

## LECTURE 24

- Reference: Section 9.3
- Course VI Underground Guide Evaluations  
<https://sixweb.mit.edu/student/evaluate/6.041-f2010>  
<https://sixweb.mit.edu/student/evaluate/6.431-f2010>

## Outline

more classical stats

Will still be  
on final  
Note statement  
conflicted in  
review section

- Review
  - Maximum likelihood estimation
  - Confidence intervals
- Linear regression
- Binary hypothesis testing
  - Types of error
  - Likelihood ratio test (LRT)

Some ways to mis use

## Review

- Maximum likelihood estimation
  - Have model with unknown parameters:  
 $X \sim p_X(x; \theta)$
  - Pick  $\theta$  that "makes data most likely"

$$\max_{\theta} p_X(x; \theta)$$

 $\theta$  is constant we do not know

- Compare to Bayesian MAP estimation:

$$\max_{\theta} p_{\theta|X}(\theta | x) \text{ or } \max_{\theta} \frac{p_X(x|\theta)p_{\theta}(\theta)}{p_Y(y)}$$

Got x

Look at max likelihood

Wait value of  $\theta$  where  $x$  is most likely to happen

- Sample mean estimate of  $\theta = E[X]$

$$\hat{\theta}_n = (X_1 + \dots + X_n)/n$$

Collect samples - find avg

- $1 - \alpha$  confidence interval

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

 $\theta$  is RVFind value of  $\theta$ , which this is largest  
So most likely  $\theta$ 

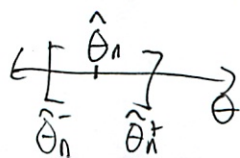
- confidence interval for sample mean

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

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

Same if prior is constant

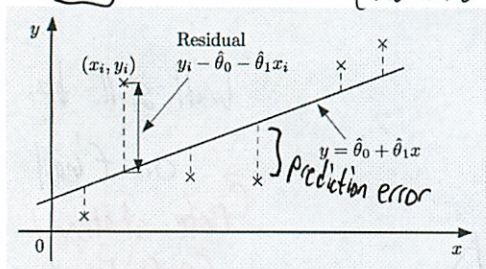
Denominator does not matter

So  $\theta$  inside interval w/ prob  $1 - \alpha$ or  $1 - \alpha$  points will fall within $P(2.1 \leq \theta \leq 3.9) \geq .95 \in \text{false w/ little } \theta \text{ it's a \#}$ So leave it as RVs - the interval - just that interval falls around  $\theta$

ie 95% of experiments will fall inside interval. (can use CLT to approx #)

$$P\left(\frac{\sqrt{n}(\bar{\theta}_n - \theta)}{\sigma} \leq z\right) \approx 2(1 - \Phi(z))$$

Regression (Remember from high school)



half of what statisticians do - most common model

Do my  $x$ s provide useful info of predicting  $y$ s

$X$  = HS GPA

$Y$  = MIT GPA

Find line that best fits data

- Data:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

- Model:  $y \approx \theta_0 + \theta_1 x$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \quad (*) \in \text{sum of square of error as small as possible}$$

- One interpretation: Relation to Prob  
 $Y_i = \theta_0 + \theta_1 x_i + W_i$ ,  $W_i \sim N(0, \sigma^2)$ , i.i.d.

- Likelihood function  $f_{X,Y|\theta}(x, y; \theta)$  is:

$$c \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \right\}$$

- same as ML w/ some assumptions

Random element  
 $\sim$  normal

- Take logs, same as (\*)

- Least sq.  $\leftrightarrow$  pretend  $W_i$  i.i.d. normal

density function of  $y$  vector

mean  
(same var  $\sigma^2$ )

try to make as large as possible to exponential  
then as small as possible

$$f_{Y_1, Y_2, \dots, Y_n}(y_1, y_2, \dots, y_n; \theta_0, \theta_1) = \prod_{i=1}^n f_{Y_i}(y_i; \theta)$$

## Linear regression

- Model  $y \approx \theta_0 + \theta_1 x$

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

derivative will be linear

2 unknowns  $\theta_0, \theta_1$

- Solution (set derivatives to zero):

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

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

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

solution intuitive

$$Y = \theta_0 + \theta_1 X + W$$

$$E[Y] = \theta_0 + \theta_1 E[X]$$

$$\text{suppose } E[X] = E[Y] = 0 \\ E[W] = 0$$

$$\bar{y} = \theta_0 + \theta_1 \bar{x} \rightarrow \hat{\theta}_0 = \bar{y} - \theta_1 \bar{x}$$

$$E[Y^2] = \theta_0 + \theta_1 E[X^2]$$

$$\hat{\theta}_1 = \frac{E[XY]}{\text{Var}(X)} = \frac{\text{cov}(X, Y)}{\text{Var}(X)}$$

$$E[XY] = \theta_0 + \theta_1 E[X^2] + E[XW]$$

- Interpretation of the form of the solution

- Assume a model  $Y = \theta_0 + \theta_1 X + W$   
 $W$  independent of  $X$  and  $Y$ , with zero mean
- Check that

$$\theta_1 = \text{cov}(X, Y) = \frac{E[(X - E[X])(Y - E[Y])]}{\text{Var}(X)}$$

- Solution formula for  $\hat{\theta}_1$  is a natural estimate of the covariance

natural estimate of var



## The world of linear regression

### Multiple linear regression:

- data:  $(x_i, x'_i, x''_i, y_i), i = 1, \dots, n$  if more than 1 variable

- model:  $y \approx \theta_0 + \theta x + \theta' x' + \theta'' x''$

- formulation:

$$\min_{\theta, \theta', \theta''} \sum_{i=1}^n (y_i - \theta_0 - \theta x_i - \theta' x'_i - \theta'' x''_i)^2$$

$X_1 = \text{HS GPA}$

$X_2 = \text{Family Income}$

$X_3 = \text{SAT score}$

$Y = \text{MIT GPA}$

explanatory variables

### Choosing the right variables

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

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

- work with data points  $(y_i, h(x_i))$

- formulation:

$$\min_{\theta} \sum_{i=1}^n (y_i - \underbrace{\theta_0 - \theta_1 h_1(x_i)}_{\text{outcome prediction}})^2$$

again set deriv = 0 and solve  
must make a choice

- just diff set of explanatory variables

## The world of regression (ctd.)

### In practice, one also reports

- Confidence intervals for the  $\theta_i$  - not explaining, but established methodology

- "Standard error" (estimate of  $\sigma$ )

-  $R^2$ , a measure of "explanatory power"

"factor explains 20% of variation"

- some noise will always be there

### Some common concerns

- Heteroskedasticity

- Multicollinearity

- Sometimes misused to conclude causal relations - biggest mistake

- etc.

Heteroskedasticity

var gets bigger  
as x gets larger

Multicollinearity

$$Y = \theta_0 + \theta_1 (\text{HS GPA}) + \theta_2 (\text{HS GPA})$$

its the same thing!  
or its closely related  
both can explain it

Fitting Line just tells

vs closely related  
- causality can go  
either way or  
none of the  
above - only  
association

- strongly affected by later data, less affected by low x data

lots  
of issues  
to watch  
out for

## Binary hypothesis testing as well as next time

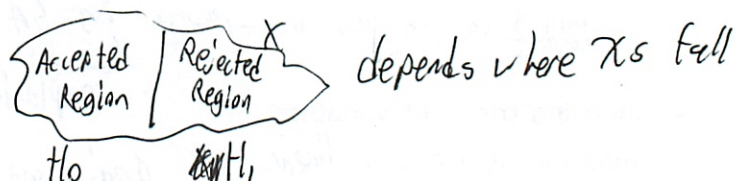
- Binary  $\theta$ ; new terminology:

- **null hypothesis  $H_0$ :** *default*  
 $X \sim p_X(x; H_0)$  [or  $f_X(x; H_0)$ ]
  - **alternative hypothesis  $H_1$ :** *new*  
 $X \sim p_X(x; H_1)$  [or  $f_X(x; H_1)$ ]
- two diff dist of X*  
*- did x come from 1 dist or the other*

- Partition the space of possible data vectors

**Rejection region  $R$ :**

reject  $H_0$  iff data  $\in R$



- Types of errors:

- **Type I (false rejection, false alarm):**

$H_0$  true, but rejected

$$\alpha(R) = P(X \in R; H_0)$$

*if shift size of region, ~~only make~~ tradeoff of errors*

- **Type II (false acceptance, missed detection):**

$H_0$  false, but accepted

$$\beta(R) = P(X \notin R; H_1)$$

*airline industry*

## Likelihood ratio test (LRT)

- Bayesian case (MAP rule): choose  $H_1$  if: *Bayes rule*  
 $P(H_1 | X = x) > P(H_0 | X = x)$   
*easy to see which error is bigger*

or

$$\frac{P(X = x | H_1)P(H_1)}{P(X = x)} > \frac{P(X = x | H_0)P(H_0)}{P(X = x)}$$

or

$$\frac{P(X = x | H_1)}{P(X = x | H_0)} > \frac{P(H_1)}{P(H_0)}$$

*decision procedure*

(likelihood ratio test)

*if odds in favor, choose  $H_1$*

- Nonbayesian version: choose  $H_1$  if

$$\frac{P(X = x; H_1)}{P(X = x; H_0)} > \xi \quad (\text{discrete case})$$

*Under  $H_1$ , how likely to observe this x*

$$\frac{f_X(x; H_1)}{f_X(x; H_0)} > \xi \quad (\text{continuous case})$$

*with x observed - is it more likely to be produced by  $H_0$  or  $H_1$*

- threshold  $\xi$  trades off the two types of error

- choose  $\xi$  so that  $P(\text{reject } H_0; H_0) = \alpha$   
 (e.g.,  $\alpha = 0.05$ )



- (a) Find the ML estimates of the linear model parameters.
- (b) Find the ML estimates of the quadratic model parameters.

Note: You may use the regression formulas and the connection with ML described in pages 478-479 of the text. However, the regression material is outside the scope of the final.

The figure below shows the data points  $(x_i, y_i)$ ,  $i = 1, \dots, 5$ , the estimated linear model

$$y = 40.53x - 65.86,$$

and the estimated quadratic model

$$y = 4.09x^2 - 3.07.$$

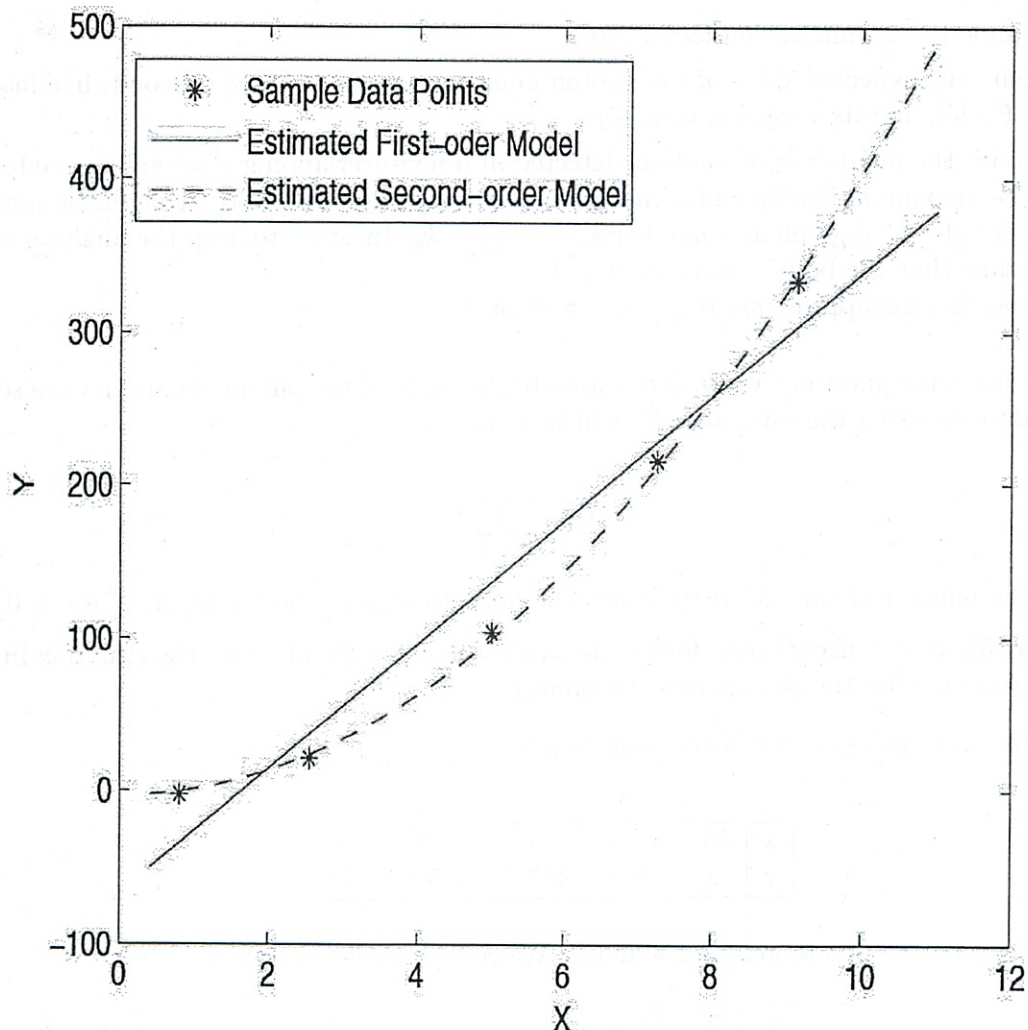


Figure 1: Regression Plot

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

---

Recitation 24  
December 7, 2010

1. A blackbody at temperature  $\theta$  radiates photons of all wavelengths, described by its characteristic spectrum. This problem will have you estimate  $\theta$ , which is fixed but unknown. The PMF for the number of photons  $K$  in a given wavelength range and a fixed very short time interval is given by,

$$p_K(k; \theta) = \frac{1}{Z(\theta)} e^{-k/\theta}, k = 0, 1, 2, \dots$$

$Z(\theta)$  is a normalization factor for the probability distribution (the physicists call it the partition function). You are given the task of determining the temperature of the body to two significant digits by photon counting in non-overlapping time intervals of duration one second. The photon emissions in non-overlapping time intervals are statistically independent from each other.

- (a) Determine the normalization factor  $Z(\theta)$ .
- (b) Compute the expected value of the photon number measured in any 1 second time interval,  $\mu_K = E_\theta[K]$ , and its variance,  $\text{var}_\theta(K) = \sigma_K^2$ .
- (c) You count the number  $k_i$  of photons detected in  $n$  non-overlapping 1 second time intervals. Find the maximum likelihood estimator,  $\hat{\theta}_n$ , for temperature  $\theta$ . Note, it might be useful to introduce the average photon number  $s_n = \frac{1}{n} \sum_{i=1}^n k_i$ . In order to keep the analysis simple we assume that the body is hot, i.e.  $\theta \gg 1$ .  
You may use the approximation:  $\frac{1}{e^{1/\theta} - 1} \approx \theta$  for  $\theta \gg 1$ .

In the following questions we wish to estimate the mean of the photon count in a one second time interval using the estimator  $\hat{K}$ , which is given by,

$$\hat{K} = \frac{1}{n} \sum_{i=1}^n K_i.$$

- (d) Find the number of samples  $n$  for which the noise to signal ratio for  $\hat{K}$ , (i.e.,  $\frac{\sigma_{\hat{K}}}{\mu_{\hat{K}}}$ ), is 0.01.
  - (e) Find a 95% confidence interval for the mean photon count estimate for the situation in part (d). (You may use the central limit theorem.)
2. Given the five data pairs  $(x_i, y_i)$  in the table below,

x	0.8	2.5	5	7.3	9.1
y	-2.3	20.9	103.5	215.8	334

we want to construct a model relating  $x$  and  $y$ . We consider a linear model

$$Y_i = \theta_0 + \theta_1 x_i + W_i, \quad i = 1, \dots, 5,$$

and a quadratic model

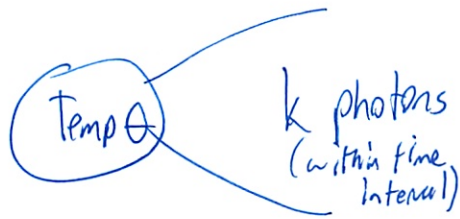
$$Y_i = \beta_0 + \beta_1 x_i^2 + V_i, \quad i = 1, \dots, 5.$$

where  $W_i$  and  $V_i$  represent additive noise terms, modeled by independent normal random variables with mean zero and variance  $\sigma_1^2$  and  $\sigma_2^2$ , respectively.

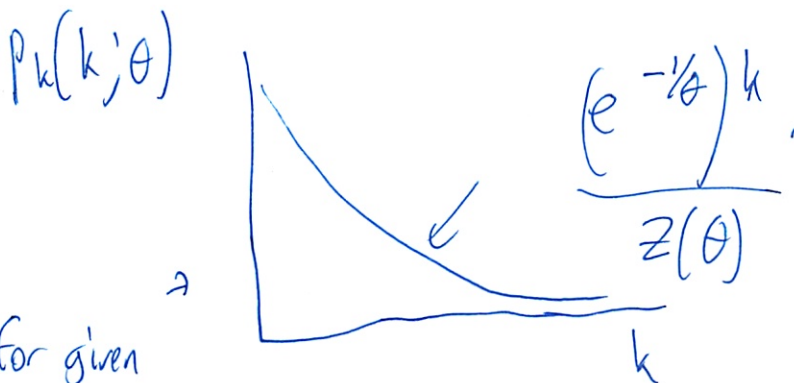


- Interference/chap 9 review

1. Material



Measure photons to estimate temp



for given  
value of

$\Theta$   $\hookrightarrow$  changes param of geometric  $\nearrow$   
but diff pick for each  $\Theta$

a) What is  $Z(\Theta)$ ?

~~Not~~

$$1 = \frac{1}{Z(\Theta)} \sum_{k=0}^{\infty} (e^{-1/\Theta})^k$$

$$= \frac{1}{Z(\Theta)} \frac{1}{1 - e^{-1/\Theta}}$$

$$Z(\Theta) = \frac{1}{1 - e^{-1/\Theta}}$$

②

b) Find mean + var  $\mu_k = E_\theta[k]$   $\sigma_k^2 = \text{var}_\theta(k)$

- is a geometric

- but is shifted ~~to the~~ to right by 1

↳ param  $p = (1 - e^{-1/\theta})$

$$\mu_k = \frac{1}{p} - 1$$

$$= \frac{1}{1 - e^{-1/\theta} - 1}$$

$$= \frac{1}{e^{1/\theta} - 1}$$

$$\text{var}_\theta(k) = \sigma_k^2 = \frac{1-p}{p^2}$$

$\underbrace{\hspace{1cm}}$   
var of  
geometric

shift does not matter

$$= \frac{e^{-1/\theta}}{(1 - e^{-1/\theta})^2}$$

$$\dots (\text{algebra})$$
$$= \mu_k^2 + \mu_k$$



3)

1) Get  $k_1, \dots, k_n$  iid samples of  $k$

- want to inter temp

- Find  $\hat{\theta}_n$

- Use max likelihood of  $\theta$  given  $k_1, \dots, k_n$

~~the~~

Likelihood  
function

$$p_{k_1, \dots, k_n}(k_1, \dots, k_n; \theta) = \frac{1}{(Z(\theta))^n} \prod_{i=1}^n e^{-k_i/\theta}$$

by indep.

- lets us take log
- so easier to take deriv
- set = 0 to maximize

$$\max_{\theta} \log \text{likelihood} = -n \log(Z(\theta)) - \frac{1}{\theta} \sum_{i=1}^n k_i$$

take deriv of both sides indep.

$$0 = \frac{d}{d\theta} \log \text{likelihood} = -n \frac{e^{-1/\theta}}{\theta^2 (1 - e^{-1/\theta})} + \frac{1}{\theta^2} \sum_{i=1}^n k_i$$

$\underbrace{\frac{e^{-1/\theta}}{\theta^2 (1 - e^{-1/\theta})}}_{\text{is } \frac{d}{d\theta} \left( \log \left( \frac{1}{1 - e^{-1/\theta}} \right) \right)} \quad \text{is } Z(\theta)$

4

$$\frac{1}{e^{1/\theta} - 1} = \frac{1}{n} \sum_{i=1}^n k_i$$

Solve for  $\theta$  to get ML

Non linear

So can't solve easier in closed form

~~then~~ trick: linearize it

$$e^{1/\theta} \approx 1 + \frac{1}{\theta} + \underbrace{O(\theta^2 + \dots)}_{\text{does not matter}}$$

$$\hat{\theta}_n \approx \frac{1}{\frac{1}{n} \sum_{i=1}^n k_i} = \hat{k}_n$$

approx = to  $\theta$  by linearization

- model's units don't agree

- but said 1 implies the other

Can do CLT, var, CI, etc now

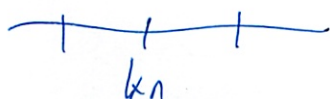
- dealing w/ sample mean

now can answer this

d) What is # of samples needed so that ~~var~~  $\hat{k}_n$

$$\frac{\sigma_{\hat{k}_n}}{\hat{k}_n} \text{ is } \leq .01$$

← like a  $\frac{\text{noise}}{\text{signal}}$  ratio.  
 $\leq 1\%$





(5)

# ~~sketch~~

$$\frac{\sigma_k / \sqrt{n}}{\mu_k} \leq .01$$

$$\sqrt{n} \geq \frac{\sqrt{\mu_k^2 + \mu_k}}{.01 \cdot \mu_k}$$

$$= 100 \sqrt{1 + \frac{1}{\mu_k}}$$

If temp is high, then # photons emitted is high

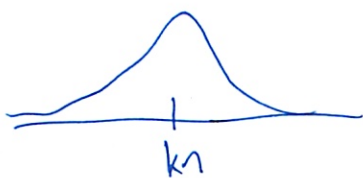
For  $\theta \gg 1$ ,  $\mu_k \gg 1$

So  $\frac{1}{\mu_k}$  goes away

And  $n$  should be  $\geq 10,000$

e) Find a 95% CI for  $\hat{k}_n$

For  $n = 10,000$



← approx normal

Want upper + lower limit so .45 of values inside  
See how many st. dev to left + right from normal table

a)

$$\begin{array}{c} \text{---} \\ | \quad \quad | \\ \hat{\mu}_n \quad \hat{\mu}_n + 1.96 \cdot \hat{\sigma} \hat{\mu}_n \\ | \\ \hat{\mu}_n - 1.96 \cdot \hat{\sigma} \hat{\mu}_n \end{array}$$

$$\left[ \hat{\mu}_n - 1.96 \cdot \frac{1}{\sqrt{n}} \hat{\sigma} \hat{\mu}_n, \hat{\mu}_n + 1.96 \cdot \frac{1}{\sqrt{n}} \hat{\sigma} \hat{\mu}_n \right]$$

- all very standard
- expect on exam
- will be at least 1 q
- ~~the~~ mid term comprehensive, but focus on new material

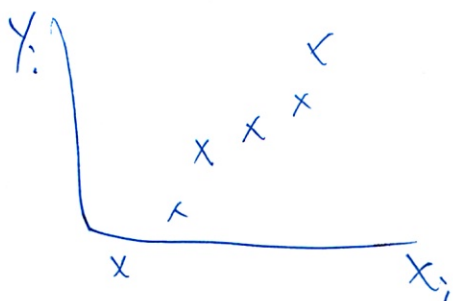
## 2. Regression

- at 1st glance no relation to prob.
- but has statistical interpretations  $\rightarrow$  ML

First do it non-stat way

Construct a model by relating  $X$  and  $Y$

$\uparrow$                        $\uparrow$   
 input                      output





⑦

Think  $x = \text{time}$   
 $y = \text{position along 1 axis}$

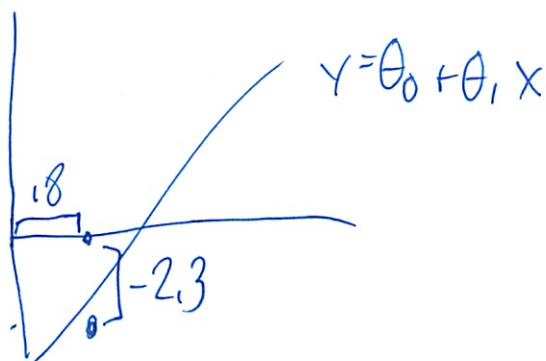
Want speed/velocity of car

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

$\uparrow$        $\uparrow$        $\uparrow$        $\uparrow$   
 pos    initial pos    velocity    time

Given data points on a table

$x$	1.8	...
$y$	2.3	...



Want vertical discrepancies to be as small as possible

Minimize

$$\min \sum_{i=1}^n \overbrace{(y_i - \theta_0 - \theta_1 x_i)^2}^{\text{that vertical discrepancy}}$$

simple regression model  
linear

8

take deriv

Set = to 0

$$\frac{\partial}{\partial \theta_0} \sum ( )^2 = 0 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)$$

$$\frac{\partial}{\partial \theta_1} \sum ( )^2 = 0 = \sum_{i=1}^n x_i (y_i - \theta_0 - \theta_1 x_i)$$

Now system of 2 eq w/ 2 unknowns

If know  $\theta_0$  - plug in, ...

If know  $\theta_1$ , <sup>↑ initial pos</sup> convert to  $\theta_0$  w/

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

↳ where  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  = output sample mean

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

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

Divide top + bottom by n

Like a sample cov

$$\approx \frac{\text{cov}(x, y)}{\text{var}(x)}$$

← that's where this ratio comes in  
- coefficient for x - in other things we've seen

9

For 5 data points

$$\bar{X} = 4.94$$

$$\bar{Y} = 134.38$$

$$\hat{\theta}_1 = 40.53$$

$$\hat{\theta}_0 = -65.86$$

} coefficients for the straight line

Represents the graph on pg 2 recitation handout

$$Y = 40.53 X - 65.86$$

---

Now try to connect w/ stats

- ~~construct~~ use ML

- to do dist, CI, etc

ML Interpretation

$$Y_i = \theta_0 + \theta_1 X_i + W_i$$

$\uparrow$   
not RV

$\uparrow$  noise  
 $\sim \text{normal}(0, \sigma^2)$

$\uparrow$  mean  $\uparrow$  var

$Y_i$  is a sample of  $Y$



(10)

Likelihood function of  $y_1, \dots, y_n$

$$= \text{Constant} \cdot e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2}$$

Maximize likelihood of  $\theta_0, \theta_1$

= ~~minimizing~~ the exponent  
max

$$= \underset{\theta_0, \theta_1}{\text{minimizing}} \text{ the expression } \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

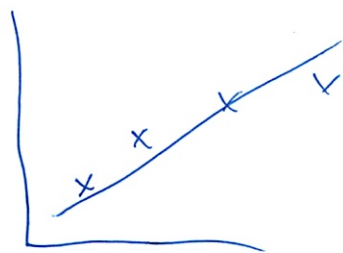
What is the advantage of doing it the states way?

$\hat{\theta}_0, \hat{\theta}_1$  become RVs which get a distribution

But also their var can be estimated

↳ to form CIs

Many other things about regression



Sometimes better w/ higher  
order polynomial

Instead of assuming linear

- assume quadratic

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

Get a better fit

⑪

Same formulas

- but w/  $(X_i)^2$  not  $(x_i)$
- so just  $^2$  wherever there is an  $X_i$

Or also non 0 slope at y axis

$$Y = \theta_0 + \theta_1 X + \theta_2 X^2$$

- can do, but not w/ our formulas
- more complex
- must differentiate 3 parts
- solve 3 eq, 3 unknowns
- Computer can do it

nonlinear  
regression

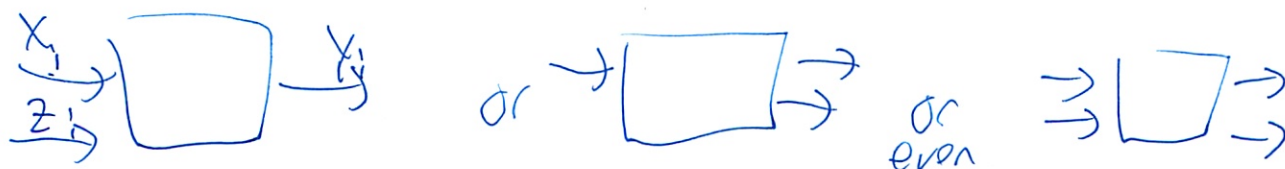
Called multi-parameter regression

Now



Multiple regression

But can also have



$$Y_i = \theta_0 + \theta_1 X_{i1} + \theta_2 X_{i2}$$

# LECTURE 25

## Outline

- Reference: Section 9.4

last lecture

- Course VI Underground Guide Evaluations

<https://sixweb.mit.edu/student/evaluate/6.041-f2010>

<https://sixweb.mit.edu/student/evaluate/6.431-f2010>

more classical stats

- Review of simple binary hypothesis tests

test hypothesis

— examples

- Testing composite hypotheses

— is my coin fair?

— is my die fair?

— goodness of fit tests

## Simple binary hypothesis testing

No prior beliefs about hyp

- null hypothesis  $H_0$ :

$$X \sim p_X(x; H_0) \quad [\text{or } f_X(x; H_0)]$$

- vector of several RVs

- alternative hypothesis  $H_1$ :

$$X \sim p_X(x; H_1) \quad [\text{or } f_X(x; H_1)]$$

- Choose a rejection region  $R$ ;  
reject  $H_0$  iff data  $\in R$



depends where data falls

- Likelihood ratio test: reject  $H_0$  if

$$\frac{p_X(x; H_1)}{p_X(x; H_0)} > \xi \quad \text{or} \quad \frac{f_X(x; H_1)}{f_X(x; H_0)} > \xi$$

if you design acceptance region - have a choice

Two types of errors

- false negatives

- false positive

- fix false rejection probability  $\alpha$ ; (e.g.,  $\alpha = 0.05$ )

- choose  $\xi$  so that  $P(\text{reject } H_0; H_0) = \alpha$

compare likelihood ratio to some constant threshold param

- but how chose it?

- Set a ~~parameter~~ criteria

- then solve to get  $\xi$



### Example (test for normal mean)

- $n$  data points, i.i.d.

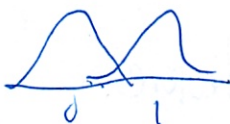
$$H_0: X_i \sim N(0, 1)$$

$$H_1: X_i \sim N(1, 1)$$

two hypotheses

both normal

but different means



- Likelihood ratio test; rejection region:

$$\frac{(1/\sqrt{2\pi})^n \exp\{-\sum_i (X_i - 1)^2/2\}}{(1/\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/2\}} > \xi$$

look at likelihood ratio

- algebra: reject  $H_0$  if:  $\sum_i X_i > \xi'$

Write down joint densities under  $H_1$  top  
right - evaluated at observed values  $x$   
 $H_0$  bottom

- Find  $\xi'$  such that

$$P\left(\sum_{i=1}^n X_i > \xi'; H_0\right) = \alpha$$

- use normal tables

Compare what we got w/ threshold

$$-\sum_i (X_i - 1)^2 + \sum_i X_i^2 > \log \xi$$

$$2\sum_i X_i - n > \log \xi$$

If sum of  $X_i$ 's big - than  
choose  $H_1$  - ie reject  $H_0$

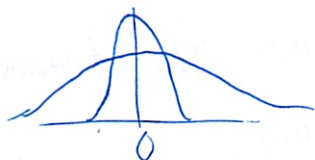
### Example (test for normal variance)

- $n$  data points, i.i.d.

$$H_0: X_i \sim N(0, 1)$$

$$H_1: X_i \sim N(0, 4)$$

here same mean



Simple since both hyp  
are fully specified  
probabilistic models

- Likelihood ratio test; rejection region:

$$\frac{(1/2\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/(2 \cdot 4)\}}{(1/\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/2\}} > \xi$$

same cookbook procedure

- algebra: reject  $H_0$  if  $\sum_i X_i^2 > \xi'$

take log, like before

- Find  $\xi'$  such that

$$P\left(\sum_{i=1}^n X_i^2 > \xi'; H_0\right) = \alpha$$

prob of rejecting hypothesis

- the distribution of  $\sum_i X_i^2$  is known  
(derived distribution problem)

- "chi-square" distribution; form not all that important  
tables are available

look up 95%<sup>th</sup> tile to find threshold

## Composite hypotheses

- Got  $S = 472$  heads in  $n = 1000$  tosses; is the coin fair?  
 $H_0: p = 1/2$  versus  $H_1: p \neq 1/2$   
*coin fair*
- Pick a "statistic" (e.g.,  $S$ )
- Pick shape of rejection region (e.g.,  $|S - n/2| > \xi$ )  
*reject if very different*
- Choose significance level (e.g.,  $\alpha = 0.05$ )  
*threshold*
- Pick critical value  $\xi$  so that:

$$P(\text{reject } H_0; H_0) = \alpha$$

Using the CLT:

$$P(|S - 500| \leq 31; H_0) \approx 0.95; \xi = 31$$

- In our example:  $|S - 500| = 28 < \xi$   
 $H_0$  not rejected (at the 5% level)

*delebrate choice of words*

Can't say  $H_1$  is accepted

And will not be exactly .500 often

But are the obs compatible (close) to  $H_0$   
 so  $H_0$  could be true (not rejected)  
Is my die fair?

$H_0$  is fully specified  
 = not fully specified model  $H_1$ ,  
 lots of alternative models

- $p = .6$
- $p = .7$
- $p = .499$

Can't write likelihood ratio anymore

- Since many candidate distributions

Build test on statistic

- Function of data  $X_i$ s

- Compress into 1 #

- here  $S = 472$  # of heads

- not always clear what to use

- Broader philosophical: Theories can't be proven correct - only that has not been rejected

If ~~rejection~~

$$H_0 = p = 1/2$$

$$H_1 = p > 1/2$$

then rejection region  
 if  $S \geq \frac{n}{2} + \frac{n}{2}$

[A] R

(current region

, R [A] R

- Hypothesis  $H_0$ :  
 $P(X = i) = p_i = 1/6, i = 1, \dots, 6$

- Observed occurrences of  $i$ :  $N_i$

- Choose form of rejection region;  
 chi-square test:

$$\text{reject } H_0 \text{ if } T = \sum_i \frac{(N_i - np_i)^2}{np_i} > \xi$$

- Choose  $\xi$  so that:

$$P(\text{reject } H_0; H_0) = 0.05$$

$$P(T > \xi; H_0) = 0.05$$

- Need the distribution of  $T$ :  
 (CLT + derived distribution problem)

- for large  $n$ ,  $T$  has approximately a chi-square distribution
- available in tables

how many times each was observed

Here what should  $S$  be?

✓ value  $i$  shows up  $\frac{n}{6}$  times supposedly

- what actually happens?

↑ take squares

divide by factor - different weights depending on  $P_i$

- but does not matter here - constant

↑ if big values are far away  
 choose like before

Formula or table

- if  $P_i$  is equal - what type of dist is  $T$ ?

-  $N_i$  is binomial ( $1/6$ )


- CLT is approx same as normal - mean of normal

- take square of a normal - sum is chi squared (earlier today)



More complicated versions of this

### Do I have the correct pdf?

Hyp:  $i$  RVs distributed by this PDF 

- Partition the range into bins
  - $np_i$ : expected incidence of bin  $i$  (from the pdf)
  - $N_i$ : observed incidence of bin  $i$
  - Use chi-square test (as in die problem)

- Split into bins

- PDF ~~say~~ predicts prob for each bin

- Count how often got a result in ~~each~~ each bin

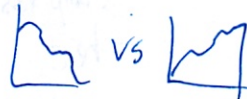
- Compare w/ expected 1/f of obsence

- add over is

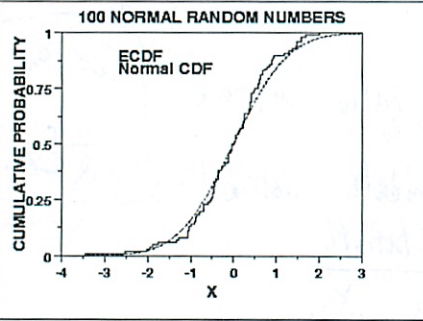
$$\sum_i (N_i - np_i)^2$$

that w/  $np_i$   
is chi squared

- Can't distinguish b/w same PDF for a certain region



- Only tells you if compatible w/ hyp - survived so far  
- not if true



(<http://www.itl.nist.gov/div898/handbook/>)

- $D_n = \max_x |F_X(x) - \hat{F}_X(x)|$
- $P(\sqrt{n}D_n \geq 1.36) \approx 0.05$  ✓ prob false rejection

Other more sophisticated methods  
- if want no binning  
- compare CDF  
- what fraction of data below some #  
- much approach prob of falling below that #  
- need to find threshold value  
- Once have it, know what test must be

### What else is there?

In statistics

- Systematic methods for coming up with shape of rejection regions likelihood ratio test
- Methods to estimate an unknown PDF (e.g., form a histogram and "smooth" it out) where coming from?
- Efficient and recursive signal processing can you do it quickly?
- Methods to select between less or more complex models
- (e.g., identify relevant "explanatory variables" in regression models)
- Methods tailored to high-dimensional unknown parameter vectors and huge number of data points (data mining)
- etc. etc....

want informative 2D summary

want to be able to analyze

- i.e. form tables

- no unique correct ans - room for judgement - art  
if get results you want - you stick w/ it

- he hopes the better methods will arrive

- also need it to be computationally implementable

~~can~~ always find very complex model that fits data exactly  
- but that does not match reality

where work is today  
takes many parameters to explain data point  
huge # of data points  
- advance class next semester



Previous chap: Bayesian approach to inference

- Unknown params modeled as RVs

but here  $\theta$  is known (not random)

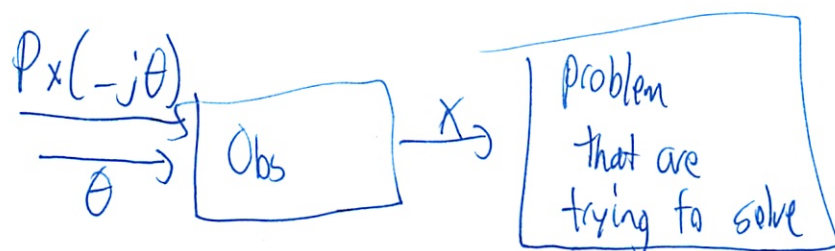
Obs  $X$  is random

$P_X(x; \theta)$  or  $f_X(x; \theta)$  depends on value of  $\theta$

↳ multiple candidate models

One for each value of  $\theta$

So want to find ~~best~~ best (at of all candidate models)



$E_\theta[h(x)]$  = expected value of RV  $h(x)$  as function of  $\theta$

$P_\theta(A)$  = prob of event  $A$

↳ functional dependence - not conditioning

Types of problems

- parameter estimation - estimates that are nearly correct  
under any possible value of unknown param.

- hypothesis testing - unknown param takes finite #  
of values  $m$  ( $m \geq 2$ ) - want to pick one w/  
least error

(2)

- Significance testing - want to accept/reject 1 hypothesis while keeping prob of false rejection small

## Inference Methods

- Maximum Likelihood (ML) - parameter that makes observed data most likely - max. chance of obtaining it equivalent  
MAP
- Linear Regression - linear relation that minimizes sum of square of error b/w line + data
- Likelihood Ratio Test - Given two <sup>hyp</sup> cat<sup>ns</sup>, compare their relative (ratio) of likelihood
- Significance testing - Given a hyp, reject it obs data falls outside of certain rejection region

## 9.1 Classical Parameter Estimation

$\theta$  is unknown constant

ML is classical equivalent of MAP

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

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

Distribution of  $X$  depends on  $\theta$ , also dist of  $\hat{\theta}$  depends of  $\theta$

$\theta$  = estimate (the actual value)

③

Sometimes interested in ~~the~~ ~~at~~ when  $n \neq \infty$ ,  $n = \#$  of obs

$\hat{\theta}_n$  = estimator

↳ <sup>actually</sup> seq of estimators - one for each  $n$

↳ mean =  $E_{\theta}(\hat{\theta}_n)$  ) # functions of  $\theta$

↳ var =  $\text{var}_{\theta}(\hat{\theta}_n)$

$\tilde{\theta}_n$  = estimation error =  $\hat{\theta}_n - \theta$

bias =  $b_{\theta}(\hat{\theta}_n)$  = expected value of estimation error

$$= E_{\theta}[\hat{\theta}_n] - \theta$$

↳  $E[\text{bias}]$  depend on  $\theta$

↳  $\text{var}(\text{bias})$  " " "

↳ estimation error also depends on observations  $X_1, X_2, \dots, X_n$

---

$\hat{\theta}_n$  is unbiased if  $E_{\theta}[\hat{\theta}_n] = \theta$  for every possible  $\theta$

$\hat{\theta}_n$  is asymptotically unbiased if  $\lim_{n \rightarrow \infty} E_{\theta}[\hat{\theta}_n] = \theta$   
for every possible  $\theta$

$\hat{\theta}_n$  is consistent if seq  $\hat{\theta}_n$  converges to true value of  $\theta$   
in prob for every possible value of  $\theta$



④

Estimator will not be  $=$  to  $\theta$  exactly

↳ so estimation error will be non 0

But if estimation error avg is 0  $\rightarrow$  unbiased

If only  $\uparrow$  as  $n \uparrow$ , then asymptotically unbiased

Also interested in size of estimation error

↳ mean squared error  $E_{\theta}[\hat{\theta}_n^2]$

$$E_{\theta}[\hat{\theta}_n^2] = b_{\theta}^2(\hat{\theta}_n) + \text{var}_{\theta}(\hat{\theta}_n)$$

In many problems tradeoff b/w the two terms on right hand side  
- ie if  $\downarrow$  var ~~var~~ bias?

(can I see a graphical example of that?)

- goal is ~~to~~ to keep both small

ML Estimator

$X = X_1, \dots, X_n$  = vector of obs

describe w/ joint PMF  $p_X(X; \theta)$

- form depends on an unknown (scalar or vector)  $\theta$

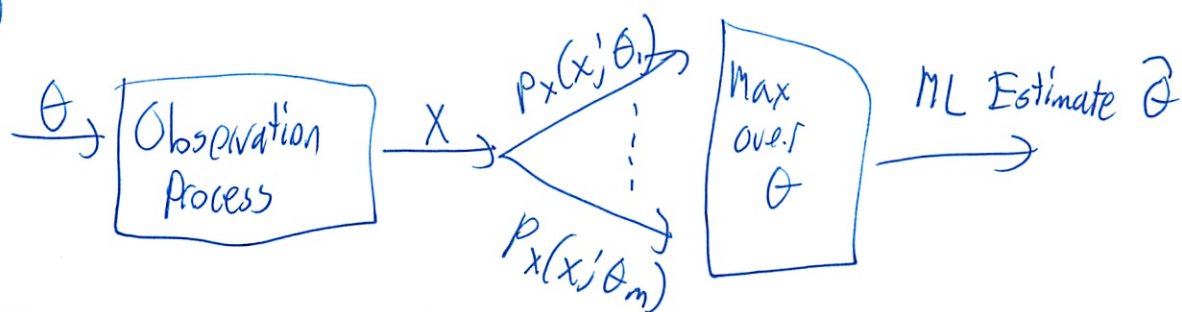
Observe a certain value  $x = (x_1, \dots, x_n)$

Maximize  $p_X(x_1, \dots, x_n; \theta)$  over all  $\theta$

$$\hat{\theta}_n = \arg \max_{\theta} p_X(x_1, \dots, x_n; \theta)$$

or  $f$  if continuous

5



In many cases observations  $X_i$  are assumed to be ind  
So likelihood function of the form

$$P_X(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n p_{X_i}(x_i; \theta)$$

But easier to max log  $\rightarrow$  log-likelihood function over  $\theta$

$$\begin{aligned} \log P_X(x_1, \dots, x_n; \theta) &= \log \prod_{i=1}^n p_{X_i}(x_i; \theta) \\ &= \sum_{i=1}^n \log p_{X_i}(x_i; \theta) \end{aligned}$$

likelihood  $\rightarrow$  is not parameter that unknown param =  $\theta$

is instead prob that observed value  $x$  can arise when parameter is = to  $\theta$

"What is the value of  $\theta$  under which the obs we have seen are most likely to arise"

with a flat prior they are the same

if one-to-one function - just apply function to estimate  
(skipping examples)

# ⑥ Estimation of Mean and Var of RV

- simple, but important problem of estimating mean + var
- we don't need to know about distribution

- Sample mean - most natural estimator of  $\theta$

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

- unbiased since  $E_\theta[M_n] = E_\theta[X] = \theta$

- MSF = var =  $\frac{V}{n}$  ← common var of  $X_i$

- does not depend on  $\theta$

- by WLLT  $\rightarrow$  converges to  $\theta$  in probability = consistent

- Sample mean not necessarily w/ smallest var estimator

- two var estimators

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

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

- coincides w/ ML if  $X_i$ s are normal

- biased, but asymptotically unbiased

- unbiased

for large  $n$  the two estimators coincide



⑦

(?but which to use?)

- one is based on ~~MM~~ ML
- second is scaled to be unbiased

## Confidence Interval

Have estimator  $\hat{\theta}_n$  of unknown param  $\theta$

So confidence  $\theta$  is within interval

Confidence level  $1 - \alpha$

Set bands

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

↙ confidence interval

(I think I get it conceptually - just need practice)

for every possible value of  $\theta$

↑ (this kinda confuses me how they worded it)

Remember in classical stats its the ~~conf~~ interval that is random

$\theta$  is fixed

Well 95% of intervals will include  $\theta$

Usually constructed by forming an interval around an estimator - want smallest possible width - but this depends on  $\theta$  - but this is usually asymptotically normal and asymptotically unbiased

⑧

So 
$$\frac{\hat{\theta}_n - \theta}{\sqrt{\text{Var}(\hat{\theta}_n)}}$$

CDF ~~of  $\hat{\theta}_n$~~  approaches st normal as  $n \uparrow$  for every  $\theta$

### Estimator Var Approx

- if var known, just use that
- if not known, need to estimate
- Using the unbiased estimator

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

- Estimate var  $\frac{V}{n}$  of sample mean by  $\frac{S_n^2}{n}$
- So for a given  $\alpha$ , use these estimates + CLT to construct approx  $1-\alpha$  confidence interval

- So 
$$\left[ \hat{\theta}_n - z \frac{\hat{S}_n}{\sqrt{n}}, \hat{\theta}_n + z \frac{\hat{S}_n}{\sqrt{n}} \right]$$

Where  $z$  is obtained from

$$\Phi(z) = 1 - \frac{\alpha}{2}$$

and normal tables

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

↪ backwards in table

9

- So get

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

- two diff approx in effect

- treating  $\hat{\theta}_n$  as if normal

- replacing var w/ estimate

- ~~is~~ always only an approx (I

$\hat{s}_n^2$  only approx to true var  $V$

- The RV  $T_n = \frac{\sqrt{n}(\hat{\theta}_n - \theta)}{\hat{s}_n}$  is not normal

- but PDF of  $T_n$  does not depend on  $\theta$  or  $V$

- so can compute explicitly

- t-distribution w/  $n-1$  degrees of freedom

- symmetric + bell shaped - but more spread out w/  
heavier tails

- When  $X_i$  normal,  $n$  relatively small

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

-  $z$  obtained from relation

$$P_{n-1}(z) = 1 - \frac{\alpha}{2}$$

(10)

So when  $n$  is large ( $n \geq 50$ )  $t$ -distribution is close to normal  $\rightarrow$  so can use normal tables

(did not do in class I think)

Otherwise we are  $t$  tables  $\rightarrow$  CDF  $\psi_{n-1}(z)$   
w/ # of degrees of freedom  
desired tail prob  
and  $z$  value

Diff estimators for var possible

- if Bernoulli  $v = \theta(1-\theta)$  of  $x$

$$\text{sn } \text{var} = \hat{\theta}_n(1-\hat{\theta}_n)$$

- as  $n \uparrow$  to  $\infty$   ~~$\hat{\theta}_n(1-\hat{\theta}_n) \rightarrow \theta(1-\theta)$~~

$$\hat{\theta}_n \rightarrow \theta \text{ in prob}$$

$$\text{sn } \hat{\theta}_n(1-\hat{\theta}_n) \rightarrow v$$

- Or say that  $\theta(1-\theta) \leq 1/4$  for  $\theta \in [0,1]$   
and use  $1/4$  as conservative var estimate



(11)

## 9.2 Linear Regressions

- building model of relation b/w 2 or more variables of interest
- can explain it simply
- or under guise of probability  
(did today in recitation)

Will start off w/ two variables

- have  $x$  = years of edu  
 $y$  = income
- have data pairs  $(x_i, y_i)$   
 $i = 1, \dots, n$

$x_i$  = years of edu,  $y_i$  = income of  $i$ th person

~~Rep~~ Represent linearly

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

$\uparrow$   $\uparrow$   
unknown params to be estimated

Given estimates  $\hat{\theta}_0, \hat{\theta}_1$  of resulting params,  $y_i$  corresponds to  $x_i$   
as predicted by the model

$$\hat{y}_i = \hat{\theta}_0 + \hat{\theta}_1 x_i$$

(12)

The "error"  $y_i$  is given by

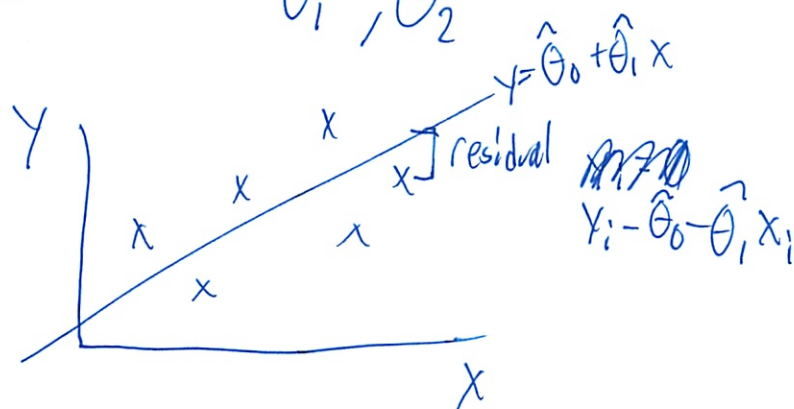
$$\tilde{y}_i = y_i - \hat{y}_i \quad \text{called the } i^{\text{th}} \text{ residual}$$

Want to choose estimates so have a small residual

So minimize the square of the residuals

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

Over all  $\theta_1, \theta_2$



The postulated linear model may or may not be true

- if relation is actually nonlinear
- in order to use this model we assumed linearity

To derive and get  $\hat{\theta}_0$  and  $\hat{\theta}_1$  - take partial derivatives + set = 0

$$\theta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

(13)

(can be justified many ways  
(skipping))

## Bayesian Linear Regression

12/8

Can explain linear regression w/ prob Basian Framework

$X_1, \dots, X_n$  = given  $\#$

$Y_1, \dots, Y_n$  = Observed values of Vector  $Y = (Y_1, \dots, Y_n)$   
of RVs that obey linear relation

$$Y_i = \theta_0 + \theta_1 X_i + W_i$$

$\theta = (\theta_0, \theta_1)$  is parameter to be estimated

$W_1, \dots, W_n$  are i.i.d RVs  
w/ mean 0 and known var  $\sigma^2$

$\theta_0, \theta_1, W_1, \dots, W_n$  are indep.

$\theta_0, \theta_1$  have mean 0 and var  $\sigma_0^2, \sigma_1^2$

Derive a Bayesian estimator based on MAP

Assume  $\theta_0, \theta_1, W_1, \dots, W_n$  are <sup>normal</sup> RV

Maximize ~~the~~  $\theta_0, \theta_1$  over posterior PDF

$$\epsilon_0(\theta_0, \theta_1) f_{Y|\theta}(Y_1, \dots, Y_n | \theta_0, \theta_1)$$



(14)

Divided by a positive normalization constant that does not depend on  $(\theta_0, \theta_1)$

This is

$$C \cdot \exp\left\{-\frac{\theta_0^2}{2\sigma_0^2}\right\} \cdot \exp\left\{-\frac{\theta_1^2}{2\sigma_1^2}\right\} \cdot \prod_{i=1}^n \exp\left\{-\frac{(y_i - \theta_0 - x_i \theta_1)^2}{2\sigma^2}\right\}$$

where  $C =$  normalizing constant

Equivalently, ~~normalize over~~ minimize over  $\theta_0, \theta_1$

$$\frac{\theta_0^2}{2\sigma_0^2} + \frac{\theta_1^2}{2\sigma_1^2} + \sum_{i=1}^n \frac{(y_i - \theta_0 - x_i \theta_1)^2}{2\sigma^2}$$

$\uparrow$  Same as w/ classical

$\uparrow$  would be identical if  $\sigma_0, \sigma_1$  so large, these terms disappear  
Take partial derivs

Set  $= 0$

$$\hat{\theta}_1 = \frac{\sigma_1^2}{\sigma^2 + \sigma_1^2 \sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\hat{\theta}_0 = \frac{n\sigma_0^2}{\sigma^2 + n\sigma_0^2} (\bar{y} - \hat{\theta}_1 \bar{x})$$

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

(15)

Some remarks - skipping

## Multiple Linear Regression

- so far have had single explanatory variable  $x$ 
  - ↳ simple regression
- but often multiple underlying or explanatory variables
  - ↳ multiple regression

if had triplets of data  $(x_i, y_i, z_i)$  and wanted  $\theta_j$   
$$y \approx \theta_0 + \theta_1 x + \theta_2 z$$

then minimize  $\sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i - \theta_2 z_i)^2$

no limit on possible explanatory variables  
 $\uparrow \quad \uparrow \quad \uparrow$   
~~minimize~~ minimize over

## Non Linear Regression

- if we assume model is nonlinear

$$y \approx h(x; \theta)$$
  
 $\uparrow$  given function       $\searrow$  param to be estimate

- given data pairs  $(x_i, y_i) \quad i = 1, \dots, n$
- seek value of  $\theta$  that minimizes sm of squared residuals

$$\sum_{i=1}^n (y_i - h(x_i; \theta))^2$$

16

But no general closed form solution

But can try w/ enough computational power

- ~~we~~ can do via ML estimation

$$Y_i = h(x_i; \theta) + W_i$$

iid normal  
w/ 0 mean

- likelihood function takes form

$$f_Y(Y; \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(Y_i - h(x_i; \theta))^2}{2\sigma^2}\right\}$$

- Heteroskedasticity

- var of  $W_i$  could vary w/  $X_i$

- so noise at the end overwhelms noise at the beginning

- can adjust weighting to avoid

- Nonlinearity

- Sometimes pick a non linear model

- Multicollinearity

- if  $x, z$  closely related model may not be able to distinguish b/w the two

- Overfitting

- fits data well, but model useful

- if your polynomial has too high of an order

- must be 5-10 x data points than parameter being estimated

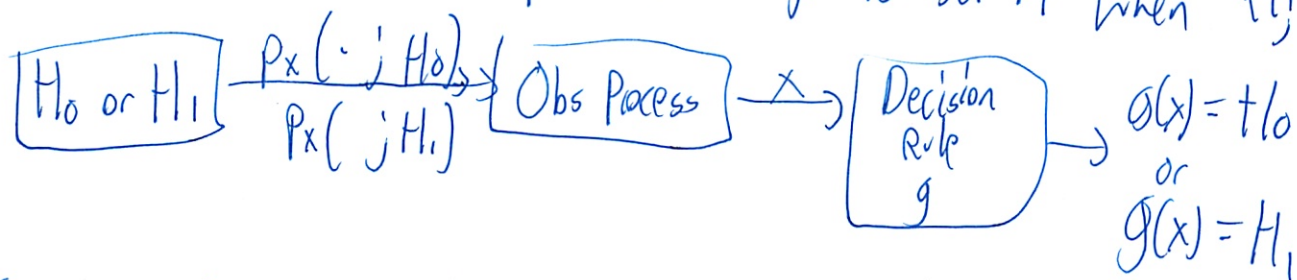
Causality

the holy grail of science - causation vs correlation

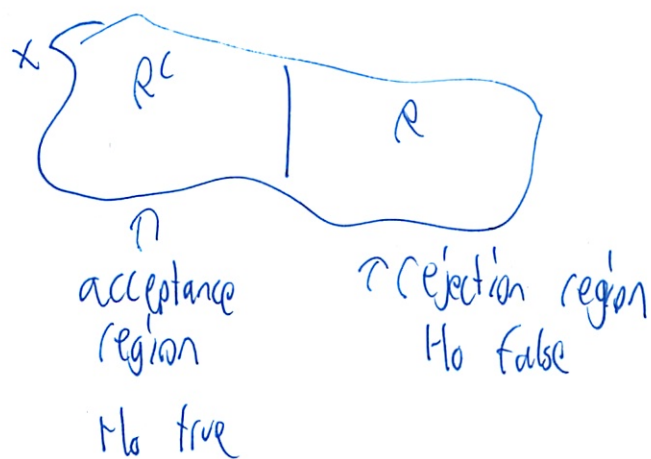


## (17) 9.3 Binary Hypothesis Testing

- back to choosing b/w 2 hyp
  - but unlike 8.2 section - no prior prob
  - like an inference problem where  $\theta$  takes 2 values
    - but call  $H_0 \equiv$  null / default  $\leftarrow$  prove or disprove
    - $H_1 \equiv$  alternate
  - Observation is vector  $X = (X_1, \dots, X_n)$  of RV
    - dist depends on hyp
- $P(X \in A; H_j)$  to denote prob  $X$  belongs to set  $A$  when  $H_j$  is true



- Can partition observations into 2 subsets



18

Two possible types of errors

a) False positive / Type I / false rejection

Reject  $H_0$ , even though  $H_0$  is true

Prob of happening

$$\alpha(R) = P(X \in R; H_0)$$

b) False negative / Type II / False acceptance

Accept  $H_0$  even though  $H_0$  is false

$$\beta(R) = P(X \notin R; H_1)$$

To decide shape of region, minimize prob of error w/ MAP rule!

Given observed value  $x$  of  $X$ , declare  $\Theta = \Theta_1$  to be true if

$$p_{\Theta}(\Theta_0) p_{X|\Theta}(x|\Theta_0) < p_{\Theta}(\Theta_1) p_{X|\Theta}(x|\Theta_1)$$

Rewrite as Likelihood Ratio  $L(x)$

$$L(x) = \frac{p_{X|\Theta}(x|\Theta_1)}{p_{X|\Theta}(x|\Theta_0)}$$

Declare  $\Theta = \Theta_1$  to be true if realized value  $x$  of obs vector  $X$  satisfies

$$L(x) \geq \xi$$

(19)

Where  $\epsilon$  is the critical value

$$\epsilon = \frac{P_0(\theta_0)}{P_0(\theta_1)}$$

If  $X$  is continuous, same but w/  $f$

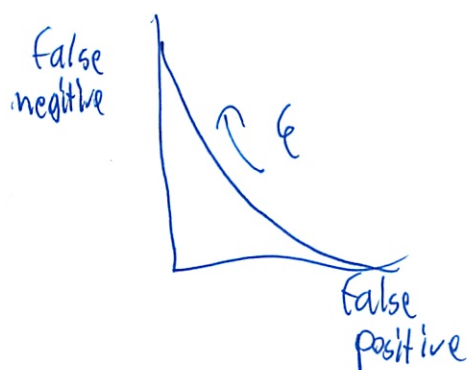
So consider rejection regions of the form

$$R = \{x \mid L(x) \geq \epsilon\}$$

$\epsilon$  chosen through other considerations

$\epsilon = 1$  is the ML rule

$\epsilon$  is a tradeoff b/w both types of error



So pick it using the Likelihood Ratio Test (LRT)

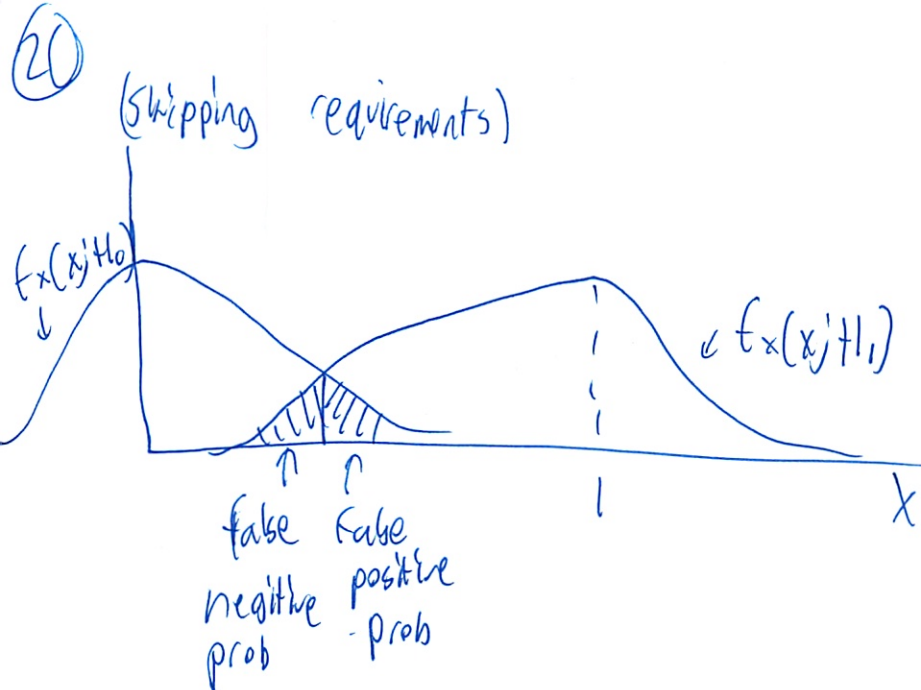
- start w/ target  $\alpha$  for false positive prob

~~Not~~ - Pick  $\epsilon$  so false ~~prob~~ pos prob = to  $\alpha$

$$P(L(x) \geq \epsilon \mid H_0) = \alpha$$

- once  $x$  observed, reject  $H_0$  if  $L(x) \geq \epsilon$





(I don't really get this chart)

(? If it is  $H_0$  will fall - but if in false pos territory will be said to be  $H_1$ , when it is  $H_0$ )

if falls in normal  $H_0$  or false negative prob territory say is  $H_0$

(kinda cool how they represent that)

## Neyman-Pearson Lemma

Consider a particular choice of  $c$  in the LRT which results in

$$P(L(x) \geq c; H_0) = \alpha \quad \text{and} \quad P(L(x) \leq c; H_1) = \beta$$

Suppose that some other test w/ rejection region  $R$  achieves smaller or = false positive prob

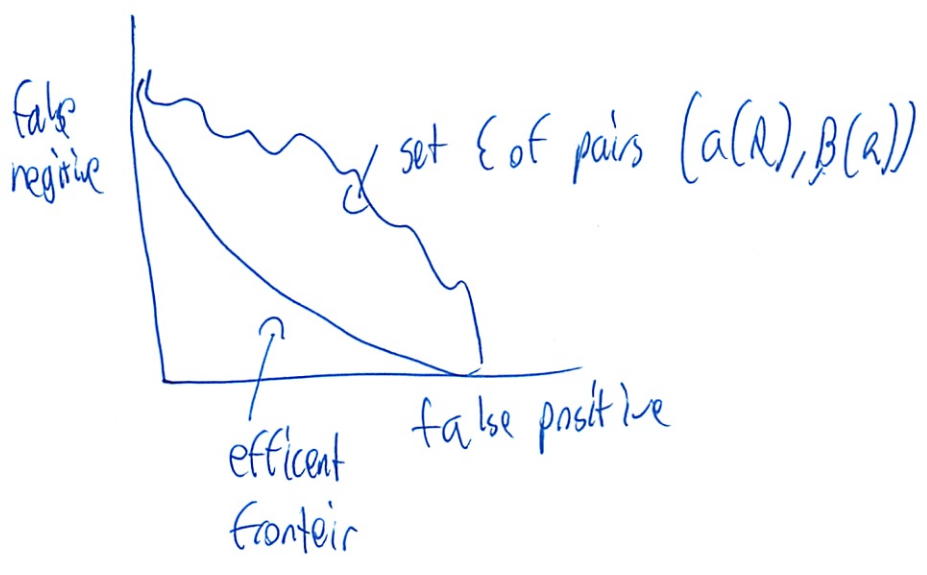
$$P(x \in R; H_0) \leq \alpha$$

Then  $P(x \in R; H_1) \geq \beta$

w/ strict inequality  $P(x \in R; H_1) > \beta$  when  $P(x \in R; H_0) < \alpha$

21

(I don't get that at all!)



Lemma states all pairs  $(a(\xi), B(\xi))$  lie on efficient frontier  
(still don't get - don't think we got into in class)

---

## 9.4 Significance Testing

Hypothesis testing problems do not always have 2 well specified alternatives, so can't use 9.3

This is more general

A lot more "art" / judgement here  
(more complex questions)

---

default hyp =  $H_0$  = null hypothesis

Want to determine based on observations  $X = (X_1, \dots, X_n)$  whether hyp should be accepted or rejected

(22)

Will restrict discussion to models w/ following characteristics

a) Parametric models - obs have dist governed by joint PMF/PDF

completely determined by  $\theta$ , belongs to  $M$  set of parameters

b) Simple Null hypothesis - ~~the~~ asserts true value of  $\theta$  is = to a given element  $\theta_0$  of  $M$

c) Alternate hypothesis -  $H_1$  that  $H_0$  is not true

(went over this in lecture today)

Coin toss example

coin tossed  $n=1000$  times

$\theta$  = unknown probability of each toss

set of all possible params  $M = [0, 1]$

$H_0$  = null hyp = the coin is fair  $\theta = \frac{1}{2}$

alt hyp =  $\theta \neq \frac{1}{2}$

Observe  $X_1, \dots, X_n$  tosses

$\hookrightarrow X_i = 0$  or  $\underset{\text{heads}}{1}$

$S = X_1 + \dots + X_n$

Use decision rule

reject  $H_0$  if  $|S - \frac{n}{2}| > c$

critical value to be determined



(23)

We have defined  $R$  (Rejection region)

↳ set of data vectors that lead to

rejection of null hyp  
Choose  $\epsilon$  so prob false positive / fake rejection =  $\alpha$

$$P(\text{reject } H_0 | H_0) = \alpha$$

↑ significance level, here  $\alpha = .05$

Now need to make some choices. Some prob calculations needed to determine critical value  $\epsilon$

Under null hyp,  $S$  is binomial w/  $p = \frac{1}{2}$ ,  $n = 1000$

Use normal approx to binomial + normal tables

$$\epsilon = 31$$

If observed  $S = s = 472$

$$|s - 500| = |472 - 500| = 28 \leq 31$$

and  $H_0$  is not rejected at 5% significance level

Only say that observed  $S$  does not provide strong evidence against hypothesis  $H_0$



(24)

## Significance Testing Methodology

A stat. test of hyp  $H_0 = \theta = \theta^*$  is performed based on obs  $X_1, \dots, X_n$

1. Choose a statistic  $S$  that is a scalar RV  
 - choose some ~~function~~ function  $h: R^n \rightarrow R$   
 resulting in statistic  $S = h(X_1, \dots, X_n)$

2. Determine shape of rejection region

- Specify set of values of  $S$  for which  $H_0$  will be ~~rejected~~ rejected as a function of critical value  $c$

3. Choose the significance level

- ie desired prob of a false rejection of  $H_0$

4. Choose the critical value  $c$  so prob of false positive  $\approx \alpha$   
 Rejection region is determined

5. Once  $X_1, \dots, X_n$  are observed

i) Calculate statistic  $s = h(x_1, \dots, x_n)$

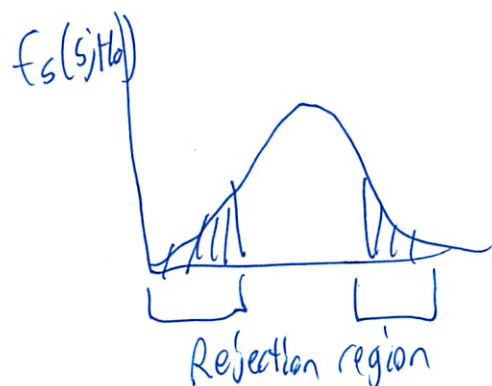
ii) Reject hyp  $H_0$  if  $s$  is in rejection region

(25)

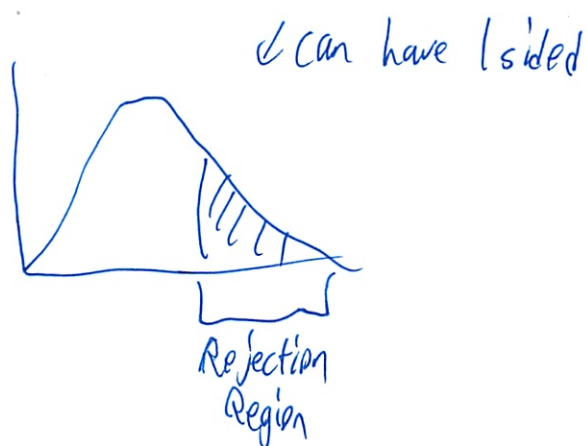
1. No right way to choose  $S$

- Sometimes obvious
- sometimes need to generalize likelihood ratios
- need to make sure it is easy to calculate

2. ~~Set~~ <sup>values</sup> of  $S$  under which  $H_0$  not rejected is usually an interval  
Surrounding peak dist of  $S$  under  $H_0$



or



3. Pick  $\alpha$  to balance tradeoff

4. Step 4 (previous page) only place prob. calculations used

- requires dist of  $L(x)$  ~~from~~
- Sometimes  $\log L(x)$
- ~~if~~ often  $S$  can not be found in closed form
- often need approx like CLT
  - but only if  $n$  is large
- may estimate  $S$  by simulation

Say  $H_0$  is rejected at the  $\alpha$  significance level



(26)

- Does not mean prob of  $H_0$  being true is less than  $\alpha$
- Instead will have false rejection/false positive a fraction  $\alpha$  of the time

Statisticians often replace step 3+4 (2 pages ago) w/ p-value

$$\text{p-value} = \min \{ \alpha \mid H_0 \text{ would be rejected at the } \alpha \text{ sig level} \}$$

- <sup>at</sup> the threshold b/w rejection + non-rejection
- ie null hyp would be rejected at 5% sig level if

$$\text{p-value} < .05$$

(skipping examples)

## Generalized Likelihood Ratio + Goodness of Fit Tests

- test whether a given ~~PMF~~ PMF conforms w/ observed data
- = goodness of fit
- Use it as an introduction to general methodology for Significance testing in face of composite alt<sup>n</sup> hypothesis
- Consider RV that takes values in finite set  $\{1, \dots, m\}$
- Let  $\theta_k$  be prob of outcome  $k$
- Dist of RV is described by vector param  $\theta = (\theta_1, \dots, \theta_m)$
- Consider hyp  $H_0: \theta = (\theta_1^*, \dots, \theta_m^*)$   $H_1: \theta \neq (\theta_1^*, \dots, \theta_m^*)$

(27)

$\theta_k^*$  are non neg # that sum to 1

Draw  $n$  ind. samples of RV

$N_k = \#$  of samples that result in outcome  $k$

~~Dist of RV~~

Obs  $X = (N_1, \dots, N_m)$

denote realized value by  $x = (n_1, \dots, n_m)$

$$N_1 + \dots + N_m = n_1 + \dots + n_m = n$$

Do generalized likelihood ratio test

1. estimate a model by ML

- ie determine a param vector  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_m)$   
that maximizes the likelihood fn  $p_x(x; \theta)$  over all  
vectors  $\theta$

2. Carry out a LRT that compares likelihood  $p_x(x; \theta^*)$   
under  $H_0$  to likelihood  $p_x(x; \hat{\theta})$  corresponding to  
estimated model. I.e form

$$\frac{p_x(x; \hat{\theta})}{p_x(x; \theta^*)}$$

and if it exceeds critical value  $c$ , reject  $H_0$   
Choose  $c$  so prob false rejection/positive  $\approx \alpha$

(28)

(we did not do this in class I believe)

Basically we are asking

- is there a model compatible w/  $H_1$  that provides a better explanation for the observed data than that provided by model corresponding to  $H_0$ ?
- To ans: compare likelihood under  $H_0$  to largest possible likelihood under models compatible w/  $H_1$ .

First step: ML estimation  $\rightarrow$  involves a maximization over the set of prob. distributions  $(\theta_0, \dots, \theta_m)$   
PMF of obs vector  $x$  is multipomical

$$p_x(x; \theta) = c \theta_1^{n_1} \dots \theta_m^{n_m}$$

normalizing constant

Easier to work w/ Log likelihood

$$\log p_x(x; \theta) = \log c + n_1 \log \theta_1 + \dots + n_{m-1} \log \theta_{m-1} + n_m \log (1 - \theta_1 - \dots - \theta_{m-1})$$

Set derivs of  $\theta_1, \dots, \theta_{m-1}$  = to 0

$$\frac{n_k}{\hat{\theta}_k} = \frac{n_m}{1 - \hat{\theta}_1 - \dots - \hat{\theta}_{m-1}} \text{ for } k = 1, \dots, m-1$$

$\tau$  all the ratios must be =

$$\hat{\theta}_k = \frac{n_k}{n} \quad k = 1, \dots, m$$



(29)

These are correct ML estimates even if some of the  $n_k = 0$   
 - so corresponding  $\hat{\theta}_k$  are also 0

Generalized form

$$\text{reject } H_0 \text{ if } \frac{p_X(x; \hat{\theta})}{p_X(x; \theta^*)} = \prod_{k=1}^m \frac{(n_k/n)^{n_k}}{(\theta_k^*)^{n_k}} > \epsilon$$

Take log to simplify

$$\text{reject } H_0 \text{ if } \sum_{k=1}^m n_k \log\left(\frac{n_k}{n\theta_k^*}\right) > \log \epsilon$$

Need to determine  $\epsilon$  by taking into account req. sig level

$$P(S > \log \epsilon; H_0) = \alpha$$

where

$$S = \sum_{k=1}^m N_k \log\left(\frac{N_k}{n\theta_k^*}\right)$$

But dist of  $S$  under  $H_0$  not readily available - can only simulate

But if  $n$  is large  $\rightarrow$  can simplify

- observed freq  $\hat{\theta}_k = \frac{n_k}{n}$  will be close to  $\theta_k^*$  under  $H_0$   
 w/ high prob

Second order Taylor series expansion shows that our  
 statistic  $S$  can be app. well by  $I_2$ , where

$$T = \sum_{k=1}^m \frac{(n_k - n\theta_k^*)^2}{n\theta_k^*}$$

(30)

When  $n$  is large, it is known under  $H_0$  the dist of  $\hat{T}$  (and dist  $2S$ ) approaches a so-called " $\chi^2$ " dist w/  $m-1$  degrees of freedom

- dist available in tables
- approx correct values of  $P(T > r; H_0)$  or  $P(2S > r; H_0)$  can be obtained from  $\chi^2$  table
- can use to determine a suitable critical value that corresponds to given significance level  $\alpha$

So have test for large values of  $n$

### Chi Square Test

- Use statistic, ~~or related statistic~~  
$$S = \sum_{k=1}^m N_k \log \left( \frac{N_k}{n \theta_k^*} \right)$$

(or related statistic  $T$ )  
and rejection region

reject  $H_0$  if  $2S > r$

(or  $T > r$ )

- Critical value  $r$  is determined from CDF tables for  $\chi^2$  dist w/  $m-1$  degrees of freedom so that  
$$P(2S > r; H_0) = \alpha \text{ given stat. level}$$

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

---

**Problem Set 11**  
**Never Due**  
**Covered on Final Exam**

**1. Problem 7, page 509 in textbook**

Derive the ML estimator of the parameter of a Poisson random variable based of i.i.d. observations  $X_1, \dots, X_n$ . Is the estimator unbiased and consistent?

2. Caleb builds a particle detector and uses it to measure radiation from far stars. On any given day, the number of particles  $Y$  that hit the detector is conditionally distributed according to a Poisson distribution conditioned on parameter  $x$ . The parameter  $x$  is unknown and is modeled as the value of a random variable  $X$ , exponentially distributed with parameter  $\mu$  as follows.

$$f_X(x) = \begin{cases} \mu e^{-\mu x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Then, the conditional PDF of the number of particles hitting the detector is,

$$p_{Y|X}(y | x) = \begin{cases} \frac{e^{-x} x^y}{y!} & y = 0, 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find the MAP estimate of  $X$  from the observed particle count  $y$ .
- (b) Our goal is to find the conditional expectation estimator for  $X$  from the observed particle count  $y$ .
- i. Show that the posterior probability distribution for  $X$  given  $Y$  is of the form

$$f_{X|Y}(x | y) = \frac{\lambda^{y+1}}{y!} x^y e^{-\lambda x}, \quad x > 0$$

and find the parameter  $\lambda$ . You may find the following equality useful (it is obviously true if the equation above describes a true PDF):

$$\int_0^\infty a^{y+1} x^y e^{-ax} dx = y! \quad \text{for any } a > 0$$

- ii. Find the conditional expectation estimate of  $X$  from the observed particle count  $y$ .  
*Hint:* you might want to express  $x f_{X|Y}(x | y)$  in terms of  $f_{X|Y}(x | y + 1)$ .
- (c) Compare the two estimators you constructed in part (a) and part (b).
3. Consider a Bernoulli process  $X_1, X_2, X_3, \dots$  with unknown probability of success  $q$ . 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  $\{t_1, t_2, t_3, \dots\}$ . In problem set 10, we solved the problem using Bayesian inference. Our focus here will be on classical estimation.



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

---

We assume that  $q$  is an unknown parameter in the interval  $(0, 1]$ . Denote the true parameter by  $q^*$ . Denote by  $\hat{Q}_k$  the maximum likelihood estimate (MLE) of  $q$  given  $k$  recordings,  $T_1 = t_1, \dots, T_k = t_k$ .

- (a) Compute  $\hat{Q}_k$ . Is this different from the MAP estimate you found in problem set 10?
- (b) Show that for all  $\epsilon > 0$

$$\lim_{k \rightarrow \infty} \mathbf{P} \left( \left| \frac{1}{\hat{Q}_k} - \frac{1}{q^*} \right| > \epsilon \right) = 0$$

- (c) Assume  $q^* \geq 0.5$ . Give a lower bound on  $k$  such that

$$\mathbf{P} \left( \left| \frac{1}{\hat{Q}_k} - \frac{1}{q^*} \right| \leq 0.1 \right) \geq 0.95$$

4. A body at temperature  $\theta$  radiates photons at a given wavelength. This problem will have you estimate  $\theta$ , which is fixed but unknown. The PMF for the number of photons  $K$  in a given wavelength range and a fixed time interval of one second is given by,

$$p_K(k; \theta) = \frac{1}{Z(\theta)} e^{-\frac{k}{\theta}}, k = 0, 1, 2, \dots$$

$Z(\theta)$  is a normalization factor for the probability distribution (the physicists call it the partition function). You are given the task of determining the temperature of the body to two significant digits by photon counting in non-overlapping time intervals of duration one second. The photon emissions in non-overlapping time intervals are statistically independent from each other.

- (a) Determine the normalization factor  $Z(\theta)$ .
- (b) Compute the expected value of the photon number measured in any 1 second time interval,  $\mu_K = \mathbf{E}_\theta[K]$ , and its variance,  $\text{var}_\theta(K) = \sigma_K^2$ .
- (c) You count the number  $k_i$  of photons detected in  $n$  non-overlapping 1 second time intervals. Find the maximum likelihood estimator,  $\hat{\Theta}_n$ , for temperature  $\Theta$ . Note, it might be useful to introduce the average photon number  $s_n = \frac{1}{n} \sum_{i=1}^n k_i$ . In order to keep the analysis simple we assume that the body is hot, i.e.  $\theta \gg 1$ .  
You may use the approximation:  $\frac{1}{e^{\frac{1}{\theta}} - 1} \approx \theta$  for  $\theta \gg 1$ .

In the following questions we wish to estimate the mean of the photon count in a one second time interval using the estimator  $\hat{K}$ , which is given by,

$$\hat{K} = \frac{1}{n} \sum_{i=1}^n K_i.$$

- (d) Find the number of samples  $n$  for which the noise to signal ratio for  $\hat{K}$ , (i.e.,  $\frac{\sigma_{\hat{K}}}{\mu_{\hat{K}}}$ ), is 0.01.
  - (e) Find the 95% confidence interval for the mean photon count estimate for the situation in part (d). (You may use the central limit theorem.)
5. The RandomView window factory produces window panes. After manufacturing, 1000 panes were loaded onto a truck. The weight  $W_i$  of the  $i$ -th pane (in pounds) on the truck is modeled as a random variable, with the assumption that the  $W_i$ 's are independent and identically distributed.

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

---

- (a) Assume that the measured weight of the load on the truck was 2340 pounds, and that  $\text{var}(W_i) \leq 4$ . Find an approximate 95 percent confidence interval for  $\mu = \mathbf{E}[W_i]$ , using the Central Limit Theorem.
- (b) Now assume instead that the random variables  $W_i$  are i.i.d., with an exponential distribution with parameter  $\theta > 0$ , i.e., a distribution with PDF

$$f_W(w; \theta) = \theta e^{-\theta w}.$$

What is the maximum likelihood estimate of  $\theta$ , given that the truckload has weight 2340 pounds?

6. Given the five data pairs  $(x_i, y_i)$  in the table below,

x	0.8	2.5	5	7.3	9.1
y	-2.3	20.9	103.5	215.8	334

we want to construct a model relating  $x$  and  $y$ . We consider a linear model

$$Y_i = \theta_0 + \theta_1 x_i + W_i, \quad i = 1, \dots, 5,$$

and a quadratic model

$$Y_i = \beta_0 + \beta_1 x_i^2 + V_i, \quad i = 1, \dots, 5.$$

where  $W_i$  and  $V_i$  represent additive noise terms, modeled by independent normal random variables with mean zero and variance  $\sigma_1^2$  and  $\sigma_2^2$ , respectively.

- (a) Find the ML estimates of the linear model parameters.
- (b) Find the ML estimates of the quadratic model parameters.

Note: You may use the regression formulas and the connection with ML described in pages 478-479 of the text. However, the regression material is outside the scope of the final.

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

---

**Problem Set 11 Solutions**

1. Check book solutions on Stellar.
2. (a) To find the MAP estimate, we need to find the value  $x$  that maximizes the conditional density  $f_{X|Y}(x | y)$  by taking its derivative and setting it to 0.

$$\begin{aligned} f_{X|Y}(x | y) &= \frac{p_{Y|X}(y | x) \cdot f_X(x)}{p_Y(y)} \\ &= \frac{e^{-x} x^y}{y!} \cdot \mu e^{-\mu x} \cdot \frac{1}{p_Y(y)} \\ &= \frac{\mu}{y! p_Y(y)} \cdot e^{-(\mu+1)x} x^y \end{aligned}$$

$$\begin{aligned} \frac{d}{dx} f_{X|Y}(x | y) &= \frac{d}{dx} \left( \frac{\mu}{y! p_Y(y)} \cdot e^{-(\mu+1)x} x^y \right) \\ &= \frac{\mu}{y! p_Y(y)} x^{y-1} e^{-(\mu+1)x} (y - x(\mu + 1)) \end{aligned}$$

Since the only factor that depends on  $x$  which can take on the value 0 is  $(y - x(\mu + 1))$ , the maximum is achieved at

$$\hat{x}_{\text{MAP}}(y) = \frac{y}{1 + \mu}$$

It is easy to check that this value is indeed maximum (the first derivative changes from positive to negative at this value).

- (b) i. To show the given identity, we need to use Bayes' rule. We first compute the denominator,  $p_Y(y)$

$$\begin{aligned} p_Y(y) &= \int_0^\infty \frac{e^{-x} x^y}{y!} \mu e^{-\mu x} dx \\ &= \frac{\mu}{y! (1 + \mu)^{y+1}} \int_0^\infty (1 + \mu)^{y+1} x^y e^{-(1+\mu)x} dx \\ &= \frac{\mu}{(1 + \mu)^{y+1}} \end{aligned}$$

Then, we can substitute into the equation we had derived in part (a)

$$\begin{aligned} f_{X|Y}(x | y) &= \frac{\mu}{y! p_Y(y)} x^y e^{-(\mu+1)x} \\ &= \frac{\mu (1 + \mu)^{y+1}}{y! \mu} x^y e^{-(\mu+1)x} \\ &= \frac{(1 + \mu)^{y+1}}{y!} x^y e^{-(\mu+1)x} \end{aligned}$$

Thus,  $\lambda = 1 + \mu$ .



ii. We first manipulate  $xf_{X|Y}(x | y)$ :

$$\begin{aligned} xf_{X|Y}(x | y) &= \frac{(1 + \mu)^{y+1}}{y!} x^{y+1} e^{-(\mu+1)x} \\ &= \frac{y+1}{1+\mu} \frac{(1 + \mu)^{y+2}}{(y+1)!} x^{y+1} e^{-(\mu+1)x} \\ &= \frac{y+1}{1+\mu} f_{X|Y}(x | y+1) \end{aligned}$$

Now we can find the conditional expectation estimator:

$$\begin{aligned} \hat{x}_{\text{CE}}(y) &= \mathbf{E}[X|Y = y] = \int_0^\infty xf_{X|Y}(x | y) dx \\ &= \int_0^\infty \frac{y+1}{1+\mu} f_{X|Y}(x | y+1) dx = \frac{y+1}{1+\mu} \end{aligned}$$

(c) The conditional expectation estimator is always higher than the MAP estimator by  $\frac{1}{1+\mu}$ .

3. (a) The likelihood function is

$$\prod_{i=1}^k P_{T_i}(T_i = t_i | Q = q) = q^k (1 - q)^{\sum_{i=1}^k t_i - k}.$$

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

$$kq^{k-1}(1 - q)^{\sum_{i=1}^k t_i - k} - \left(\sum_{i=1}^k t_i - k\right)q^k(1 - q)^{\sum_{i=1}^k t_i - k - 1} = 0,$$

or equivalently

$$k(1 - q) - \left(\sum_{i=1}^k t_i - k\right)q = 0,$$

which yields  $\hat{Q}_k = \frac{k}{\sum_{i=1}^k t_i}$ . This is not different from the MAP estimate found before. Since the MAP estimate is calculated using a uniform prior, the likelihood function is a ‘scaled’ version of posterior probability and they can be maximized at the same value of  $q$ .

(b) Since  $\frac{1}{\hat{Q}_k} = \frac{\sum_{i=1}^k T_i}{k}$ , and that each  $T_i$  is independent identically distributed, it follows that  $\frac{1}{\hat{Q}_k}$  is actually a sample mean estimator. The weak law of large numbers says that, when the number of samples increases to infinity, the sample mean estimator converges to the actual mean, which is  $\frac{1}{q^*}$  in this case. So we can write the limit of probability as

$$\lim_{k \rightarrow \infty} \mathbf{P} \left( \left| \frac{1}{\hat{Q}_k} - \frac{1}{q^*} \right| > \epsilon \right) = \lim_{k \rightarrow \infty} \mathbf{P} \left( \left| \frac{\sum_{i=1}^k T_i}{k} - \mathbf{E}[T_1] \right| > \epsilon \right) = 0.$$

(c) Chebyshev inequality states that

$$\mathbf{P} \left( \left| \frac{\sum_{i=1}^k T_i}{k} - \mathbf{E}[T_1] \right| \geq \epsilon \right) \leq \frac{\text{var}(T_1)}{k\epsilon^2}.$$

So we have

$$\begin{aligned} \mathbf{P} \left( \left| \frac{1}{\hat{Q}_k} - \frac{1}{q^*} \right| \leq 0.1 \right) &= \mathbf{P} \left( \left| \frac{\sum_{i=1}^k T_i}{k} - \frac{1}{q^*} \right| \leq 0.1 \right) \\ &= 1 - \mathbf{P} \left( \left| \frac{\sum_{i=1}^k T_i}{k} - \mathbf{E}[T_1] \right| \geq 0.1 \right) \geq 1 - \frac{\text{var}(T_1)}{k * 0.1^2} \end{aligned}$$

To ensure the above probability to be greater than 0.95, we need that

$$1 - \frac{\text{var}(T_1)}{k * 0.1^2} = 1 - \frac{\frac{1-q}{q^2}}{k * 0.1^2} \geq 0.95,$$

or

$$k \geq 2000 \text{var}(T_1) = 2000 \frac{1-q}{q^2}$$

The number of observations  $k$  needed depends on the variance of  $T_1$ . For  $q$  close to 1, the variance is close to 0, and the required number of observations is very small (close to 0). For  $q = 1/2$ , the variance is maximum ( $\text{var}(T_1) = 2$ ), and we require  $k = 4000$ . Thus, to guarantee the required accuracy and confidence for all  $q$ , we need that,

$$k \geq 4000.$$

4. (a) Normalization of the distribution requires:

$$1 = \sum_{k=0}^{\infty} p_K(k; \theta) = \sum_{k=0}^{\infty} \frac{e^{-\frac{k}{\theta}}}{Z(\theta)} = \frac{1}{Z(\theta)} \sum_{k=0}^{\infty} e^{-\frac{k}{\theta}} = \frac{1}{Z(\theta) \cdot (1 - e^{-\frac{1}{\theta}})},$$

$$\text{so } Z(\theta) = \frac{1}{1 - e^{-\frac{1}{\theta}}}.$$

(b) Rewriting  $p_K(k; \theta)$  as:

$$p_K(k; \theta) = \left( e^{-\frac{1}{\theta}} \right)^k \left( 1 - e^{-\frac{1}{\theta}} \right), \quad k = 0, 1, \dots$$

the probability distribution for the photon number is a geometric probability distribution with probability of success  $p = 1 - e^{-\frac{1}{\theta}}$ , and it is shifted with 1 to the left since it starts with  $k = 0$ . Therefore the photon number expectation value is

$$\mu_K = \frac{1}{p} - 1 = \frac{1}{1 - e^{-\frac{1}{\theta}}} - 1 = \frac{1}{e^{\frac{1}{\theta}} - 1}$$

and its variance is

$$\sigma_K^2 = \frac{1-p}{p^2} = \frac{e^{-\frac{1}{\theta}}}{(1 - e^{-\frac{1}{\theta}})^2} = \mu_K^2 + \mu_K.$$

- (c) The joint probability distribution for the  $k_i$  is

$$p_K(k_1, \dots, k_n; \theta) = \frac{1}{Z(\theta)^n} \prod_{i=1}^n e^{-k_i/\theta} = \frac{1}{Z(\theta)^n} e^{-\frac{1}{\theta} \sum_{i=1}^n k_i}.$$

The log likelihood is  $-n \cdot \log Z(\theta) - 1/\theta \sum_{i=1}^n k_i$ .

We find the maxima of the log likelihood by setting the derivative with respect to the parameter  $\theta$  to zero:

$$\frac{d}{d\theta} \log p_K(k_1, \dots, k_n; \theta) = -n \cdot \frac{e^{-\frac{1}{\theta}}}{\theta^2(1 - e^{-\frac{1}{\theta}})} + \frac{1}{\theta^2} \sum_{i=1}^n k_i = 0$$

or

$$\frac{1}{e^{\frac{1}{\theta}} - 1} = \frac{1}{n} \sum_{i=1}^n k_i = s_n.$$

For a hot body,  $\theta \gg 1$  and  $\frac{1}{e^{\frac{1}{\theta}} - 1} \approx \theta$ , we obtain

$$\theta \approx \frac{1}{n} \sum_{i=1}^n k_i = s_n.$$

Thus the maximum likelihood estimator  $\hat{\Theta}_n$  for the temperature is given in this limit by the sample mean of the photon number

$$\hat{\Theta}_n = \frac{1}{n} \sum_{i=1}^n K_i.$$

- (d) According to the central limit theorem, the sample mean approaches for large  $n$  a Gaussian distribution with standard deviation our root mean square error

$$\sigma_{\hat{\Theta}_n} = \frac{\sigma_K}{\sqrt{n}}.$$

To allow only for 1% relative root mean square error in the temperature, we need  $\frac{\sigma_K}{\sqrt{n}} < 0.01\mu_K$ . With  $\sigma_K^2 = \mu_K^2 + \mu_K$  it follows that

$$\sqrt{n} > \frac{\sigma_K}{0.01\mu_K} = 100 \frac{\sqrt{\mu_K^2 + \mu_K}}{\mu_K} = 100 \sqrt{1 + \frac{1}{\mu_K}}.$$

In general, for large temperatures, i.e. large mean photon numbers  $\mu_K \gg 1$ , we need about 10,000 samples.

- (e) The 95% confidence interval for the temperature estimate for the situation in part (d), i.e.

$$\sigma_{\hat{\Theta}_n} = \frac{\sigma_K}{\sqrt{n}} = 0.01\mu_K,$$

is

$$[\hat{K} - 1.96\sigma_{\hat{K}}, \hat{K} + 1.96\sigma_{\hat{K}}] = [\hat{K} - 0.0196\mu_K, \hat{K} + 0.0196\mu_K].$$



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

---

5. (a) The sample mean estimator  $\hat{\Theta}_n = \frac{W_1 + \dots + W_n}{n}$  in this case is

$$\hat{\Theta}_{1000} = \frac{2340}{1000} = 2.34.$$

From the standard normal table, we have  $\Phi(1.96) = 0.975$ , so we obtain

$$\mathbf{P} \left( \frac{|\hat{\Theta}_{1000} - \mu|}{\sqrt{\text{var}(W_i)/1000}} \leq 1.96 \right) \approx 0.95.$$

Because the variance is less than 4, we have

$$\mathbf{P} \left( \hat{\Theta}_{1000} - \mu \leq 1.96 \sqrt{\text{var}(W_i)/1000} \right) \leq \mathbf{P} \left( \hat{\Theta}_{1000} - \mu \leq 1.96 \sqrt{4/1000} \right),$$

and letting the right-hand side of the above equation  $\approx 0.95$  gives a 95% confidence, i.e.,

$$\left[ \hat{\Theta}_{1000} - 1.96 \sqrt{4/1000}, \hat{\Theta}_{1000} + 1.96 \sqrt{4/1000} \right] = \left[ \hat{\Theta}_{1000} - 0.124, \hat{\Theta}_{1000} + 0.124 \right] = [2.216, 2.464]$$

- (b) The likelihood function is

$$f_W(w; \theta) = \prod_{i=1}^n f_{W_i}(w_i; \theta) = \prod_{i=1}^n \theta e^{-\theta w_i},$$

And the log-likelihood function is

$$\log f_W(w; \theta) = n \log \theta - \theta \sum_{i=1}^n w_i,$$

The derivative with respect to  $\theta$  is  $\frac{n}{\theta} - \sum_{i=1}^n w_i$ , and by setting it to zero, we see that the maximum of  $\log f_W(w; \theta)$  over  $\theta \geq 0$  is attained at  $\hat{\theta}_n = \frac{n}{\sum_{i=1}^n w_i}$ . The resulting estimator is

$$\hat{\Theta}_n^{mle} = \frac{n}{\sum_{i=1}^n W_i}.$$

In this case,

$$\hat{\Theta}_n^{mle} = \frac{1000}{2340} = 0.4274.$$

6. (a) Using the regression formulas of Section 9.2, we have

$$\hat{\theta}_1 = \frac{\sum_{i=1}^5 (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^5 (x_i - \bar{x})^2}, \quad \hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x},$$

where

$$\bar{x} = \frac{1}{5} \sum_{i=1}^5 x_i = 4.94, \quad \bar{y} = \frac{1}{5} \sum_{i=1}^5 y_i = 134.38.$$

The resulting ML estimates are

$$\hat{\theta}_1 = 40.53, \quad \hat{\theta}_0 = -65.86.$$

(b) Using the same procedure as in part (a), we obtain

$$\hat{\theta}_1 = \frac{\sum_{i=1}^5 (x_i^2 - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^5 (x_i^2 - \bar{x})^2}, \quad \hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x},$$

where

$$\bar{x} = \frac{1}{5} \sum_{i=1}^5 x_i^2 = 33.60, \quad \bar{y} = \frac{1}{5} \sum_{i=1}^5 y_i = 134.38.$$

which for the given data yields

$$\hat{\theta}_1 = 4.09, \quad \hat{\theta}_0 = -3.07.$$

Figure 1 shows the data points  $(x_i, y_i)$ ,  $i = 1, \dots, 5$ , the estimated linear model

$$y = 40.53x - 65.86,$$

and the estimated quadratic model

$$y = 4.09x^2 - 3.07.$$

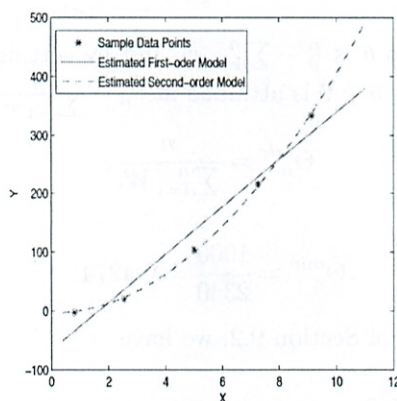


Figure 1: Regression Plot

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

---

Recitation 25  
December 9, 2010  
Based on Spring 10 Final exam

Question 2 An atom of the radioactive element Vestium decays to an atom of Hockfieldium after a time that is an exponential random variable with parameter  $\lambda$ . Hockfieldium is a stable element; i.e., it is not radioactive. Each radioactive atom decays independently of any other atoms.

(a) Suppose a box has  $n$  atoms of Vestium at time 0, where  $n$  is a positive integer. Let  $V$  be the remaining atoms of Vestium in the box at time  $t$ , where  $t$  is a positive real number. Find the PMF of  $V$ .

(b) An atom of Vestium can itself be the product of the radioactive decay of an atom of Grayon. The decay of any one atom of Grayon to an atom of Vestium occurs after a time that is an exponential random variable with parameter  $\mu$ .

Suppose a box initially contains two atoms of Grayon and nothing else. Find the expected time until the box is no longer radioactive, i.e., it contains neither Grayon nor Vestium—only Hockfieldium.

Question 4 **Breaking a stick more than twice.** We start with a stick of length  $\ell$ . We break at a point which is chosen randomly and uniformly over its length, and keep the piece that contains the left end of the stick. We then repeat the same process several times on the piece that we were left with. Denote by  $X_n$  the length of the piece we are left with after breaking  $n$  times.

(a) Find  $E[X_n]$ .

(b) After breaking the stick  $n$  times, we randomly pick one of the  $n + 1$  pieces, each of the pieces being equally likely to be picked. Calculate the expected length of the chosen piece.

(c) Does the sequence  $X_1, X_2, \dots$  converge in probability to a number? If so, to what value? Prove.

Question 5 Let  $W_1, W_2$ , and  $W_3$  be independent, continuous random variables each uniformly distributed over  $[0, 1]$ . Let  $X = W_1 + W_2$  and  $Y = X + W_3$ .

(a) Find the linear least mean squares (LLMS) estimator of  $X$  from  $Y$ .

(b) Find the maximum a posteriori probability (MAP) estimator of  $X$  from  $Y$ .





## Recitation 25

12/9

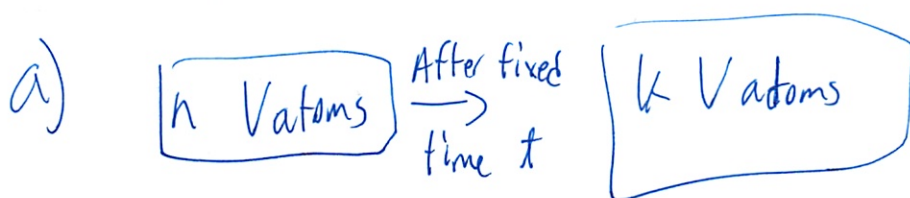
- focus on 2nd half of course

- this recitation samples everything except Markov

2. Grayson  $\rightarrow$  Vestim  $\rightarrow$  Hochfeldm

G atom  $\rightarrow$  V atom after decay time  $\sim \exp(-\lambda t)$

~~V~~ V  $\rightarrow$  H  $\sim \exp(-\lambda t)$



Want PDF of  $k$

- each one decays int.

- if it does not decay = success

$P_n(k) \sim \text{Binomial}(p)$

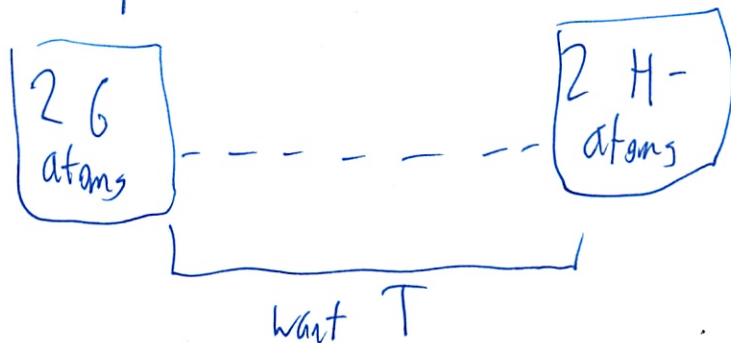
$$P(\text{No decay by time } t) = 1 - e^{-\lambda t}$$



$$= \binom{n}{k} (e^{-\lambda t})^k (1 - e^{-\lambda t})^{n-k} \quad k = 0, 1, \dots, n$$

②

b) More poisson like

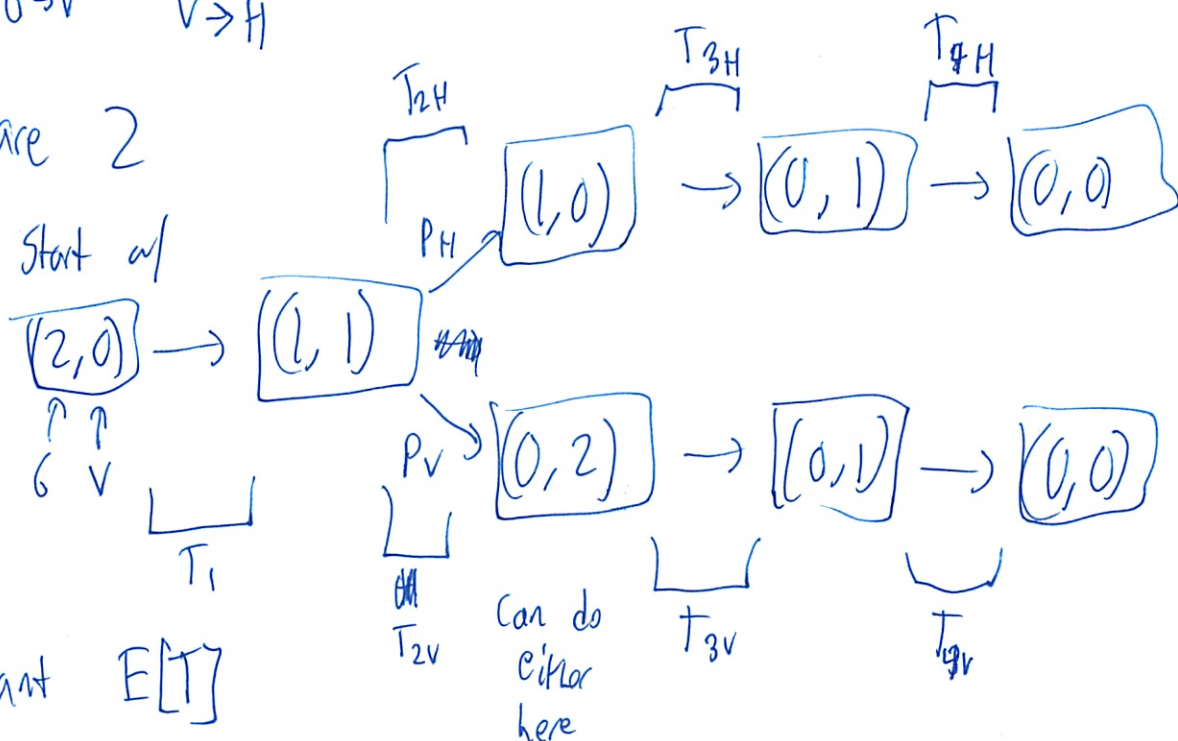


If only 1

$$\frac{1}{\mu_{G \rightarrow V}} + \frac{1}{\lambda_{V \rightarrow H}}$$

But are 2

Start w/



Want  $E[T]$

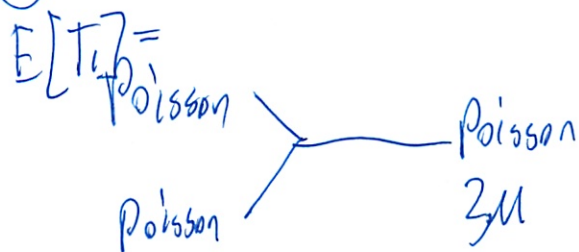
- calculate T upper  
lower

- then expected value theorem

$$E[T] = P_H \cdot E[\text{time of upper band}] + P_V \cdot E[\text{time of lower band}]$$



③



What is time of 1st arrival in  $2\mu$  Poisson merged process?

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

$$E[T_{2H}] = \frac{1}{\lambda}$$

$$E[T_{3H}] = \frac{1}{\mu} \quad \text{time for } G \rightarrow V$$

$$E[T_{4H}] = \frac{1}{\lambda} \quad \text{time } V \rightarrow H$$

$$E[T_{2V}] = \frac{1}{\mu}$$

$$E[T_{3V}] = \frac{1}{2\lambda} \quad \leftarrow \text{from merged } 2\lambda \text{ process}$$

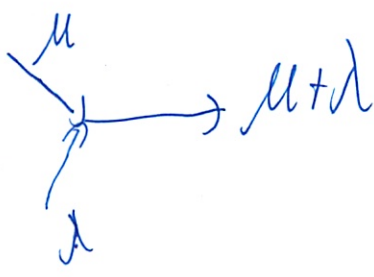
$$E[T_{4V}] = \frac{1}{\lambda} \quad \leftarrow \text{either one can leave}$$

$$E[T] = P_H \cdot \left( \frac{1}{2\mu} + \frac{1}{\lambda} + \frac{1}{\mu} + \frac{1}{\lambda} \right) + P_V \cdot \left( \frac{1}{2\mu} + \frac{1}{\mu} + \frac{1}{2\lambda} + \frac{1}{\lambda} \right)$$

Note the two branches are different  
bottom branch is faster

(This all seems so easy - can I think of it?)

4)



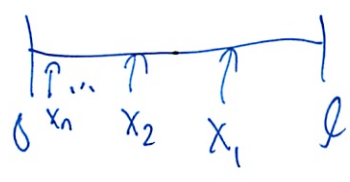
$$P_H = \frac{\lambda}{\mu + \lambda}$$

$$P_V = \frac{\mu}{\mu + \lambda}$$

$$E[T] = \frac{\lambda}{\mu + \lambda} \cdot \left( \frac{1}{2\mu} + \frac{1}{\lambda} + \frac{1}{\lambda} + \frac{1}{\lambda} \right) + \frac{\mu}{\mu + \lambda} \left( \frac{1}{2\mu} + \frac{1}{\mu} + \frac{1}{2\lambda} + \frac{1}{\lambda} \right)$$

4. Stick breaking problem

(try to get intuition to do it - don't just memorize example)



$X_n$  = length of left most piece after  $n$  breaks

$$E[X_n] = ?$$

$$E[X_1] = \frac{l}{2} \quad \text{c since break according to uniform dist}$$

$$E[X_2] = \text{conditioning + law of iterated expectation}$$

$$= E[E[X_n | X_{n-1}]]$$

$$= E[\frac{X_{n-1}}{2}]$$

5)

$$= E[X_{n-1}]$$

$$E[X_n] = \frac{2}{2^n} \text{ generalize}$$

$$= \frac{l}{2^n}$$

b) After  $n$  breaks, of lengths  $L_0, L_1, \dots, L_n$   
have  $n+1$  pieces

$M_i$ : length of randomly chosen piece

↳ prob  $\frac{1}{n+1}$

$$E[M] = \sum_{i=0}^n P(\text{ith piece chosen}) E[M | \text{ith piece chosen}]$$

$$\quad \quad \quad \uparrow \quad \quad \quad \uparrow$$

$$\quad \quad \quad \frac{1}{n+1} \quad \quad \quad E[L_i]$$

$$= \frac{1}{n+1} \sum_{i=0}^n E[L_i]$$

$$= \frac{1}{n+1} E\left[\sum_{i=0}^n L_i\right]$$

$$= \frac{l}{n+1} \text{ } \leftarrow \text{all of the pieces sum up to } l$$



6

c) As  ~~$x_n$~~  get <sup>larger,  $L_n$  gets</sup> smaller + smaller + converges to 0

Does  $\{X_n\}$  converge in prob and to what limit?  
- to 0

Does it converge?

- fancy math

- but straightforward - only 1 path for proof

For any  $\epsilon$ , find  $P(|X_n - 0| \geq \epsilon)$   
 $\rightarrow 0$

Show converges  $\rightarrow 0$  as  $n \rightarrow \infty$

Markov inequality

$$P(X_n \geq \epsilon) \leq \frac{E[X_n]}{\epsilon}$$

$$= \frac{1}{\epsilon \cdot 2^n}$$

$\rightarrow 0$  as  $n \rightarrow \infty$

~~Q.E.D.~~ Q.E.D.  
Proved  $\rightarrow 0$  as  $n \rightarrow \infty$  for any  $\epsilon > 0$

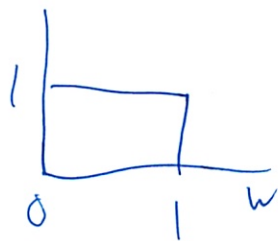
⑦

5. Inference problem

Bayesian - MAP, LMS

Classical - ML

Have 3 RV  $W_1, W_2, W_3 \rightarrow \text{iid}$



Uniformly distributed

$$X = W_1 + W_2$$

$$Y = X + W_3 \quad \hookrightarrow X, Y \text{ correlated}$$

If know  $Y$ , can make an inference about  $X$

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

Find these values + plug in

$$\begin{aligned} E[X] &= E[W_1] + E[W_2] \\ &= \frac{1}{2} + \frac{1}{2} \\ &= 1 \end{aligned}$$

$$\begin{aligned} E[Y] &= E[X] + E[W_3] \\ &= 1 + \frac{1}{2} \\ &= \frac{3}{2} \end{aligned}$$

ind. does not matter

(8)

$$\text{Var}(x) = \text{sum of two ind } \overset{\text{RV}}{\cancel{\text{var}}} = \text{sum of var}$$

$$= \cancel{\text{Var}}(w_1) + \text{Var}(w_2)$$

$$= \frac{1}{12} + \frac{1}{12}$$

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

$$\text{Var}(y) = \text{var}(x) + \text{var}(w_3)$$

$$= \frac{1}{6} + \frac{1}{12}$$

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

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

↑  
need cross  
2nd moment

$$E[XY] = E[X(X + w_3)]$$

$$= \text{both ind so } E[XY] = E[X] E[Y]$$

$$= E[X^2] + E[X] E[w_3]$$

$$\cancel{E[XY]} = \text{Var}(x) + (E[X])^2 + E[X] E[w_3]$$

$$= \frac{1}{6} + 1^2 + 1 \cdot \frac{1}{2}$$

$$= \frac{5}{3}$$



9)

$$\text{cov}(X, Y) = \frac{5}{3} - 1 \cdot \frac{2}{3}$$

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

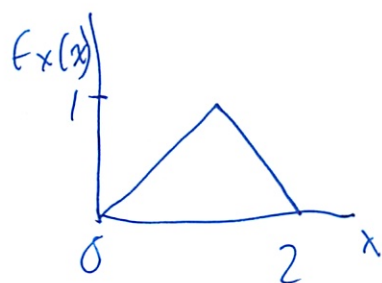
$$\hat{X}_{\text{LMS}} = 1 + \frac{1/6}{1/4} \left( Y - \frac{3}{2} \right)$$

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

$$= \frac{2}{3} Y$$

b) MAP estimator of  $X$  based on  $Y$

$f_X(x)$  = triangular density through lengthy convolution process



~~MAP is max posterior of here =  $\hat{X}_{\text{MAP}}$~~

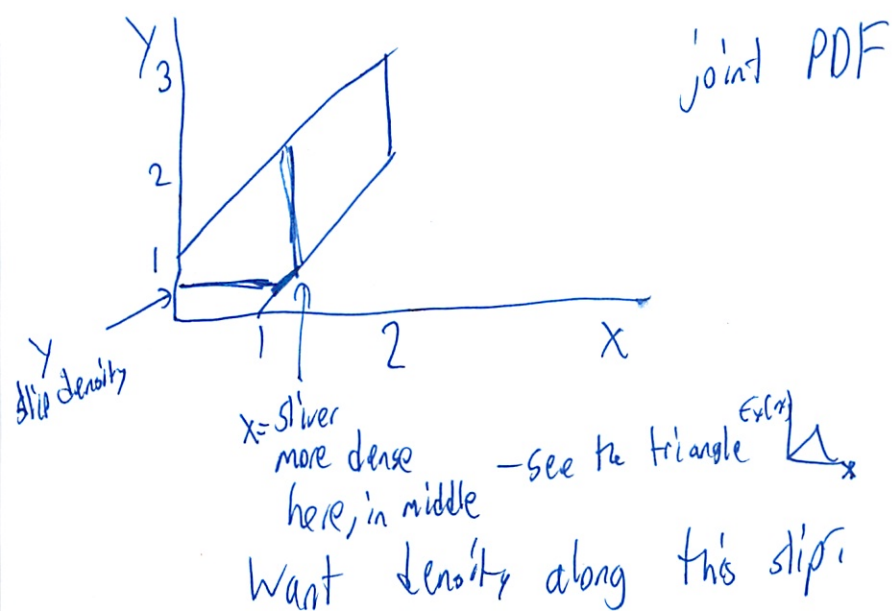
$$\text{MAP} \rightarrow \max_x f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \quad \leftarrow \text{can calculate easily}$$

denominator does not matter

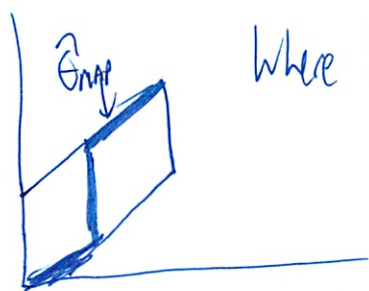
$$= \max_x f_{X,Y}(x,y)$$

(10)

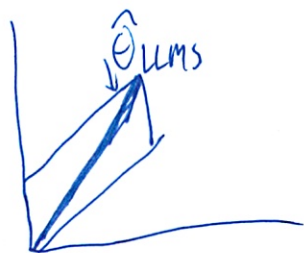
$$Y = X + \underbrace{[0,1]}_{\text{Uniform}}$$



$$\hat{X}_{MAP} = \begin{cases} 1 & \text{if } 1 \leq y \leq 2 \\ y & \text{if } 0 \leq y \leq 1 \\ y-1 & \text{if } 2 \leq y \leq 3 \end{cases}$$



Where is joint maximized?



Conditional expectation estimator

— integrate along slice —

②

Must do calculation

In this case same as LMS  
    <sup>↑</sup> not true in general



③

## 9.5 Summary

Have not discussed

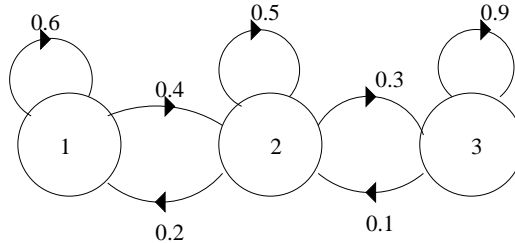
- ~~the~~ estimation in time-varying cov
- non parametric estimation
- further developments in linear + non-linear regression
- methods for designing statistical experiments
- methods for validating conclusions of stat. study
- computational methods
- etc

end of book!

## 6.041/6.431 Fall 2010 Final Exam Solutions

### Wednesday, December 15, 9:00AM - 12:00noon.

**Problem 1. (32 points)** Consider a Markov chain  $\{X_n; n = 0, 1, \dots\}$ , specified by the following transition diagram.



1. **(4 points)** Given that the chain starts with  $X_0 = 1$ , find the probability that  $X_2 = 2$ .

**Solution:** The two-step transition probability is:

$$\begin{aligned}
 r_{12}(2) &= p_{11} \cdot p_{12} + p_{12} \cdot p_{22} \\
 &= 0.6 \cdot 0.4 + 0.4 \cdot 0.5 \\
 &= 0.44.
 \end{aligned}$$

2. **(4 points)** Find the steady-state probabilities  $\pi_1, \pi_2, \pi_3$  of the different states.

**Solution:** We set up the balance equations of a birth-death process and the normalization equation as such:

$$\begin{aligned}
 \pi_1 p_{12} &= \pi_2 p_{21} \\
 \pi_2 p_{23} &= \pi_3 p_{32} \\
 \pi_1 + \pi_2 + \pi_3 &= 1.
 \end{aligned}$$

Solving the system of equations yields the following steady state probabilities:

$$\begin{aligned}
 \pi_1 &= 1/9 \\
 \pi_2 &= 2/9 \\
 \pi_3 &= 6/9.
 \end{aligned}$$

*In case you did not do part 2 correctly, in **all** subsequent parts of this problem you can just use the symbols  $\pi_i$ : you do not need to plug in actual numbers.*

3. **(4 points)** Let  $Y_n = X_n - X_{n-1}$ . Thus,  $Y_n = 1$  indicates that the  $n$ th transition was to the right,  $Y_n = 0$  indicates it was a self-transition, and  $Y_n = -1$  indicates it was a transition to the left. Find  $\lim_{n \rightarrow \infty} \mathbf{P}(Y_n = 1)$ .

**Solution:** Using the total probability theorem and steady state probabilities,

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \mathbf{P}(Y_n = 1) &= \sum_{i=1}^3 \pi_i \cdot \mathbf{P}(Y_n = 1 \mid X_{n-1} = i) \\
 &= \pi_1 p_{12} + \pi_2 p_{23} \\
 &= 1/9.
 \end{aligned}$$

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

---

4. (4 points) Is the sequence  $Y_n$  a Markov chain? Justify your answer.

**Solution:** No. Assume the Markov process is in steady state. To satisfy the Markov property,

$$\mathbf{P}(Y_n = 1 \mid Y_{n-1} = 1, Y_{n-2} = 1) = \mathbf{P}(Y_n = 1 \mid Y_{n-1} = 1).$$

For large  $n$ ,

$$\mathbf{P}(Y_n = 1 \mid Y_{n-1} = 1, Y_{n-2} = 1) = 0,$$

since it is not possible to move upwards 3 times in a row. However in steady state,

$$\begin{aligned} \mathbf{P}(Y_n = 1 \mid Y_{n-1} = 1) &= \frac{\mathbf{P}(\{Y_n = 1\} \cap \{Y_{n-1} = 1\})}{\mathbf{P}(Y_{n-1} = 1)} \\ &= \frac{\pi_1 p_{12} p_{23}}{\pi_1 p_{12} + \pi_2 p_{23}} \\ &\neq 0. \end{aligned}$$

Therefore, the sequence  $Y_n$  is not a Markov chain.

5. (4 points) Given that the  $n$ th transition was a transition to the right ( $Y_n = 1$ ), find the probability that the previous state was state 1. (You can assume that  $n$  is large.)

**Solution:** Using Bayes' Rule,

$$\begin{aligned} \mathbf{P}(X_{n-1} = 1 \mid Y_n = 1) &= \frac{\mathbf{P}(X_{n-1} = 1) \mathbf{P}(Y_n = 1 \mid X_{n-1} = 1)}{\sum_{i=1}^3 \mathbf{P}(X_{n-1} = i) \mathbf{P}(Y_n = 1 \mid X_{n-1} = i)} \\ &= \frac{\pi_1 p_{12}}{\pi_1 p_{12} + \pi_2 p_{23}} \\ &= 2/5. \end{aligned}$$

6. (4 points) Suppose that  $X_0 = 1$ . Let  $T$  be defined as the first *positive time* at which the state is again equal to 1. Show how to find  $\mathbf{E}[T]$ . (It is enough to write down whatever equation(s) needs to be solved; you do not have to actually solve it/them or to produce a numerical answer.)

**Solution:** In order to find the mean recurrence time of state 1, the mean first passage times to state 1 are first calculated by solving the following system of equations:

$$\begin{aligned} t_2 &= 1 + p_{22}t_2 + p_{23}t_3 \\ t_3 &= 1 + p_{32}t_2 + p_{33}t_3. \end{aligned}$$

The mean recurrence time of state 1 is then  $t_1^* = 1 + p_{12}t_2$ .

Solving the system of equations yields  $t_2 = 20$  and  $t_3 = 30$  and  $t_1^* = 9$ .

7. (4 points) Does the sequence  $X_1, X_2, X_3, \dots$  converge in probability? If yes, to what? If not, just say "no" without explanation.

**Solution:** No.



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

---

8. (4 points) Let  $Z_n = \max\{X_1, \dots, X_n\}$ . Does the sequence  $Z_1, Z_2, Z_3, \dots$  converge in probability? If yes, to what? If not, just say “no” without explanation.

**Solution:** Yes. The sequence converges to 3 in probability.

For the original Markov chain, states  $\{1, 2, 3\}$  form one single recurrent class. Therefore, the Markov process will eventually visit each state with probability 1. In this case, the sequence  $Z_n$  will, with probability 1, converge to 3 once  $X_n$  visits 3 for the first time.

**Problem 2. (68 points)** Alice shows up at an Athena cluster at time zero and spends her time exclusively in typing emails. The times that her emails are sent are a Poisson process with rate  $\lambda_A$  per hour.

1. (3 points) What is the probability that Alice sent exactly three emails during the time interval  $[1, 2]$ ?

**Solution:** The number of emails Alice sends in the interval  $[1, 2]$  is a Poisson random variable with parameter  $\lambda_A$ . So we have:

$$\mathbf{P}(3, 1) = \frac{\lambda_A^3 e^{-\lambda_A}}{3!}.$$

2. Let  $Y_1$  and  $Y_2$  be the times at which Alice's first and second emails were sent.

- (a) (3 points) Find  $\mathbf{E}[Y_2 | Y_1]$ .

**Solution:** Define  $T_2$  as the second inter-arrival time in Alice's Poisson process. Then:

$$Y_2 = Y_1 + T_2$$

$$\mathbf{E}[Y_2 | Y_1] = \mathbf{E}[Y_1 + T_2 | Y_1] = Y_1 + \mathbf{E}[T_2] = Y_1 + 1/\lambda_A.$$

- (b) (3 points) Find the PDF of  $Y_1^2$ .

**Solution:** Let  $Z = Y_1^2$ . Then we first find the CDF of  $Z$  and differentiate to find the PDF of  $Z$ :

$$F_Z(z) = \mathbf{P}(Y_1^2 \leq z) = \mathbf{P}(-\sqrt{z} \leq Y_1 \leq \sqrt{z}) = \begin{cases} 1 - e^{-\lambda_A \sqrt{z}} & z \geq 0 \\ 0 & z < 0. \end{cases}$$

$$\begin{aligned} f_Z(z) = \frac{dF_Z(z)}{dz} &= \lambda_A e^{-\lambda_A \sqrt{z}} \left( \frac{1}{2} z^{-1/2} \right) & (z \geq 0) \\ f_Z(z) &= \begin{cases} \frac{\lambda_A}{2\sqrt{z}} e^{-\lambda_A \sqrt{z}} & z \geq 0 \\ 0 & z < 0. \end{cases} \end{aligned}$$

- (c) (3 points) Find the joint PDF of  $Y_1$  and  $Y_2$ .

**Solution:**

$$\begin{aligned} f_{Y_1, Y_2}(y_1, y_2) &= f_{Y_1}(y_1) f_{Y_2|Y_1}(y_2|y_1) \\ &= f_{Y_1}(y_1) f_{T_2}(y_2 - y_1) \\ &= \lambda_A e^{-\lambda_A y_1} \lambda_A e^{-\lambda_A (y_2 - y_1)} & y_2 \geq y_1 \geq 0 \\ &= \begin{cases} \lambda_A^2 e^{-\lambda_A y_2} & y_2 \geq y_1 \geq 0 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

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

---

3. You show up at time 1 and you are told that Alice has sent exactly one email so far. (Only give answers here, no need to justify them.)

- (a) **(3 points)** What is the conditional expectation of  $Y_2$  given this information?

**Solution:** Let  $A$  be the event {exactly one arrival in the interval  $[0,1]$ }. Looking forward from time  $t = 1$ , the time until the next arrival is simply an exponential random variable ( $T$ ). So,

$$\mathbf{E}[Y_2 \mid A] = 1 + \mathbf{E}[T] = 1 + 1/\lambda_A.$$

- (b) **(3 points)** What is the conditional expectation of  $Y_1$  given this information?

**Solution:** Given  $A$ , the times in this interval are equally likely for the arrival  $Y_1$ . Thus,

$$\mathbf{E}[Y_1 \mid A] = 1/2.$$

4. Bob just finished exercising (without email access) and sits next to Alice at time 1. He starts typing emails at time 1, and fires them according to an independent Poisson process with rate  $\lambda_B$ .

- (a) **(5 points)** What is the PMF of the total number of emails sent by the two of them together during the interval  $[0, 2]$ ?

**Solution:** Let  $K$  be the total number of emails sent in  $[0, 2]$ . Let  $K_1$  be the total number of emails sent in  $[0, 1]$ , and let  $K_2$  be the total number of emails sent in  $[1, 2]$ . Then  $K = K_1 + K_2$  where  $K_1$  is a Poisson random variable with parameter  $\lambda_A$  and  $K_2$  is a Poisson random variable with parameter  $\lambda_A + \lambda_B$  (since the emails sent by both Alice and Bob after time  $t = 1$  arrive according to the merged Poisson process of Alice's emails and Bob's emails). Since  $K$  is the sum of independent Poisson random variables,  $K$  is a Poisson random variable with parameter  $2\lambda_A + \lambda_B$ . So  $K$  has the distribution:

$$p_K(k) = \frac{(2\lambda_A + \lambda_B)^k e^{-(2\lambda_A + \lambda_B)}}{k!} \quad k = 0, 1, \dots$$

- (b) **(5 points)** What is the expected value of the total typing time associated with the email that Alice is typing at the time that Bob shows up? (Here, "total typing time" includes the time that Alice spent on that email both before and after Bob's arrival.)

**Solution:** The total typing time  $Q$  associated with the email that Alice is typing at the time Bob shows up is the sum of  $S_0$ , the length of time between Alice's last email or time 0 (whichever is later) and time 1, and  $T_1$ , the length of time from 1 to the time at which Alice sends her current email.  $T_1$  is exponential with parameter  $\lambda_A$ . and  $S_0 = \min\{T_0, 1\}$ , where  $T_0$  is exponential with parameter  $\lambda_A$ .

Then,

$$Q = S_0 + T_1 = \min\{T_0, 1\} + T_1$$

and

$$\mathbf{E}[Q] = \mathbf{E}[S_0] + \mathbf{E}[T_1].$$

We have:  $\mathbf{E}[T_1] = 1/\lambda_A$ .

We can find  $\mathbf{E}[S_0]$  via the law of total expectations:

$$\begin{aligned}
 \mathbf{E}[S_0] = \mathbf{E}[\min\{T_0, 1\}] &= \mathbf{P}(T_0 \leq 1)\mathbf{E}[T_0 \mid T_0 \leq 1] + \mathbf{P}(T_0 > 1)\mathbf{E}[1 \mid T_0 > 1] \\
 &= (1 - e^{-\lambda_A}) \int_0^1 t f_{T|T_0 \leq 1}(t) dt + e^{-\lambda_A} \\
 &= (1 - e^{-\lambda_A}) \int_0^1 t \frac{\lambda_A e^{-\lambda_A t}}{(1 - e^{-\lambda_A})} dt + e^{-\lambda_A} \\
 &= \int_0^1 t \lambda_A e^{-\lambda_A t} dt + e^{-\lambda_A} \\
 &= \frac{1}{\lambda_A} \int_0^1 t \lambda_A^2 e^{-\lambda_A t} dt + e^{-\lambda_A} \\
 &= \frac{1}{\lambda_A} (1 - e^{-\lambda_A} - \lambda_A e^{-\lambda_A}) + e^{-\lambda_A} \\
 &= \frac{1}{\lambda_A} (1 - e^{-\lambda_A})
 \end{aligned}$$

where the above integral is evaluated by manipulating the integrand into an Erlang order 2 PDF and equating the integral of this PDF from 0 to 1 to the probability that there are 2 or more arrivals in the first hour (i.e.  $\mathbf{P}(Y_2 < 1) = 1 - \mathbf{P}(0, 1) - \mathbf{P}(1, 1)$ ). Alternatively, one can integrate by parts and arrive at the same result.

Combining the above expectations:

$$\mathbf{E}[Q] = \mathbf{E}[S_0] + \mathbf{E}[T_1] = \frac{1}{\lambda_A} (1 - e^{-\lambda_A}) + \frac{1}{\lambda_A} = \frac{1}{\lambda_A} (2 - e^{-\lambda_A}).$$

- (c) **(5 points)** What is the expected value of the time until each one of them has sent at least one email? (Note that we count time starting from time 0, and we take into account any emails possibly sent out by Alice during the interval  $[0, 1]$ .)

**Solution:** Define  $U$  as the time from  $t = 0$  until each person has sent at least one email.

Define  $V$  as the remaining time from when Bob arrives (time 1) until each person has sent at least one email (so  $V = U - 1$ ).

Define  $S$  as the time until Bob sends his first email after time 1.

Define the event  $A = \{\text{Alice sends one or more emails in the time interval } [0, 1]\} = \{Y_1 \leq 1\}$ , where  $Y_1$  is the time Alice sends her first email.

Define the event  $B = \{\text{After time 1, Bob sends his next email before Alice does}\}$ , which is equivalent to the event where the next arrival in the merged process from Alice and Bob's original processes (starting from time 1) comes from Bob's process.

We have:

$$\mathbf{P}(A) = \mathbf{P}(Y_1 \leq 1) = 1 - e^{-\lambda_A}$$

$$\mathbf{P}(B) = \frac{\lambda_B}{\lambda_A + \lambda_B}.$$



Then,

$$\begin{aligned}
 \mathbf{E}[U] &= \mathbf{P}(A)\mathbf{E}[U | A] + \mathbf{P}(A^c)\mathbf{E}[U | A^c] \\
 &= (1 - e^{-\lambda_A})(1 + \mathbf{E}[V | A]) + e^{-\lambda_A}(1 + \mathbf{E}[V | A^c]) \\
 &= (1 - e^{-\lambda_A})(1 + \mathbf{E}[V | A]) + e^{-\lambda_A}(1 + \mathbf{P}(B | A^c)\mathbf{E}[V | B \cap A^c] + \mathbf{P}(B^c | A^c)\mathbf{E}[V | B^c \cap A^c]) \\
 &= (1 - e^{-\lambda_A})(1 + \mathbf{E}[V | A]) + e^{-\lambda_A}(1 + \mathbf{P}(B)\mathbf{E}[V | B \cap A^c] + \mathbf{P}(B^c)\mathbf{E}[V | B^c \cap A^c]) \\
 &= (1 - e^{-\lambda_A})(1 + \mathbf{E}[V | A]) + e^{-\lambda_A} \left( 1 + \frac{\lambda_B}{\lambda_A + \lambda_B} \mathbf{E}[V | B \cap A^c] + \frac{\lambda_A}{\lambda_A + \lambda_B} \mathbf{E}[V | B^c \cap A^c] \right).
 \end{aligned}$$

Note that  $\mathbf{E}[V | B^c \cap A^c]$  is the expected value of the time until each of them sends one email after time 1 (since, given  $A^c$ , Alice did not send any in the interval  $[0, 1]$ ) and given Alice sends an email before Bob. Then this is the expected time until an arrival in the merged process followed by the expected time until an arrival in Bob's process. So,  $\mathbf{E}[V | B^c \cap A^c] = \frac{1}{\lambda_A + \lambda_B} + \frac{1}{\lambda_B}$ .

Similarly,  $\mathbf{E}[V | B \cap A^c]$  is the time until each sends an email after time 1, given Bob sends an email before Alice. So  $\mathbf{E}[V | B \cap A^c] = \frac{1}{\lambda_A + \lambda_B} + \frac{1}{\lambda_A}$ .

Also,  $\mathbf{E}[V | A]$  is the expected time it takes for Bob to send his first email after time 1 (since, given  $A$ , Alice already sent an email in the interval  $[0, 1]$ ). So  $\mathbf{E}[V | A] = \mathbf{E}[S] = 1/\lambda_B$ . Combining all of this with the above, we have:

$$\begin{aligned}
 \mathbf{E}[U] &= (1 - e^{-\lambda_A})(1 + 1/\lambda_B) \\
 &\quad + e^{-\lambda_A} \left( 1 + \frac{\lambda_B}{\lambda_A + \lambda_B} \left( \frac{1}{\lambda_A + \lambda_B} + \frac{1}{\lambda_A} \right) + \frac{\lambda_A}{\lambda_A + \lambda_B} \left( \frac{1}{\lambda_A + \lambda_B} + \frac{1}{\lambda_B} \right) \right).
 \end{aligned}$$

- (d) **(5 points)** Given that a total of 10 emails were sent during the interval  $[0, 2]$ , what is the probability that exactly 4 of them were sent by Alice?

**Solution:**

$$\begin{aligned}
 \mathbf{P}(\text{Alice sent 4 in } [0, 2] \mid \text{total 10 sent in } [0, 2]) &= \frac{\mathbf{P}(\text{Alice sent 4 in } [0, 2] \cap \text{total 10 sent in } [0, 2])}{\mathbf{P}(\text{total 10 sent in } [0, 2])} \\
 &= \frac{\mathbf{P}(\text{Alice sent 4 in } [0, 2] \cap \text{Bob sent 6 in } [0, 2])}{\mathbf{P}(\text{total 10 sent in } [0, 2])} \\
 &= \frac{\left( \frac{(2\lambda_A)^4 e^{-2\lambda_A}}{4!} \right) \left( \frac{(\lambda_B)^6 e^{-\lambda_B}}{6!} \right)}{\frac{(2\lambda_A + \lambda_B)^{10} e^{-2\lambda_A - \lambda_B}}{10!}} \\
 &= \binom{10}{4} \left( \frac{2\lambda_A}{2\lambda_A + \lambda_B} \right)^4 \left( \frac{\lambda_B}{2\lambda_A + \lambda_B} \right)^6.
 \end{aligned}$$

As the form of the solution suggests, the problem can be solved alternatively by computing the probability of a single email being sent by Alice, given it was sent in the interval  $[0, 2]$ . This can be found by viewing the number of emails sent by Alice in  $[0, 2]$  as the number of arrivals arising from a Poisson process with twice the rate ( $2\lambda_A$ ) in an interval of half the duration (particularly, the interval  $[1, 2]$ ), then merging this process with Bob's process. Then the probability that an email sent in the interval  $[0, 2]$  was sent by Alice is the probability that an arrival in this new merged process came from the newly constructed  $2\lambda_A$  rate process:

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

---

$$p = \frac{2\lambda_A}{2\lambda_A + \lambda_B}.$$

Then, out of 10 emails, the probability that 4 came from Alice is simply a binomial probability with 4 successes in 10 trials, which agrees with the solution above.

5. **(5 points)** Suppose that  $\lambda_A = 4$ . Use Chebyshev's inequality to find an upper bound on the probability that Alice sent at least 5 emails during the time interval  $[0, 1]$ . Does the Markov inequality provide a better bound?

**Solution:**

Let  $N$  be the number of emails Alice sent in the interval  $[0, 1]$ . Since  $N$  is a Poisson random variable with parameter  $\lambda_A$ ,

$$\mathbf{E}[N] = \text{var}(N) = \lambda_A = 4.$$

To apply the Chebyshev inequality, we recognize:

$$\mathbf{P}(N \geq 5) = \mathbf{P}(N - 4 \geq 1) \leq \mathbf{P}(|N - 4| \geq 1) \leq \frac{\text{var}(N)}{1^2} = 4.$$

In this case, the upper-bound of 4 found by application of the Chebyshev inequality is uninformative, as we already knew  $\mathbf{P}(N \geq 5) \leq 1$ .

To find a better bound on this probability, use the Markov inequality, which gives:

$$\mathbf{P}(N \geq 5) \leq \frac{\mathbf{E}[N]}{5} = \frac{4}{5}.$$

6. **(5 points)** You do not know  $\lambda_A$  but you watch Alice for an hour and see that she sent exactly 5 emails. Derive the maximum likelihood estimate of  $\lambda_A$  based on this information.

**Solution:**

$$\begin{aligned}\hat{\lambda}_A &= \arg \max_{\lambda} \log(p_N(5; \lambda)) \\ &= \arg \max_{\lambda} \log\left(\frac{\lambda^5 e^{-\lambda}}{5!}\right) \\ &= \arg \max_{\lambda} -\log(5!) + 5\log(\lambda) - \lambda.\end{aligned}$$

Setting the first derivative to zero

$$\begin{aligned}\frac{5}{\lambda} - 1 &= 0 \\ \hat{\lambda}_A &= 5.\end{aligned}$$

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

---

7. **(5 points)** We have reasons to believe that  $\lambda_A$  is a large number. Let  $N$  be the number of emails sent during the interval  $[0, 1]$ . Justify why the CLT can be applied to  $N$ , and give a precise statement of the CLT in this case.

**Solution:** With  $\lambda_A$  large, we assume  $\lambda_A \gg 1$ . For simplicity, assume  $\lambda_A$  is an integer. We can divide the interval  $[0, 1]$  into  $\lambda_A$  disjoint intervals, each with duration  $1/\lambda_A$ , so that these intervals span the entire interval from  $[0, 1]$ . Let  $N_i$  be the number of arrivals in the  $i$ th such interval, so that the  $N_i$ 's are independent, identically distributed Poisson random variables with parameter 1. Since  $N$  is defined as the number of arrivals in the interval  $[0, 1]$ , then  $N = N_1 + \dots + N_{\lambda_A}$ . Since  $\lambda_A \gg 1$ , then  $N$  is the sum of a large number of independent and identically distributed random variables, where the distribution of  $N_i$  does not change as the number of terms in the sum increases. Hence,  $N$  is approximately normal with mean  $\lambda_A$  and variance  $\lambda_A$ .

If  $\lambda_A$  is not an integer, the same argument holds, except that instead of having  $\lambda_A$  intervals, we have an integer number of intervals equal to the integer part of  $\lambda_A$  ( $\bar{\lambda}_A = \text{floor}(\lambda_A)$ ) of length  $1/\lambda_A$  and an extra interval of a shorter length  $(\lambda_A - \bar{\lambda}_A)/\lambda_A$ .

Now,  $N$  is a sum of  $\lambda_A$  independent, identically distributed Poisson random variables with parameter 1 added to another Poisson random variable (also independent of all the other Poisson random variables) with parameter  $(\lambda_A - \bar{\lambda}_A)$ . In this case,  $N$  would need a small correction to apply the central limit theorem as we are familiar with it; however, it turns out that even without this correction, adding the extra Poisson random variable does not preclude the distribution of  $N$  from being approximately normal, for large  $\lambda_A$ , and the central limit theorem still applies.

To arrive at a precise statement of the CLT, we must “standardize”  $N$  by subtracting its mean then dividing by its standard deviation. After having done so, the CDF of the standardized version of  $N$  should converge to the standard normal CDF as the number of terms in the sum approaches infinity (as  $\lambda_A \rightarrow \infty$ ).

Therefore, the precise statement of the CLT when applied to  $N$  is:

$$\lim_{\lambda_A \rightarrow \infty} \mathbf{P} \left( \frac{N - \lambda_A}{\sqrt{\lambda_A}} \leq z \right) = \Phi(z)$$

where  $\Phi(z)$  is the standard normal CDF.

8. **(5 points)** Under the same assumption as in last part, that  $\lambda_A$  is large, you can now pretend that  $N$  is a normal random variable. Suppose that you observe the value of  $N$ . Give an (approximately) 95% confidence interval for  $\lambda_A$ . State precisely what approximations you are making.

*Possibly useful facts:* The cumulative normal distribution satisfies  $\Phi(1.645) = 0.95$  and  $\Phi(1.96) = 0.975$ .

**Solution:** We begin by estimating  $\lambda_A$  with its ML estimator  $\hat{\lambda}_A = N$ , where  $\mathbf{E}[N] = \lambda_A$ . With  $\lambda_A$  large, the CLT applies, and we can assume  $N$  has an approximately normal distribution. Since  $\text{var}(N) = \lambda_A$ , we can also approximate the variance of  $N$  with ML estimator for  $\lambda_A$ , so  $\text{var}(N) \approx N$ , and  $\sigma_N \approx \sqrt{N}$ .

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

---

To find the 95% confidence interval, we find  $\beta$  such that:

$$\begin{aligned} 0.95 &= \mathbf{P}(|N - \lambda_A| \leq \beta) \\ &= \mathbf{P}\left(\frac{|N - \lambda_A|}{\sqrt{N}} \leq \frac{\beta}{\sqrt{N}}\right) \\ &\approx 2\Phi\left(\frac{\beta}{\sqrt{N}}\right). \end{aligned}$$

So, we find:

$$\beta \approx \sqrt{N}\Phi^{-1}(0.975) = 1.96\sqrt{N}.$$

Thus, we can write:

$$\mathbf{P}(N - 1.96\sqrt{N} \leq \lambda_A \leq N + 1.96\sqrt{N}) \approx 0.95.$$

So, the approximate 95% confidence interval is:  $[N - 1.96\sqrt{N}, N + 1.96\sqrt{N}]$ .

9. You are now told that  $\lambda_A$  is actually the realized value of an exponential random variable  $\Lambda$ , with parameter 2:

$$f_\Lambda(\lambda) = 2e^{-2\lambda}, \quad \lambda \geq 0.$$

- (a) **(5 points)** Find  $\mathbf{E}[N^2]$ .

**Solution:**

$$\begin{aligned} \mathbf{E}[N^2] &= \mathbf{E}[\mathbf{E}[N^2 \mid \Lambda]] = \mathbf{E}[\text{var}(N \mid \Lambda) + (\mathbf{E}[N \mid \Lambda])^2] \\ &= \mathbf{E}[\Lambda + \Lambda^2] \\ &= \mathbf{E}[\Lambda] + \text{var}(\Lambda) + (\mathbf{E}[\Lambda])^2 \\ &= \frac{1}{2} + \frac{2}{2^2} \\ &= 1. \end{aligned}$$

- (b) **(5 points)** Find the linear least squares estimator of  $\Lambda$  given  $N$ .

**Solution:**

$$\hat{\Lambda}_{\text{LLMS}} = \mathbf{E}[\Lambda] + \frac{\text{cov}(N, \Lambda)}{\text{var}(N)}(N - \mathbf{E}[N]).$$

Solving for the above quantities:

$$\mathbf{E}[\Lambda] = \frac{1}{2}$$

$$\mathbf{E}[N] = \mathbf{E}[\mathbf{E}[N \mid \Lambda]] = \mathbf{E}[\Lambda] = \frac{1}{2}.$$

$$\text{var}(N) = \mathbf{E}[N^2] - (\mathbf{E}[N])^2 = 1 - \frac{1}{2^2} = \frac{3}{4}.$$

$$\text{cov}(N, \Lambda) = \mathbf{E}[N\Lambda] - \mathbf{E}[N]\mathbf{E}[\Lambda] = \mathbf{E}[\mathbf{E}[N\Lambda \mid \Lambda]] - (\mathbf{E}[\Lambda])^2 = \mathbf{E}[\Lambda^2] - (\mathbf{E}[\Lambda])^2 = \text{var}(\Lambda) = \frac{1}{4}.$$



Substituting these into the equation above:

$$\begin{aligned}\hat{\Lambda}_{\text{LLMS}} &= \mathbf{E}[\Lambda] + \frac{\text{cov}(N, \Lambda)}{\text{var}(N)}(N - \mathbf{E}[N]) \\ &= \frac{1}{2} + \frac{1/4}{3/4} \left( N - \frac{1}{2} \right) \\ &= \frac{1}{3}(N + 1) .\end{aligned}$$

FINAL EXAM STATISTICS

