 - \hat{P}_n)$

$$\hat{P}_n^{\#} = \hat{P}_n + 1.96 \sqrt{\frac{1/4}{300}}$$

= a #