

Measure? ← why?

$$\text{Eg: } P[2^{\text{nd}} \text{ year}] = 10\% = \frac{10}{|S|} \leftarrow \begin{array}{l} \# \text{ of second \\ year student} \\ \leftarrow \text{all people sitting} \\ \text{in the room} \end{array}$$

You are measuring something.

### Axioms

①  $P[A]$  has to be non-negative

②  $P[\emptyset] = 0$

③  $P[S] = 1$

$$\text{Ex } P[S] = P[\text{uni students}] = 1$$

$$\text{④ } P[\text{Additive}] = P[A_1 \cup A_2 \cup \dots] = P[A_1] + P[A_2] + \dots$$

Eg  $A_1 = \text{first year}$       } disjoint because  
 $A_2 = \text{second year}$       } nothing in common

$$3 = H H H$$

$$2 = H HT$$

$$2 = HT H$$

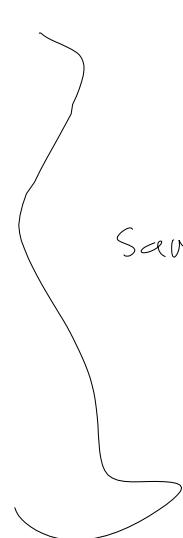
$$1 = HTT$$

$$2 = THT H$$

$$1 = TH T$$

$$1 = TTH H$$

$$0 = TTT$$



sample space

$$P[H] = \frac{1}{2}$$

$$P[T] = \frac{1}{2}$$

$$P[3H] = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$$

## Random variable

Maps sample space

$$S \rightarrow \mathbb{R} \quad , \quad X = \# \text{ of heads}$$

$X_i$	0	1	2	3
$P[X_i]$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

/ or add probabilities

↑

table called pmf or probability mass functions.

$$\text{To find } E(X) = \sum x p[X] \rightarrow \text{discrete}$$

$$\int x p[X] dx \rightarrow \text{continuous}$$

$$\text{To find } E[X^2] = \sum x^2 p[X] \rightarrow \text{discrete}$$

$$E[X^2] = \int_{-\infty}^{\infty} x^2 p[X] dx \rightarrow \text{continuous.}$$

Bernoulli:

$$p_X(x) = \begin{cases} \theta & \\ 1 - \theta & \end{cases}$$

$$\left| \begin{array}{l} E[X] / x = \# \text{ of H} \leq 0 \\ P[H] = \theta \end{array} \right.$$

$$P[T] = 1 - \varnothing$$

Poisson mean and variance is exactly the same  $\Rightarrow$  ✓

Continuous  $\quad k$

Uniform, Exp, gamma, normal, Beta.

Next class review function, mean, cdf, pdf etc...

Marginal prob example

	1	2	3	4	5	6	
H	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\vdots$
T	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$= \frac{1}{2}$

$A = H$       ;  
 $B = T$       is  $A \perp B$  ? Yes

Proof  $P[\text{Joint } (A \cap B)] = \frac{1}{12}$

$$P(A) = \frac{1}{2}$$

$$P(B) = \frac{1}{6}$$

$$P(A \cap B) = P(A) P(B) \quad \text{if } A \perp B$$

$$P[1|H] = \frac{P(1 \cap H)}{P(H)} = \frac{\frac{1}{12}}{\frac{1}{6}} = \frac{1}{2}$$

$$\begin{array}{c} x : 0 \quad 1 \quad 2 \quad 3 \\ P[x]: \frac{1}{8} \quad \frac{3}{8} \quad \frac{3}{8} \end{array}$$

## Treatment Heart Transplant

### Control group

P1 - survived 49 days after died

P28 - survived 118 days - but he is alive

P - indicator

P1 - wait time for transplant is 0 days  $\rightarrow$  survived 15 days - dead

P52 - wait time 5 days  $\rightarrow$  survived 43 days - alive

$$30 \text{ ppl} \rightarrow \text{mean}(t)$$

Treatment

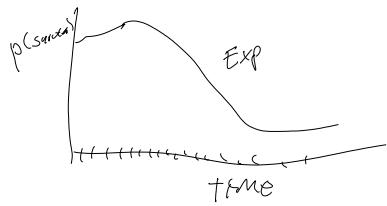
$$\text{Control}$$

P

$$52 \text{ ppl} \rightarrow \text{mean}(c)$$

compare the two values.

Assume SD is coming from a distribution  
e.g.



if mean of treatment group is bigger  $\Rightarrow$  mean of control group

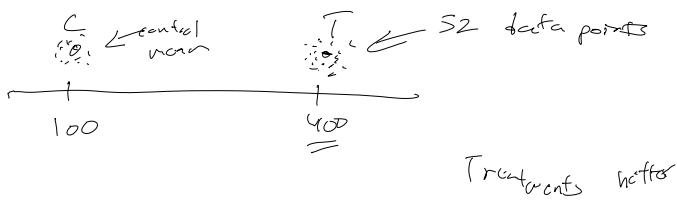
Interpretation about the population:

$\Rightarrow$  if you go through the treatment you are expected to live longer.

$\text{mean}(T) > \text{mean}(c)$  are samples, interpretation we want to make is about the entire population.

assume

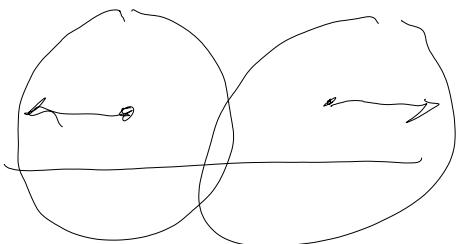
- ①  $X \sim \text{Exp}(\lambda) \rightarrow$   
 ② estimate of  $\lambda$  from these  
 fits (parameters)



In real life you never know the distribution

$$\left\{ \begin{array}{l} \text{Recall} \\ X \sim \text{Exp}(\lambda) \\ f_X(x) = \lambda e^{-\lambda x} \end{array} \right.$$

mean is the center of the distribution



radius also matters.

Data - Circle

mean - center of data

Standard deviation - radius

• look if they are overlapping  
 or far apart.

$$5.1.1 \quad \text{mean}(c) = 93.2 \\ = \frac{\text{sum (30 numbers)}}{30}$$

$$\text{mean}(t) = 356.2$$

You could just compare the mean, look at SD and then mean.

Ex 5.1.8

$$X \sim \text{Exp}(\lambda)$$

$$x_1, x_2, \dots, x_n$$

$$E(x) = \int_0^\infty x f(x) dx$$

$$= \int_0^\infty x \lambda e^{-\lambda x} dx$$

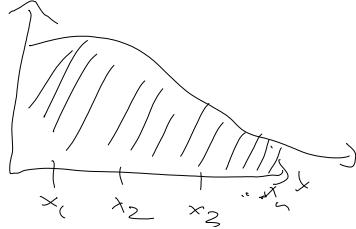
$$= \lambda \int_0^\infty x e^{-\lambda x} dx$$

difference between  $x$  and  $x_i$  is

$$f_{x|T=x} \quad \lambda = 2$$

$$f_{x|T=x} = 2e^{-2x}$$

$$f(x)$$



$$= \lambda \int_0^{\infty} x e^{-\lambda x} dx$$

$$= \frac{\lambda}{\Gamma(2)} \int_0^{\infty} x^{2-1} e^{-\lambda x} dx \cdot \frac{\lambda^2}{\Gamma(2)} \leftarrow \text{divide by } \Gamma(2)$$

gamma density = 1

$$= \frac{1}{\lambda} (1) = \frac{1}{\lambda}$$

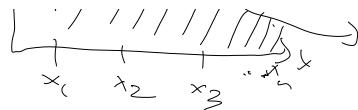
weak law of large #'s says

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \xrightarrow{\frac{1}{n}} \lambda \Rightarrow \frac{x_1 + x_2 + \dots + x_n}{n} \xrightarrow{\text{converge}} \frac{1}{\lambda} \xrightarrow{\text{approx}} \frac{1}{\bar{x}}$$

$$= \bar{x} \xrightarrow{\text{estimate}} \frac{1}{\lambda} \leftarrow \text{estimate of mean}$$

$$= \frac{1}{\bar{x}} \xrightarrow{\text{estimate}} \lambda \quad \downarrow$$

Estimate of lambda



can take any value  $X$  or  $y$ .  $x$  is the random variable that can take any value

$x_i$  is sample value  
 $x_1 = 100 \text{ days}$   
 $x_2 = 200 \text{ days}$

Recall: Gamma

$$y \sim G(\alpha, \beta)$$

$$f(y) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}, y \geq 0$$

problems

$$\bar{x} = 0.05 \quad \hat{\lambda} = 20$$

$$\bar{x} = 0.06 \quad \hat{\lambda} = 16.66$$

if  $\lambda$  is really big eg 200 to contain it  
you need a large sample size

R-vector:  $c(1, 3, 5) \Rightarrow$

1
3
5

$$c(1, 3, 5) + 2 =$$

3
5
2

variable:  $x = c(5, 7, 9, 10)$

printing:  $>x$   
 $[1] 5 7 9 10$

$$\log(x) \Rightarrow \log \begin{bmatrix} 5 \\ 7 \end{bmatrix}$$

1 2 3 ...

$$\log(x) \Rightarrow \log \begin{pmatrix} 5 \\ 7 \\ 9 \\ 10 \end{pmatrix}$$

1.60903 1.69890 ...

$\log_{10}(x)$

$\sin(x)$

$\cos(x)$

$x^2$

does it for all others

$y =$

5.1.1 /  $x \sim N(0, 1)$

$$y = x + 2x^3 - 3$$

$P(y \in (1, 2)) \leftarrow$  prob of between (1, 2)

$\chi^2 \rightarrow \text{Chi-sq}$

$\overbrace{-\infty}^{\infty}$

Ex 5.2.1

$x = \text{life length of a machine}$

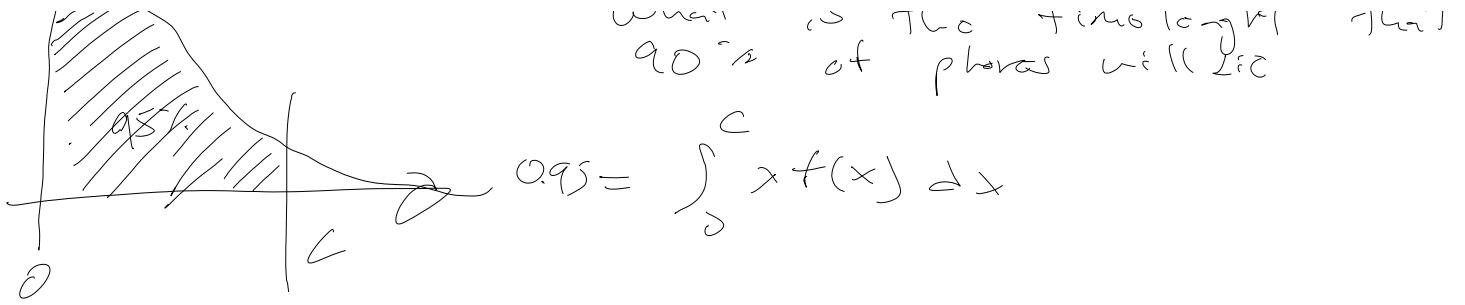
$x \sim \text{Exp}(1)$

mean life length?  $\Rightarrow$  asking for  $E(x)$

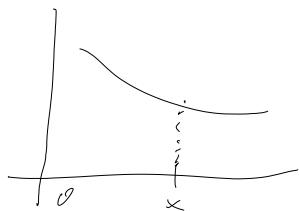
$$E \times_p(\lambda) = \frac{1}{\lambda} \Rightarrow E \times_p(1) = \frac{1}{1} = 1$$



what is the time length that 90% of phones will last



$$F(x) = \int_0^x f(x) dx \leftarrow cdf$$



$$1 - e^{-c} = 0.95$$

$$e^{-c} = \frac{0.95}{1}$$

$$\ln(e^{-c}) = \ln(0.05)$$

$$c = -\ln(0.05)$$

Conditional Distribution

X:	0	1	2	3
$P[X=x]$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

$$\text{mean} = (0)\left(\frac{1}{8}\right) + 1\left(\frac{3}{8}\right) + 2\left(\frac{3}{8}\right) + 3\left(\frac{1}{8}\right) \\ = 1.5$$

Is this valid Probability Distribution?

Yes because

- 1)  $0 \leq P(X=x) \leq 1$
- 2)  $\sum P[X=x] = 1$

Condition: cellphone survived 6 month then you won't have the same distribution

X:	0	1	2	3
$P[X=x]$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

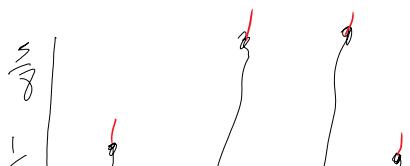
Then 0 isn't an option

This is not a valid probability distribution since it does not add up to 1.

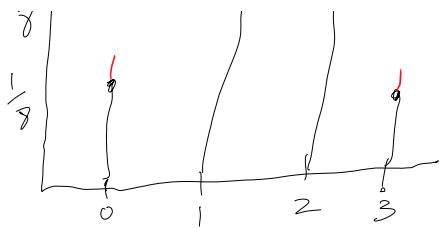
# First realize if a valid probability distribution

$$P(X=1 | X > 0) ?$$

$$\frac{P(X=1 \cap X > 0)}{P(X > 0)} = \frac{P(X=1)}{P(X > 0)} = \frac{\frac{3}{8}}{\frac{7}{8}} = \frac{3}{7}$$



# make it a bit bigger  
i.e. don't - .



if  $x \sim \text{Exp}(1)$

PDF:  $f_x(x) = e^{-x}$

$$E[x | x > 1] = \int x, \text{ conditional function}$$

$$f_{x|x>1}(x) = \frac{f_x(x)}{P[x > 1]} \leftarrow \begin{matrix} \text{condition needed} \\ \text{for } \sum \text{ to } = 1 \end{matrix}$$

$$= \frac{e^{-x}}{\int_1^\infty e^{-x} dx}$$

$$\Rightarrow = \int_1^\infty x e^{-(x-1)} dx$$

### Review S.1

Example  $\rightarrow$  Samples

$\hookrightarrow$  Distribution

$\hookrightarrow$  Inference  $\rightarrow$  commenting on the calculation

\* If  $x$  = life length

what is mean/expected life length?

$\hookrightarrow$  asking for expected value  $E(x)$

$\hookrightarrow$  by what time 95% of the machine

↳ by what time 95% of the machines will fail.

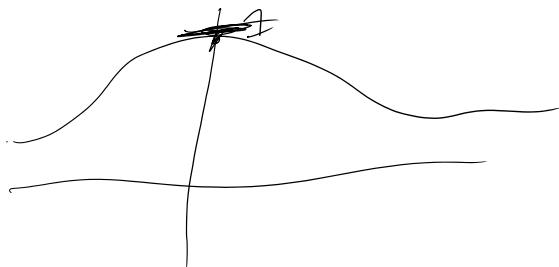
# calculate cdf, equal it to 0.95 or whatever and solve it.

2nd question: conditional, given  $x > 1$

↳ mean / expected life?

$$E(x | x > 1)$$

product life time distribution; mean or mode



mode = highest point

To find # failure distribution  
set it = 0

max min

5.2.7 example where you have to calculate the mode.

5.2.8  $x \sim p(x)$

\* Predict the future values?  $\begin{cases} \text{mean} \\ \text{mode} \end{cases}$

mode

$$P[X=x] = \frac{\lambda^x e^{-\lambda}}{x!}, \text{ if discrete no derivative}$$

$$P[X=x+1] = \frac{\lambda^{x+1} e^{-\lambda}}{(x+1)!}$$

$$\frac{P[Y=x+1]}{P[X=x]} = \frac{\cancel{\lambda}^{x+1} \cancel{e^{-\lambda}}}{(x+1)!} \div \frac{\cancel{\lambda} \cancel{e^{-\lambda}}}{x!}$$

$$= \lambda \frac{x!}{(x+1)!} = \frac{\lambda}{x+1}$$

$$\Rightarrow \frac{P[Y=x+1]}{P[X=x]} = \frac{1}{x+1}$$

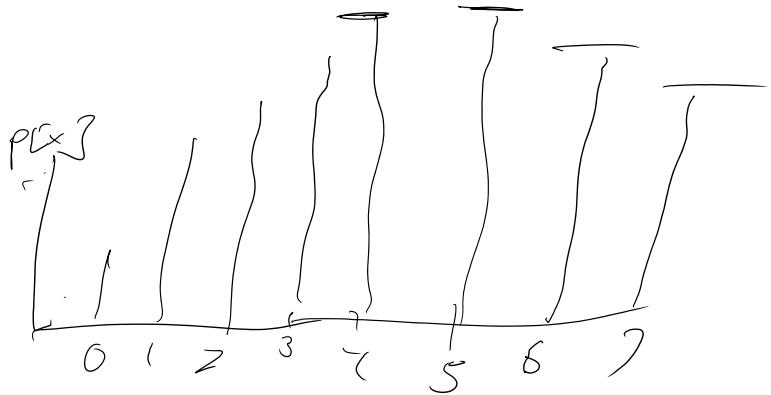
$$X=0 \Rightarrow \frac{P[X=1]}{P[X=0]} = \lambda \rightarrow \text{if } \lambda = 5 \text{ then its } S_X \text{ ratio}$$

$$X=1 \Rightarrow \frac{P[X=2]}{P[X=1]} = \frac{\lambda}{2} \rightarrow \text{if } \lambda = 5 \text{ then its } S_X \text{ ratio}$$

$$\frac{P[X=3]}{P[X=2]} = \frac{5}{3} = 1.666$$

$$\frac{P[X=4]}{P[X=3]} = \frac{5}{4} = 1.25$$

$$\frac{P[x=5]}{P[x=4]} = \frac{5}{5}$$



### 5.3 Statistical Models

① Population: a combination of all the subjects in your subspace. All the outcome.

Sample: A small subset of the population  
 $\text{Sample} \subseteq \text{Population}$

Parameter: Any characteristic of a population is a parameter

populations are too big samples are easier to work with.

# Study sample  $\Rightarrow$  make inference about population. That is statistics

$$\text{Eg } f_{\lambda}(x) \quad , \quad p_{\theta}(x)$$

$$f(x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad x \text{ is r.v., } \lambda \text{ is the parameter.}$$

Statistical model

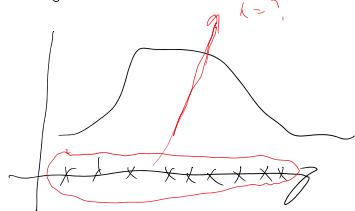
$\left\{ \begin{array}{l} p_{\theta}: \Omega \rightarrow \mathbb{R} \\ \text{function} \\ \text{e.g. poisson} \\ \text{normal} \end{array} \right.$	$\theta$ can take any value in the sample space
---	--

parameter space

$$\frac{e^{-5} 5^x}{x!} \sim p_5(x)$$

$$\frac{e^{-100} 100^x}{x!} \sim p_{100}(x)$$

$$\text{Eg } x \sim p_{\lambda}(\lambda = ?)$$



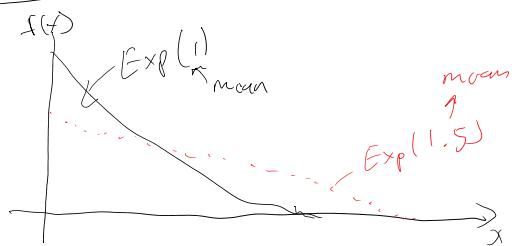
$$\bar{x} \rightarrow \frac{1}{\lambda}$$

we use  $\bar{x}$  to estimate  $\frac{1}{\lambda}$ . This is an example of point estimation.

## Interval estimation - Review

Example: Pg 264

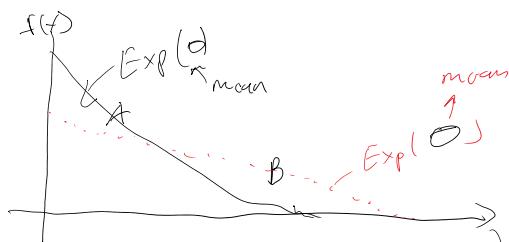
Ex 5.3.2



$$P_\theta(x) = \underbrace{f_\theta(x)}_{\text{discrete}} \quad \begin{cases} \text{continuous} \\ \text{continuous} \end{cases}$$

Sample  $(x_1, x_2, \dots, x_5)$

In this case, we know the parameters as the mean is already there.



← we should be able to distinguish between if a function by a mean of 1 or 1.5.

$$\textcircled{1} (x_1, \dots, x_5) = (5.0, 3.5, 3.3, 4.1, 2.8)$$

$$\textcircled{2} (x_1, \dots, x_5) = (2.0, 2.5, 3.0, 3.1, 3.0)$$

which one belongs with which sample?

$$\textcircled{1} \Rightarrow \text{Exp}(1.5)$$

$$\textcircled{2} \Rightarrow \text{Exp}(1)$$

Parameter space in this case is

$$\Pi = \{1, 1.5\}$$

$$\Omega = \{A, B\}$$

Ex 5.3.1

$$\textcircled{1} \text{ Exp} \quad \leftarrow \text{know function}$$

$$\textcircled{2} \Omega = \{1, 1.5\} \quad \leftarrow \text{know parameter space}$$

$$\textcircled{3} \text{ Sample} \leftarrow \text{select sample and pick}$$

Real life

- ① observe sample
- ② find the distribution
  - known or assumed

② Sample  $\leftarrow$  "small sample" and pick  $\theta$  - known or assumed

③ Parameter - Use sample to make inference about parameter

5.3.2 parameter

$\theta$	$f_\theta(x=1)$	$f_\theta(x=2)$	$f_\theta(x=3)$
A	0.5	0.5	0
B	0	0.5	0.5

$\rightarrow$  space  $\{A, B\}$

b) if we observe  $x=1$  then  $x$  is coming from A

if we observe  $x=3$  then  $x$  is coming from B

if we observe  $x=2$  then either A or B.

parameter

$\theta$	$f_\theta(x=1)$	$f_\theta(x=2)$	$f_\theta(x=3)$
A	$\frac{1}{2}$	$\frac{1}{2}$	0
B	0	0.5	0.5

what if  $x=2$ ?

probably A

### Notation

$$\{f_\theta : \theta \in \Omega\}$$

for one <sup>sample</sup> value what is the statistic model

$$\text{Ex } \hat{P}(\lambda)(x)$$

$$\text{Ex } x \sim \text{Exp}(\theta)$$

$$x = L$$

$$\sum \frac{\hat{P}(\lambda)^x \cdot \lambda^x}{x!}$$

$$f_\theta(x) = \theta e^{-\theta x}$$

$$x=5 = \theta e^{-\theta(5)} \leftarrow \text{statistical model of one sample.}$$

Joint density of sample

$$\text{Ex } P[x_1, x_2, y_3, \dots, x_n] \text{ if I break it down}$$

$$\Rightarrow f_\theta(x_1) f_\theta(x_2) \dots f_\theta(x_n)$$

I is an assumption

$$x = \{2, 4, 9\}$$

$$X = \{2, 4, 9\}$$

$$\theta e^{2\theta} \times \theta e^{-k\theta} \times \theta e^{-\theta}$$

Ex 5.3.3

$$x \sim \text{Bin}(\theta)$$

$$p_X(x) = \begin{cases} 1, & x = \theta \\ 0, & x = 1 - \theta \end{cases}$$

$$p_X(x) = \theta^x (1-\theta)^{1-x}$$

$$x_1, x_2, \dots, x_n$$

$$\theta^x (1-\theta)^{1-x_1} \times \theta^{x_2} (1-\theta)^{1-x_2} \times \dots \times \theta^{x_n} (1-\theta)^{1-x_n}$$

$$= \theta^{x_1 + x_2 + \dots + x_n} (1-\theta)^{n - (x_1 + x_2 + \dots + x_n)}$$

$$= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$$

Ex 5.3.4

$$(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$$

$$\theta = (\mu, \sigma^2) \in \mathbb{R}^1 \times \mathbb{R}^+, \quad \mathbb{R}^+ = (0, \infty)$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_1-\mu)^2} \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_2-\mu)^2} \times \dots \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_n-\mu)^2}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \left[ -\frac{1}{2\sigma^2} \sum (x_i - \mu)^2 \right]$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\begin{aligned}
 &= \sum_{i=1}^n (x_i - \mu)^2 \quad \begin{matrix} a-b \\ a-c < b \\ a-c < -b \end{matrix} \\
 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2 \\
 &= \sum_{i=1}^n (x_i - \bar{x})^2 + 2(\bar{x} - \mu) \sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n (\bar{x} - \mu)^2 \\
 &\quad \begin{matrix} \cancel{\sum_{i=1}^n (x_i - \bar{x})} \\ = \sum x_i - n\bar{x} \end{matrix} \quad \text{sum of deviation from mean} = 0 \\
 &= n\bar{x} - n\bar{x} = 0
 \end{aligned}$$

$$\sum_{i=1}^n (\bar{x} - \mu)^2 = n(\bar{x} - \mu)^2$$

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \sim \text{sample variance}$$

$x_i$	$\bar{x}$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
2	-2	4	16
4	0	4	16
6	2	4	16
	0	0	0
	8	0	0

$$\begin{aligned}
 s^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\
 &= \frac{8}{2} = 4
 \end{aligned}$$

$\frac{1}{Var} = \text{precision}$

$$\lambda = \frac{1}{6^2}$$

precision

$$\begin{aligned}
 &(2\pi\sigma^2)^{-\frac{1}{2}} \exp\left[-\frac{1}{2\sigma^2} (x - \mu)^2\right] \\
 &= (2\pi)^{-\frac{1}{2}} (\lambda)^{\frac{1}{2}} \exp\left[-\frac{1}{2\lambda} \lambda (x - \mu)^2\right] \quad \leftarrow \text{still normal but in precision instead of variance}
 \end{aligned}$$

Reparameterization - only do it if it's a one-to-one function. Precision is  $\frac{1}{\text{variance}}$ ) the change old parameter to new parameter goes to no one to one

old param  $\xrightarrow{\text{func}}$  new param.

Tues Oct 3, 3.5.

## Population

Population CDF

$$F_X(x) = \frac{\text{Count } x \leq X}{N}$$

Ex 5.1

$$N = 20$$

$$\underline{\text{min}} = 3$$

$$P(X \leq 3) = 0$$

$$P(X \leq 4) = \frac{3}{20}$$

$$P(X \leq 4) - P(X \leq 3) = P(X=3)$$

5.4 = Population = finite  
= Sample — distribution known  
→ Unknown distribution  
→ Parameter

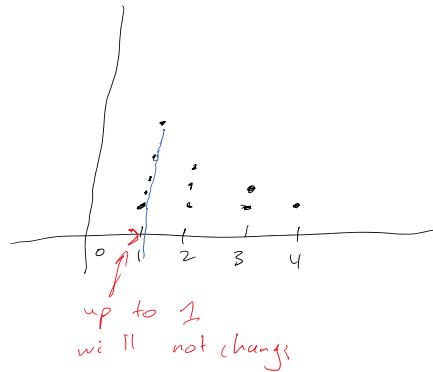
If a population is finite, do you need money/time to study them?

Ex

i	1	2	3	4	5	6	7	8	9	10
$X(\pi_i)$	1	1	2	1	2	3	3	1	2	4

Sort

i	1	1	2	1	2	3	3	1	2	4
$X(\pi_i)$	1	1	1	2	2	3	3	4		



$$P[X \leq 0] = 0$$

$$P[X \leq 0.999] = 0$$

$$P[X \leq 1] = \frac{4}{10} = 0.4$$

$$P[X \leq 1.999] = \frac{4}{10} = 0.4$$

$$P[X \leq 2] = \frac{7}{10} = 0.7$$

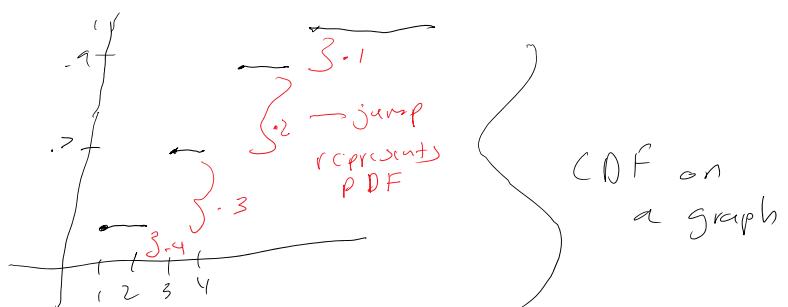
$$P[X \leq 2.999] = \frac{7}{10} = 0.7$$

$$P[X \leq 3] = \frac{9}{10} = 0.9$$

$$P[X \leq 3.999] = 0.9$$

$$P[X \leq 4] = 1$$

$$CDF = F_X(x) = \begin{cases} 0 & x < 1 \\ .4 & 1 \leq x \\ .7 & 2 \leq x < 3 \\ .9 & 3 \leq x \leq 4 \\ 1 & x \leq 4 \end{cases}$$



PMF:

$$f_X(x) = \begin{cases} .4, & x=1 \\ .3, & x=2 \\ .2, & x=3 \\ .1, & x=4 \\ 0, & \text{otherwise} \end{cases} \Rightarrow$$

X	1	2	3	4
$P[X]$	.4	.3	.2	.1

To calculate small f, pdf just calculate the proportion.

$$\text{CDF} = F_x(x) = \begin{cases} 0 & x < 1 \\ .1 & 1 \leq x \\ .7 & 2 \leq x < 3 \\ .9 & 3 \leq x \leq 4 \\ 1 & x \geq 4 \end{cases}$$

Population

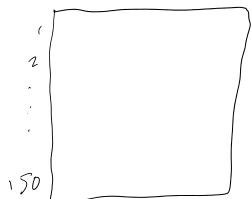
Same calculation but if it's sampled it's called Empirical distribution function

Empirical distribution: Same calculation but put  $\hat{F}_x(x)$

Ex/ do it on your own drug.

### Simple random Sampling:

Ex, in a class, blindly picking #'s



To do it in R:  
`sample(1:150, size=1)`  
 for a random sample between 1-150

SA, B, C, D, E, Z

$N=5$  - draw 2 samples  $n=2$

$$P(A) = \frac{1}{5}$$

$$P(A \text{ being selected}) = \frac{1}{N} = \frac{1}{5}$$

$$P(A \text{ being selected} | b \text{ is selected}) = \frac{1}{N-1} = \frac{1}{4} = 0.25$$

Samples are not  $\perp$  because one being in the sample changes prob of the other one.

After picking with replacement — pick and put it back  $\rightarrow$  samples are  $\perp$

$$\frac{1}{N} \quad \frac{1}{N-1}$$

$$N = 100,000$$

0.2

0.25

0.00060 |  $\approx$  next numbers  
some things  
are closer.  
like that

large sample even though they are dependent because N is large change is insignificant, so we treat as 1

Conditions

$N \rightarrow$  large

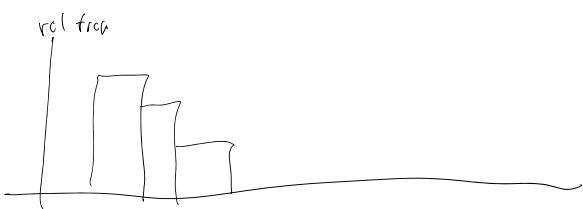
$n \rightarrow$  small relative to  $N$

$\therefore \hat{F}_x(x) \rightarrow F_x(x)$   
 $\underbrace{\text{cdf}}$   
calculated  
based on  
sample

Histogram

Ex height

$(h_1, h_2] (h_2, h_3] (h_3, h_4] \dots$   
5 5:6 5:6 6 6 6:5



relative frequency = proportion

Density histogram



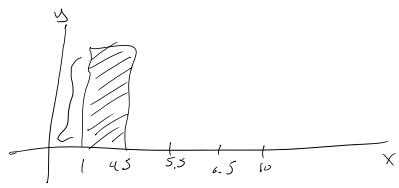
$$h_x(k) = \frac{\text{proportion}}{\text{length}}$$

5.4.5

$x \in [1.2 \ 1.8 \ 2.3 \ 2.5 \ 3.1 \ 3.4 \ 3.7 \ 3.2 \ 3.9 \ 4.3 \ 4.4 \ 4.5 \ 4.5]$

$(4.8 \ 4.8] (5.6 \ 5.8] (6.9 \ 7.2 \ 8.5]$

$$h_x(x) = \frac{\frac{13}{20}}{(4.5 - 1)} \quad (1, 4.5]$$



$$h_x(x) = \frac{\frac{13}{20}}{(4.5 - 1)} \quad , (1, 4.5)$$

$$= \frac{13}{20}$$

Question on  
Midterm on this.

$$(4.5, 5.5] \quad , (4.5, 5.5]$$

$$h_x(x) = \frac{\frac{2}{20}}{(5.5 - 4.5)}$$

Ex

$$f_x(x) = \begin{cases} .4 & , x = 1 \\ .3 & , x = 2 \\ .2 & , x = 3 \\ .1 & , x = 4 \\ 0 & , \text{o/w} \end{cases}$$

from this list make the table

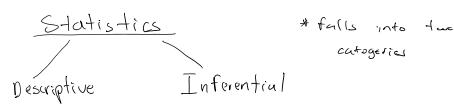
# look for unique numbers

① 2, 3, 4

# count pop by total + get proportion

# loop

Prof uploaded yr code. look at the table.



Descriptive: describes any summary  
eg mean of sample  
standard deviation  
median

Inferential: moment you use those #'s to make a inference about the data.

Recall:  $f_x(x)$  is the proportion of the population members whose  $X$  measurements equal  $x$ .

$$\Rightarrow f_x(x) = P[X=x]$$

$F_x(x)$  is the proportion of population members whose  $X$  measurements is less than or equal to  $x$

Ex  $\{1.2, -2.1, 0.4, 3.3, -2.1, 4.0, -0.3, 2.2, 1.5, 5.0\}$

# put in ascending order

$$\{-2.1, -0.3, 0.4, 1.2, 1.5, 2.1, 2.2, 3.3, 4.0, 5.0\}$$

# flag on each data point to get PDF

$$f_x(x) =$$

$$f_x(-3) = P[X=-3] = 0$$

$$f_x(-2.1) = P[X=-2.1] = \frac{1}{10}$$

$$f_x(-0.3) = P[X=-0.3] = \frac{1}{10}$$

$$f_x(5) = P[X=5] = \frac{1}{10}$$

CDF same thing but everything up to the point

$$F_x(-2.1) = P[X \leq -2.1] = \frac{1}{10}$$

$$F_x(-0.3) = P[X \leq 0.3] = \frac{2}{10} = \frac{1}{5}$$

A natural estimate of  $F_x(x)$  is given by  $\hat{F}_x(x)$

$$\hat{F}_x(x) = \sum_{i=1}^n I_{(-\infty, x]}(x_i)$$

that's the formal definition for what we did above

This is also called empirical distribution function of  $x$

$\hookrightarrow$  means sample.

### Calculating Population Quantiles

Given value calculate percentile.

$$x(60)$$

p-quantile

Note .75 quantile = 75th percentile



↳  $F(x) = \Pr[X \leq x]$

p-quantile

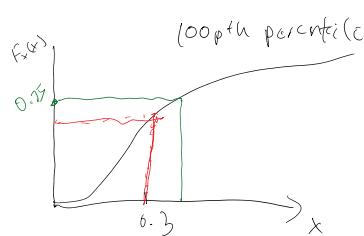
$N_{100} \rightarrow 75$  quantile  $\hat{x}$   $\approx 75^{\text{th}}$  percentile

100<sup>th</sup> percentile

$\xrightarrow{\text{100}}$

$$\frac{1}{100}$$

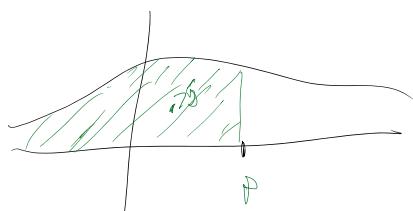
80<sup>th</sup> percentile  $\hat{x} = 0.8$  quantile



$$F_x(0.3) = \Pr[X \leq 0.3]$$

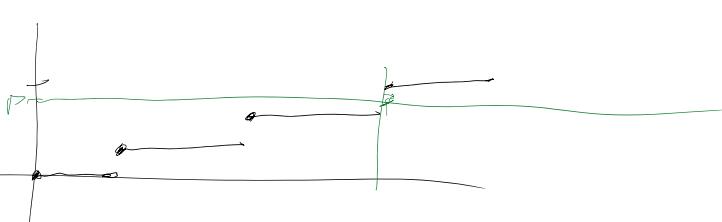
Percentile is opposite

$$F_x(x) = \Pr[X \leq x] = 0.75$$



$$F(y_p) = p$$

$y_p = ?$  if continuous  
cdf



$\{1.2, 2.1, 0.4, 3.3, -2.1, 4.0, -0.3, 2.2, 1.5, 5.0\}$

# put in ascending order

$\begin{array}{ccccccccc} x(1) & x(2) & & & & & x(n) \\ \frac{-2.1}{10} & \frac{-0.3}{10} & \frac{0.4}{10} & \frac{1.2}{10} & \frac{1.5}{10} & \frac{2.1}{10} & \frac{2.2}{10} & \frac{3.3}{10} & \frac{4.0}{10} & \frac{5.0}{10} \end{array}$

70<sup>th</sup> percentile?  $2.2$

$$\Pr[X \leq 2.2] = 0.7$$

\* if you don't have equals then go to the next one

$P = 0.75 ?$  its a # between  $2.2 - 3.3$

$2.2 (\frac{70}{100})$   $1.5 (\frac{75}{100})$   $3.3$

take  $\frac{7}{10}$

$$2.2 (\frac{5}{10} \text{ of data}) = 2.2 + \frac{(3.3 - 2.2)}{2}$$

$$\frac{i-1}{n} \leq P \leq \frac{i}{n}$$
$$\frac{7}{10} \leq P \leq \frac{8}{10}$$



Continuous one always a solution, if discrete no solution sometimes!

$$x = x_{(i-1)} + (x_i - x_{i-1}) n \left( p - \frac{i-1}{n} \right) \stackrel{\text{in this case}}{\Rightarrow} 2.2 + (3.3 - 2.2) 10 (0.75 - 0.7) \\ \Rightarrow 2.2 + (1.1)(0.5) \\ = 2.75$$

25<sup>th</sup> percentile ← midterm / final question

$$\frac{i-1}{n} < p \leq \frac{i}{n} \\ 0.2 < p \leq 0.3, i=3$$

$$x = x_{(i-1)} + (x_i - x_{i-1}) n \left( p - \frac{i-1}{n} \right) \quad \cancel{\text{unconfirmed}} \\ \approx -0.3 + (0.4 + 0.3) 10 (0.25 - 0.2) \\ \approx 0.05$$

25<sup>th</sup> percentile =  $Q_1$  = 1<sup>st</sup> quartile

75<sup>th</sup> percentile =  $Q_3$  = 3rd quartile

$$P[x \leq 1.5] = 0.5$$

$$P[x \geq 1.5] = 0.6$$

$$\{1.2, 2.1, 0.4, 3.3, -2.1, 4.0, -0.3, 2.2, 1.5, 5.0\}$$

# put in ascending order

$$\{-2.1, -0.3, 0.4, 1.2, \textcircled{1.5}, 2.1, 2.2, 3.3, 4.0, 5.0\}$$

$\uparrow$   
median

If calculate percentile use this formula.

on midterm and final

In one definition

$P[x \leq 1.5] = 0.5$  is enough to call it a median but in another definition

$P[x \geq 1.5] = 0.6$  will not since it is  $!= 0.5$ .

Result if  $n$  is odd  $\rightarrow \frac{n+1}{2}$  th term

$$\text{even} \rightarrow \frac{\frac{n}{2} \text{th} + \frac{n+1}{2} \text{th term}}{2}$$

$$P(x \geq 1.8) = 0.5 \quad \checkmark$$

$$\begin{array}{c|cc} 1.5 & & 2.1 \\ \hline 1.8 & & \end{array}$$

Two definition of median, also easiest one.

$$1, 3, 5 \text{ median is } \frac{n+1}{2} = \frac{3+1}{2} = 2 \text{nd term}$$

7

0 1

Interquartile range: (width of data)

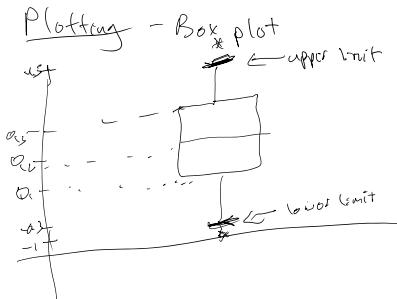
$$IQR = Q_3 - Q_1$$

replacement of SD.



extreme outliers median doesn't suffer

Skew  $\rightarrow$  median  
symmetric  $\rightarrow$  mean.



$$\text{lower limit} = Q_1 - 1.5 \times IQR$$

$$\text{upper limit} = Q_3 + 1.5 \times IQR$$

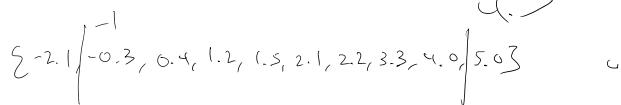
\* Box plots are for continuous cases.

$$\text{lower limit} = Q_1 - 1.5 \times (Q_3 - Q_1) = 1 \leftarrow \text{means line goes down to } -1$$

stop at -0.3      \* = outliers

$$\begin{aligned} \text{upper limit} &= Q_3 + 1.5 \times (Q_3 - Q_1) \\ &= 4.5 \end{aligned}$$

stops at 4.5



4.5 is also an outlier.

whiskers = \*

adjacent values: 4.5, -1

Ex/

- (1) CMC 0.42
- (2) Van 0.28
- (3) BS 0.22
- (4) St 0.08

} categorical variable has no order.

if categorical stuff above does not apply

Do this instead

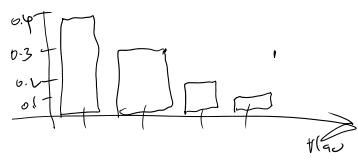


freq

# put shift in order

# determine regim





cover

# determine region

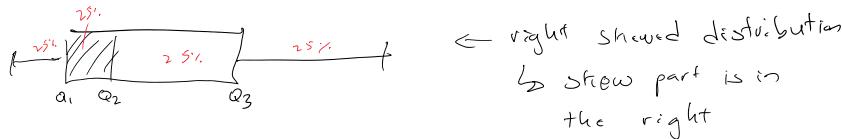
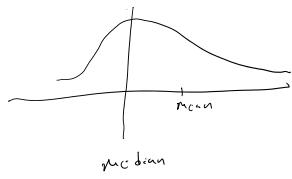
# calculate percentile



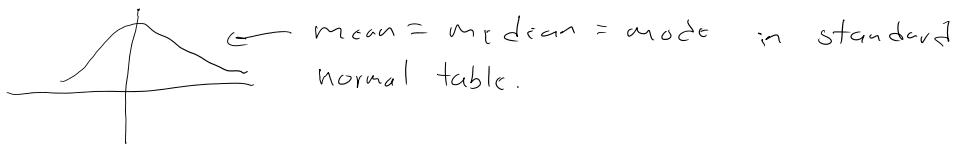
Recap

$$\tilde{x}_p = x_{(i-1)} + (x_{(i)} - x_{(i-1)}) n \left( p - \frac{i-1}{n} \right)$$

odd	$\left(\frac{n+1}{2}\right)$ th term
even	$\frac{n}{2}$ th + $\left(\frac{n}{2}+1\right)$ th

skewness:

if skewed use median and not the mean



## 5.2 $X \sim \text{Exp}(1)$

- ① Predict future value
  - ②  $P[X > 5]$
  - ③  $P[X \leq a] = 0.95$
- These are inferences about the future  $X$ . (observation)

5.3 inference is about distribution

In this chapter  $\text{Exp}(\theta) \rightarrow$  estimate  $\theta$

chi  
 $\psi(\theta) = \frac{1}{\theta}$       mean =  $\frac{1}{\theta}$       var( $\theta$ ) =  $\frac{1}{\theta^2}$

$$\psi(\theta) = \frac{1}{\theta^2} \quad \text{median} = \int_0^M f_\theta(x) dx = 0.5$$

$$\Rightarrow F(m) = 0.5 \\ \Rightarrow m = F^{-1}(0.5)$$

$\boxed{\psi(\theta) = F_\theta^{-1}(0.5)}$

median

25<sup>th</sup> percentile:  $\psi(\theta) = F_\theta^{-1}(0.25)$

Inverse ex/ Exp(1)

$$f_\theta(x) = e^{-x}$$

$$F_\theta(x) = \int_0^x e^{-r} dr = 1 - e^{-x}$$

$$\begin{aligned} F_\theta(m) &= 1 - e^{-m} = 0.5 \\ e^{-m} &= 0.5 \\ -m &= \ln(0.5) \\ \Rightarrow m &= -\ln(0.5) \\ \Rightarrow F_\theta^{-1}(0.5) & \end{aligned}$$

$$f_\theta(x) \rightarrow \text{Normal}(\mu, \sigma^2)$$

$\psi(\theta) \leftarrow$  Parameter of interest

$$\psi(\theta) = \mu$$

$x_1, x_2, \dots, x_n \leftarrow$  can you use this to get mean?

yes  $\Rightarrow \bar{x}$  use sample mean

$$T(s) = \frac{1}{n} \sum x_i$$

$T(s)$  is an estimate of  $\psi(\theta)$

Variance:  $\psi(\theta) = \sigma^2$

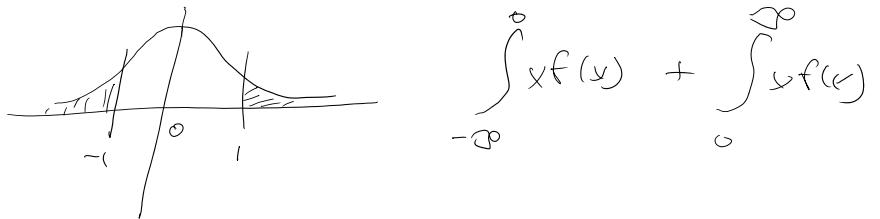
$$T(s) = \frac{\sum (x_i - \bar{x})^2}{n-1} \rightarrow \text{an estimate of } \frac{\sigma^2}{\psi(\theta)}$$

$$\begin{matrix} \bar{x} & \longrightarrow & E(x) \\ n & \longrightarrow & \end{matrix}$$

### 5.5.8

3rd moment  $N(\mu_0, \sigma_0^2)$   
mean not sigma not squared

$$\begin{aligned}
 \psi(\theta) &= E(x) = E[x - \mu + \mu]^3 \\
 &= E[(x-\mu)^3 + 3(x-\mu)^2\mu + 3(x-\mu)\mu^2 + \mu^3] \\
 &= E[(x-\mu)^3] + E[3(x-\mu)^2\mu] + E[3(x-\mu)\mu^2] + E[\mu^3] \\
 &= \underbrace{E[(x-\mu)^3]}_{\substack{\text{3rd central} \\ \text{moment}}} + \underbrace{3\mu E[(x-\mu)^2]\mu}_{\substack{\text{var}(x) \\ \sigma^2}} + 3\mu^2 E[(x-\mu)] + \mu^3
 \end{aligned}$$



$$x \leq 1 = x \geq 1$$

if the distribution is symmetric around ( $\mu$ )

$$E[x-\mu] = 0$$

$$E[(x-\mu)^3] = 0$$

$$E[(x-\mu)^5] = 0$$

$$\begin{aligned}
 &E[(x-\mu)^3] + 3\mu E[(x-\mu)^2]\mu + 3\mu^2 E[(x-\mu)] + \mu^3 \\
 &\quad \substack{\text{3rd central} \\ \text{moment}} \quad \substack{\text{var}(x) \\ \sigma^2} \quad \substack{\text{||} \\ 0} \\
 &= 3\mu\sigma^2 + \mu^3
 \end{aligned}$$

if  $N(\mu, \sigma^2)$ , what's  $\Theta$ ?  $\Theta = \{\mu, \sigma^2\}$   
 means all first of parameter

$$G(\alpha, \beta) \Rightarrow \Theta = \{\alpha, \beta\}$$

Section: 5.2

- Predicting future value

- $X \sim (\theta)$ 
  - mean
  - mode: predicting  $\hat{x}$
  - $P[X > s], P[X \leq s]$
  - $E[X | Y = y]$

Section 5.3

$S = x_1, x_2, x_3, \dots, x_n$  (Let's say you have a sample,  
 $f_\theta(s) = f_\theta(x_1, x_2, \dots, x_n)$  what is  $f_\theta(s)$  ← probability function  
 of the sample)

- reparameterization:  $\theta \rightarrow \theta$  one to one then we can reparameterizeSection: 5.4

$$F_p(x) \rightarrow f_x(x)$$

population CDF

$$\begin{array}{c} -3 \\ -2 \\ -1 \\ 0 \\ 1 \\ 2 \\ 3 \end{array} \xrightarrow{\text{CDF}} \begin{array}{c} 0 \\ 0.2 \\ 0.5 \\ 0.8 \\ 0.9 \\ 0.98 \\ 1 \end{array}$$

- Once you know the sample  $F_p$ , can calculate everything.-  $\hat{F}_n(x) \leftarrow$  empirical distribution function

- histogram

-  $\hat{f}_n(x)$ 

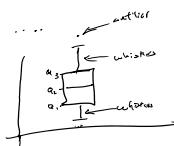
- median

- quartiles

- IQR

- Q-LS(IQR)

- box plot

-  $\hat{\theta} = \bar{x}, \hat{\sigma}^2$  $E(\hat{\theta}) = \dots$  (refer to last class notes)6.1 Likelihood Function

$X \sim \text{Bern}(\theta) \perp \theta \in [0, 1]$

$S = \{x_1, x_2, \dots, x_n\}$

$[f_\theta(S) = f_\theta(x_1) \times f_\theta(x_2) \times \dots \times f_\theta(x_n)]$

Toss coins  $S$  times:

$[1, 0, 1, 1, 0] = S$

$f_\theta(S) = \theta(1-\theta)\theta \theta(1-\theta) \dots$

$L(\theta|S) = \theta^3(1-\theta)^2$

likelihood function  
 of theta given  
 the sample.

 $* \text{Not the probability of observing } \theta$  $\text{It is the probability of observing}$  $\text{the sample for a given true } \theta.$  $\theta \rightarrow$  fixed true value,  $\theta \in \Omega$  $\theta$  any other member from  $\Omega$ below  $f_{\theta_1}(S) > f_{\theta_2}(S)$ 

The more data the smaller the likelihood.

## Ex/ 6.1.1

Suppose  $S = \{1, 2, \dots, 3\}$  and that the statistical model is  $[\theta : \theta \in \{1, 2\}]$ ,where  $P_1$  is the uniform distribution on the integers  $\{1, \dots, 10^3\}$  and $P_2$  is the uniform distribution on  $\{1, \dots, 10^6\}$ 

$$\text{Dist 1} \rightarrow f_{\theta_1} \rightarrow \text{Unif } \{1, 2, 3, \dots, 1000\}$$

$\sim \text{Unif } \{1, 2, 3, \dots, 1000000\}$

Dist 1  $\rightarrow$  100

Dist 2  $\rightarrow$  for Unit  $\{1, 2, 3, \dots, 1000000\}$

$S = \{100\}$   $\rightarrow$  comes from distribution 2. Since distribution 2 stops at 1000.

$S = \{10\}$   $\rightarrow$  # calculate likelihood for both and compare

$$L(\theta_1 | S) = \frac{1}{1000} \Rightarrow \text{sample 10 is a } 1000 \times$$

$$L(\theta_2 | S) = \frac{1}{1000000} \Rightarrow \text{more likely to come from } \theta_1$$

Example  $x \sim N(\theta, \sigma^2)$ ,  $\theta$  known?

1 unknown parameter  $\theta$ , mean

$$S = x_1, x_2, \dots, x_n$$

$$f_{\theta}(x) = \prod_{i=1}^n \underbrace{(2\pi\sigma^2)^{-\frac{1}{2}}}_{\text{constant}} \exp\left[-\frac{1}{2\sigma^2}(x_i - \theta)^2\right]$$

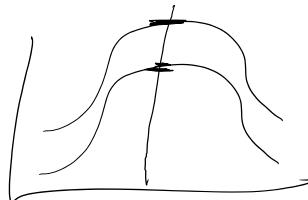
$$= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2\right] \quad \begin{array}{l} \text{Recognize this} \\ \text{expand it, Do it} \\ \text{yourself.} \end{array}$$

$$= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left[-\frac{n}{2\sigma^2} (\bar{x} - \theta)^2\right] \exp\left[-\frac{1}{2\sigma^2} (n-1)s^2\right]$$

constant sample variance

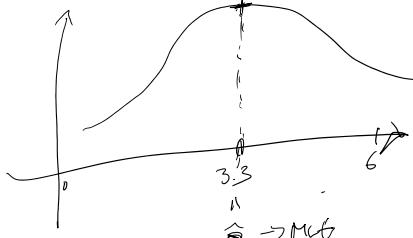
Since constant we write

$$L(\theta | S) \propto \exp\left[-\frac{n}{2\sigma^2} (\bar{x} - \theta)^2\right] \quad \Rightarrow$$



$$n = 25, \sigma^2 = 1, \bar{x} = 3.3$$

$$L(\theta | S) \propto \exp\left[-\frac{25}{2}(3.3 - \theta)^2\right]$$



$$\begin{cases} \text{vertical} \\ y = \bar{x} \\ y \propto x \end{cases} \quad \begin{array}{l} \text{proportional} \\ \text{to } x \end{array}$$

Same graph but just stretched as a result you can make inference about param

① Point of estimate of  $\theta$

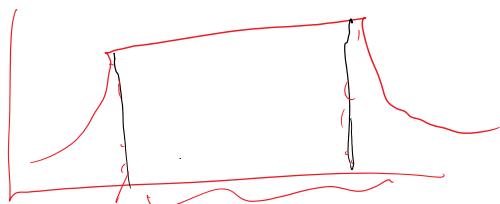
$$\hat{\theta} \rightarrow \theta$$

estimate

MLE  $\Rightarrow$  maximum likelihood estimator of  $\theta$

$$L(\hat{\theta} | S) \geq L(\theta | S)$$

Note you can have multiple MLE for  $\sigma^2$



this whole  
thing is MLE

Ex 6.2.1

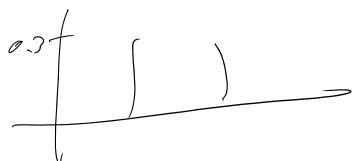
	$s=1$	$s=2$	$s=3$
$f_1$	0.3	0.4	0.3
$f_2$	0.1	0.7	0.2

$\{s=1\} \Rightarrow$  probably came from 1<sup>st</sup> one

$\{s=2\} \Rightarrow$  probably came from 2<sup>nd</sup> one

$\{s=3\} \Rightarrow$  probably came from 3<sup>rd</sup> one

	$s=1$	$s=2$	$s=3$
$f_1$	0.3	0.4	0.3
$f_2$	0	0.7	0.3

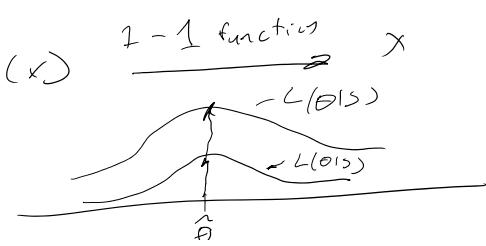


$\{s=3\} \Rightarrow$  MLE is probably  
 $\theta_1$  or  $\theta_2$

$$L(\theta(s)) = \prod_{i=1}^n f_{\theta}(x_i)$$

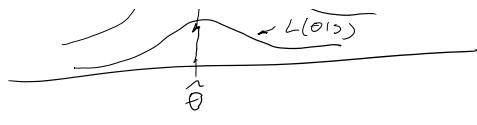
# taking the ln  
allows us to  
change to  
summation

$$\ln(L(\theta(s)))$$



$$\ln(L(\theta|s)) \quad \text{sum...}$$

$$= \sum_{i=1}^n \ln f_\theta(x_i)$$



# Now maximize

$$L(\theta|s) \propto \exp\left[\frac{-n}{2\sigma^2} (\bar{x} - \theta)^2\right]$$

proportional

$$(1) \quad l(\theta|s) = \frac{-n}{2\sigma^2} (\bar{x} - \theta)^2 \leftarrow \text{ln likelihood of } \theta$$

# Differentiate

$$(2) \quad \text{Score function} = \frac{\partial l(\theta|s)}{\partial \theta}$$

$$(3) \quad \text{Score equation} = \frac{\partial l(\theta|s)}{\partial \theta} = 0$$

$\Rightarrow \hat{\theta} = \boxed{\text{something}}$

# To check if max or min, second derivative.

$$\frac{\partial^2 l(\theta|s)}{\partial \theta^2} \Bigg|_{\theta=\hat{\theta}} < 0$$

$$l(\theta|s) = \frac{-n}{2\sigma^2} (\bar{x} - \theta)^2$$

$$\begin{aligned} \frac{\partial l(\theta|s)}{\partial \theta} &= -\frac{2n}{2\sigma^2} (\bar{x} - \theta)(-1) \\ &= \frac{n}{\sigma^2} (\bar{x} - \theta) \quad \leftarrow \text{score} \end{aligned}$$

$$\Rightarrow \frac{n}{\sigma^2} (\bar{x} - \theta) = 0$$

$$\theta = \bar{x}$$

# Check MLE

$$\text{Score} = \frac{n}{\sigma^2} (\bar{x} - \theta) \quad \sim \text{sample size}$$

$$\frac{\partial^2 l(\theta|s)}{\partial \theta^2} = \frac{n}{\sigma^2} \quad \leftarrow \text{62 positive since } n > 0$$

$$\frac{\partial^2 l(\theta|s)}{\partial \theta^2} < 0$$

or (a)

$$\therefore \frac{-n}{\theta^2} \leq 0$$

$\therefore \bar{x}$  is the MLE of  $\theta$

Do it for exp distribution

Question midterm or final or both, learn to calculate MLE of everything distribution we know.

### Invariance

$f_\theta(x) \xrightarrow{\theta} \hat{\theta}$  is MLE

$$\psi(\theta)$$

function  
of  $\theta$

how do you find MLE?

$$\psi(\theta) \xrightarrow{\theta} \psi(\hat{\theta})$$

b. 2.2 If  $(x_1, \dots, x_n)$  is a sample from  $Bir(\theta)$  distribution where  $\theta \in [0, 1]$  is unknown, then determine the MLE of  $\theta$

$x \sim Bir(\theta)$   $\leftarrow$  mle?

$$S = x_1, x_2, \dots, x_n$$

# write the distribution plug  $x_1, x_2, \dots, x_n$

$$\text{Likelihood } L(\theta | S) = \theta^{x_1} (1-\theta)^{1-x_1} \cdot \theta^{x_2} (1-\theta)^{1-x_2} \cdots \cdots \cdot \theta^{x_n} (1-\theta)^{1-x_n}$$

$$= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$$

# get mle?

In likelihood

$$\begin{aligned} \ell(\theta | s) &= \ln(\theta^{\sum_i x_i}) + \ln(1-\theta)^{n-x_i} \\ &= \sum_{i=1}^n x_i \ln \theta + (n - \sum_{i=1}^n x_i) \ln(1-\theta) \end{aligned}$$

# get first derivative

$$\frac{\partial \ell(\theta | s)}{\partial \theta} = \sum_{i=1}^n x_i \cdot \frac{1}{\theta} - (n - \sum_{i=1}^n x_i) \cdot \frac{1}{1-\theta} \quad \nexists \text{ score}$$

# set it to zero to get score equation

$$\frac{\partial \ell(\theta | s)}{\partial \theta} = \sum_{i=1}^n x_i \cdot \frac{1}{\theta} - (n - \sum_{i=1}^n x_i) \cdot \frac{1}{1-\theta} = 0$$

$$\Rightarrow \frac{\sum x_i}{\theta} = \frac{n - \sum x_i}{1-\theta}$$

$$\Rightarrow \frac{1-\theta}{\theta} = \frac{n - \sum x_i}{\sum x_i}$$

$$\Rightarrow \frac{1}{\theta} - 1 = \frac{n}{\sum x_i} - 1$$

$$\Rightarrow \frac{1}{\theta} = \frac{n}{\sum x_i}$$

$$\sum x_i = n \theta$$

$$\hat{\theta} = \frac{\sum x_i}{n} = \bar{x}$$

# check second derivative then you can say its MLE

Ex/  
 $\theta \rightarrow \bar{x}$  MLE     $\nexists$  since one to one  
 $\omega^2 \rightarrow \bar{x}^2$       or  $Ber(\theta), [0, 1]$

$$\Theta^2 \sim \bar{X}^2$$

6.2.5 Suppose  $(x_1, \dots, x_n)$

$$x \sim \text{Uniform}[0, \theta]$$

$$f_\theta(x) = \frac{1}{\theta} I[x \in [0, \theta]]$$

$$x_1, x_2, \dots, x_n$$

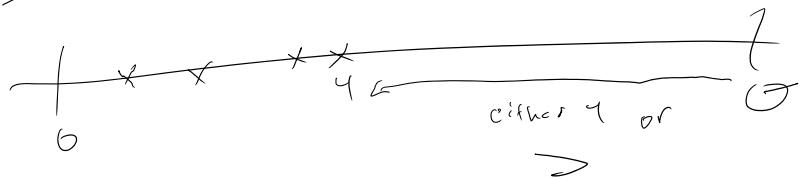
$$\begin{aligned} L(\theta | s) &= \frac{1}{\theta} \cdot \frac{1}{\theta} \cdots \frac{1}{\theta} \\ &= \frac{1}{\theta^n} \end{aligned}$$

# Differentiation part doesn't work

$\frac{1}{\theta^n} \rightarrow$  this is maximized when bottom is minimized.

$\therefore \min \theta \rightarrow \max L(\theta)$  is maximum

Ex assume 1, 2, 3, 4



minimum  $\theta = \max$  value in sample

$$\hat{\theta} = \max(s)$$

Midterm: November 3<sup>rd</sup>

1-4 pm

TC #130

Recap:

$$\text{MLE} \quad \frac{\partial \ln L}{\partial \theta} = 0 \Rightarrow \hat{\theta}$$

$$\boxed{\exists} \quad x \sim \text{Bern}(\theta)$$

$$S_1 = (1, 0, 1, 1, 0, 0, 1) \quad 1 \Rightarrow H$$

$$s_2 = (0, 1, 0, 1, 0, 1, 1)$$

$$L(\theta | s_1) = \theta(1-\theta)^{n-s} (1-\theta)^{s-1} \theta^{n-s} \quad L(\theta, s_2) = (1-\theta)^s \theta^{n-s} (1-\theta)^{s-1} \\ = (1-\theta)^s \theta^n = (1-\theta)^s \theta^{n-1}$$

Ex/  $n = 7$   
# of heads = 4

$$\text{sol} \quad (\underline{\underline{z}}) \Theta \underbrace{(-\Theta)^3}_{L(\Theta(s))}$$

All 3 give you the same likelihood functions, because

$$\sum x_i = 4, \text{ sum of random variables is } 4.$$

sufficient categories  	insufficient stuff  
---	--

$$\text{Inference space}$$

$$T(s_1) = T(s_2)$$

$$\Rightarrow L(s_1) = L(s_2)$$

may or may not have constant in the front. Constant doesn't matter when maximizing it.

$T = \text{sum of all your sample.}$

$$f_0(s) = h(s) * g_0(T(s))$$

↓      ↓  
free of  $\theta$        $\theta, T(s)$

only  $h(s)$       different state

in general       $\theta \sim \theta_0$

$$\text{Ex} \quad \begin{array}{c} u \\ (\underbrace{\quad}_{\text{NS}}) \Theta \end{array} \quad \begin{array}{c} u \\ (\underbrace{\quad}_{\text{NS}}) (-\Theta)^3 \end{array}$$

In general

$$\Theta^{\Sigma^x} \quad \begin{array}{c} u = \Sigma^x \\ (-\Theta)^{\underbrace{\quad}_{\text{NS}}} \end{array}$$

Ex  $x \sim \text{Exp}(0)$

$x_1, x_2, \dots, x_n$

$$\begin{aligned} f_{\theta}(s) &= \partial e^{-\theta x_1} \cdot \partial e^{-\theta x_2} \cdots \partial e^{-\theta x_n} \\ &= \partial e^{n - \theta \sum x_i} \\ \Rightarrow h(s) &= 1 \quad \text{because both terms are of } \theta \end{aligned}$$

$\sum x_i$  is sufficient if statistic

$$= 0 \sum x_i \quad \begin{matrix} \text{sufficient} \\ \text{of statistic} \end{matrix}$$

$$= (-\theta)^e \underbrace{w}_{w(\leq)} \underbrace{g_\theta}_{g_\theta(\cdot)}$$

Ex  $x \sim \text{Pois}(\theta)$

$$S = x_1, x_2, \dots, x_n$$

$$f_{\theta}(s) = \frac{e^{-\theta} \theta^{x_1}}{x_1!} \times \frac{e^{-\theta} \theta^{x_2}}{x_2!} \times \dots \times \frac{e^{-\theta} \theta^{x_n}}{x_n!}$$

likelihood function

$$= \frac{e^{-n\theta} \theta^{\sum x_i}}{\prod_{i=1}^n x_i!}$$

$$L(s) = \frac{1}{\prod_{i=1}^n x_i!}$$

$$g_{\theta} = (e^{-n\theta} \theta^{\sum x_i})$$

$$T(s) = \sum x_i$$

if  $[2, 3, 4, 2, 1]$  is # of accidents on  $\text{H01}$   
and follows  $\text{Po}(\theta)$  then

$$\frac{e^{-n\theta} \theta^{\sum x_i}}{\prod_{i=1}^n x_i!} = e^{-5\theta} \theta^{12}$$

take ln

, don't need to know  $\theta$ , you just need the sum

, sum is sufficient of data

If from likelihood you can get simplification of  $\times$  but  
not sample.

$T(s)$  can be reduced back from your likelihood  
is called a sufficient statistic.

Sufficient means data reduction.

$T(>)$

↑  
minimal sufficient  
that means you can't reduce further

$$\begin{aligned} \text{if } s_1 &= [x_1, x_2, \dots, x_n] \\ s_2 &= [y_1, y_2, \dots, y_n] \\ \frac{\partial \mathcal{L}}{\partial \theta} - \frac{\partial \mathcal{L}_Y}{\partial \theta} &= \frac{\partial}{\partial \theta} (\mathcal{L}_X - \mathcal{L}_Y) \\ &\Rightarrow \sum x = \sum y \\ &\Rightarrow \text{minimal statistic} \end{aligned}$$

## Likelihood Review

For some data  $S = (x_1, \dots, x_n)$  observed which follows a joint distribution  $f_\theta(s)$ , write likelihood

$$L(\theta|S) = f_\theta(S) = \prod_{i=1}^n f_\theta(x_i) \text{ when } x_i \text{ iid}$$

One of the things we can do with likelihood is sufficient statistic.

Sufficient statistic of  $\theta$  ( $T(S)$ )

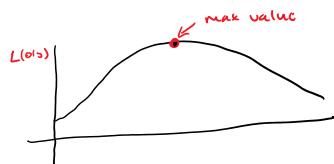
- ↳ To find use factorization theorem.
  - This helps identify  $T(S) \leftarrow$  sufficient statistic of  $\theta$

$$f_\theta(S) = h(S) g_\theta(T(S))$$

## Maximum Likelihood Estimation (MLE)

- MLE( $\theta$ ) value that maximizes  $L(\theta|S)$  namely

$$L(\hat{\theta}|S) \geq L(\theta|S) \quad \forall \theta \in \Theta$$



### Steps for finding MLE

1. write out  $L(\theta|S) \leftarrow$  likelihood function
2. take  $\ln$  of  $L(\theta|S) \leftarrow$  ln likelihood function
3. take the derivative of  $\ln(L(\theta|S))$  with respect to  $\theta$
4. set to 0 and solve for  $\theta$

Review chapter 4.6 for background

## Inference based on MLE

- for  $T(S)$  an estimate for  $\psi(\theta)$ , measure "closeness" using

$$\text{MSE}_\theta(T(S)) = E_\theta [T(S) - \psi(\theta)]^2$$

↑ mean squared error  
 ↑ sufficient statistic of  $\theta$   
 ↑ expectation

= average squared distance between  $T$  and unknown  $\psi(\theta)$

Recall, if  $T(S) = \bar{x}$  then sampling distribution of  $T(S)$  when  $S$  is  $N(\mu, \sigma^2)$  is  $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$

- Easier MSE formula

$$\text{MSE}_\theta(T(S)) = \text{Var}_\theta(T(S)) + [E_\theta(T(S)) - \psi(\theta)]^2$$

### Proof

$$\begin{aligned} & E_\theta [(T(S) - \psi(\theta))^2] \\ &= E_\theta [[T(S) - E_\theta(T(S))]^2 + [E_\theta(T(S)) - \psi(\theta)]^2], \quad \text{cheap trick prof gonna use, add and subtract same} \end{aligned}$$

\* update prof said look at the

0

value. Figure out using  
before the midterm!

ending answer  
and used only for practice

$$\begin{aligned}
 &= E[(T(\omega) - E[T(\omega)])^2] + 2E[(T(\omega) - E[T(\omega)])(E[T(\omega)] - \psi(\omega))] + (E[T(\omega)] - \psi(\omega))^2 \\
 &= \text{Var}(T(\omega)) + (E[T(\omega)] - \psi(\omega))^2 + \underbrace{\text{bias}}_{\text{bias}}
 \end{aligned}$$

expand it  
 $E(1c) = 0$   
 where  $c$  is  
 a constant

### Bias of $T$ when $E_\theta[T]$ exists

$$\text{Bias}(T) = E_\theta[T(\omega)] - \psi(\omega) = \begin{cases} 0 & \text{then } T \text{ is unbiased} \\ \text{otherwise } T \text{ is biased} \end{cases}$$

or equivalently, if  $E[T(\omega)] = \psi(\omega)$  then  $T$  is unbiased.

### Ex 6.3.1

$(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$  with unknown  $\mu$  and known  $\sigma^2$ . Want to calculate the MSE

$$\begin{aligned}
 \text{MSE}(\hat{\theta}) &= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \mu)^2 \\
 &= \text{Var}\left(\frac{1}{n} \sum x_i\right) + (E[\bar{x}] - \mu)^2 \\
 &= \frac{1}{n^2} \sum \text{Var}(x_i) + \left(\frac{1}{n} \sum E(x_i) - \mu\right)^2 \\
 &= \frac{1}{n^2} (n\sigma^2) + \left(\frac{1}{n} n\mu - \mu\right)^2 \\
 &= \frac{\sigma^2}{n} + (0)^2
 \end{aligned}$$

Previous section  
 MLE of  $\mu$   $\hat{\theta} = \bar{x}$  is  $\bar{x}$   
 mean  
 Recall from U.b  
 $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$

This implies its unbiased

### Example 6.3.4

Now  $(x_1, \dots, x_n) \sim N(\mu, \frac{\sigma^2}{n})$ ,  $(\mu, \sigma^2)$  are unknown

We have MLE for  $(\mu, \sigma^2)$  is  $(\bar{x}, \frac{n-1}{n} s^2)$

still  $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$  so  $E[\bar{x}] = \mu$

$$\text{MSE}(\bar{x}) = \text{Var}(\bar{x}) + \text{Bias}(\bar{x})$$

$$= \text{Var}(\bar{x}) + 0$$

$$= \frac{\hat{\sigma}^2}{n}$$

$$= \frac{1}{n^2} (n-1)s^2 \approx \frac{s^2}{n} \leftarrow \text{for large values of } n, \text{ we can approximate it to } \frac{s^2}{n} \text{ because as } n \uparrow \frac{n-1}{n} \text{ becomes more insignificant}$$

### Bias

Variance of statistic is often called standard error =  $\sqrt{\text{Var}}$

### Confidence Intervals

\* Copy slide 7, 8, 9  $\rightarrow$  she read off the slide really fast.

- "γ" is referred to as the confidence level of the interval  
 $\uparrow$   
 gamma

Example 6.3.6  $(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$ ,  $\mu$  unknown and  $\sigma^2$  known

# want to develop an interval for  $\mu$  using our data that isn't too wide but has high confidence.

# Start from likelihood

$$L(\mu | x) \propto \exp \left\{ \frac{-n}{2\sigma^2} (\bar{x} - \mu)^2 \right\}$$

proportional

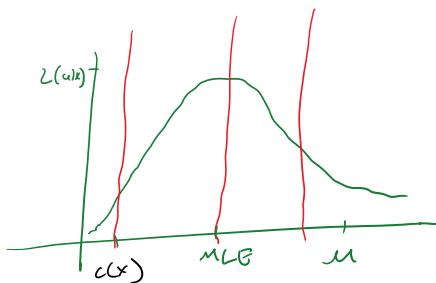
$\left\{ \begin{array}{l} \text{IR - high confidence} \\ \text{- very wide} \end{array} \right.$

Define  $C(x_1, \dots, x_n)$  such that  $\mu \in C(x)$  and

$\uparrow$   
 This is a  
 "C"

$\uparrow$   
 This is a  
 "C"

$$L(\mu_2 | x) \geq L(\mu_1 | x) \text{ then } \mu_2 \in C(x)$$



Thus, we can write

$$\begin{aligned} C(x) &= \{ \mu : L(\mu | x_1, \dots, x_n) \geq k(x_1, \dots, x_n) \} \\ &= \{ \mu : \exp \left\{ \frac{-n}{2\sigma^2} (\bar{x} - \mu)^2 \right\} \geq k(x_1, \dots, x_n) \} \\ &= \{ \mu : \frac{n}{2\sigma^2} (\bar{x} - \mu)^2 \leq -2 \frac{\sigma^2}{n} \ln k(x_1, \dots, x_n) \} \\ &= \{ \mu : \bar{x} - k^*(x) \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + k^*(x) \frac{\sigma}{\sqrt{n}} \} \end{aligned}$$

where  $k^*(x) = \sqrt{-2 \ln k(x)}$

Now we need to choose  $k(x)$  or equivalently  $k^*(x)$   
 we know

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

can rewrite  $C(x)$  as a probabilistic statement by

$$\begin{aligned} \gamma &\leq P(\mu \in C(x)) = P\left(\bar{x} - k^* \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + k^* \frac{\sigma}{\sqrt{n}}\right) \\ &= P(-k^* < \bar{x} - \mu < k^*) \end{aligned}$$

- Visually



$$= P\left(-k^* \leq \frac{\bar{x} - \mu}{\frac{\sigma_0}{\sqrt{n}}} \leq k^*\right)$$

$$= P\left(\left|\frac{\bar{x} - \mu}{\frac{\sigma_0}{\sqrt{n}}}\right| \leq k^*\right)$$

$$\Rightarrow P(|Z| \leq k^*) \geq \gamma$$

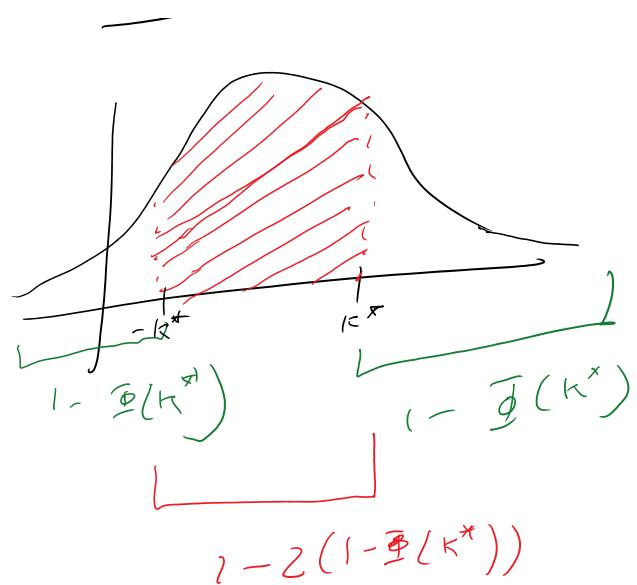
$\uparrow$   
 $N(0, 1)$

$$\Rightarrow 1 - Z(1 - \Phi(k^*)) \geq \gamma$$

$$\Rightarrow \Phi(k^*) = \frac{1 + \gamma}{2}$$

$$\Rightarrow k^* = Z\left(\frac{1 + \gamma}{2}\right)$$

from  $N(0, 1)$  table



$$C(S) = \left[ \bar{x} - Z\left(\frac{1 + \gamma}{2}\right) \frac{\sigma_0}{\sqrt{n}}, \bar{x} + Z\left(\frac{1 + \gamma}{2}\right) \frac{\sigma_0}{\sqrt{n}} \right]$$

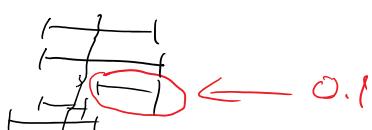
margin of error

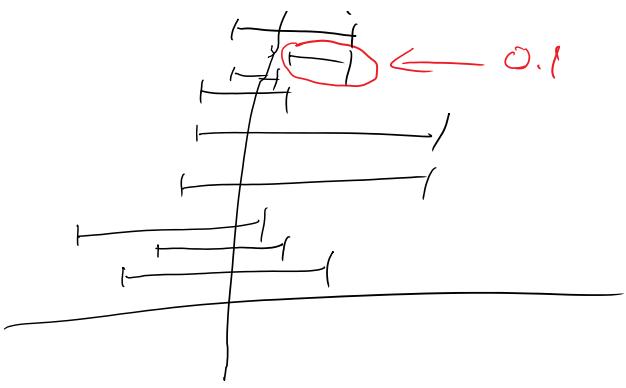
\* based on the likelihood construction, proper CI interpretation is that for repeated sampling with CI constructed each time, then 100γ% of those CIs contain true  $\psi(\theta)$

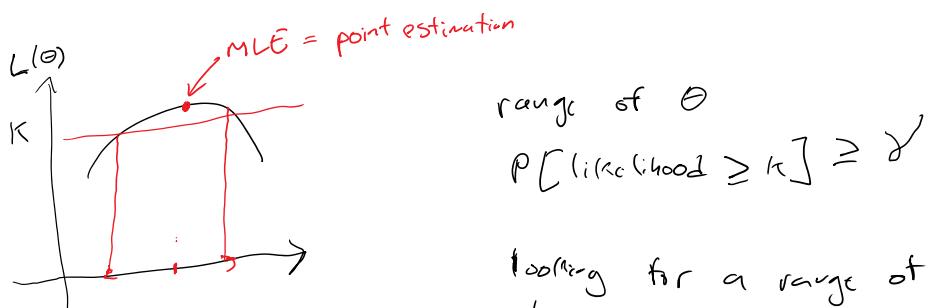
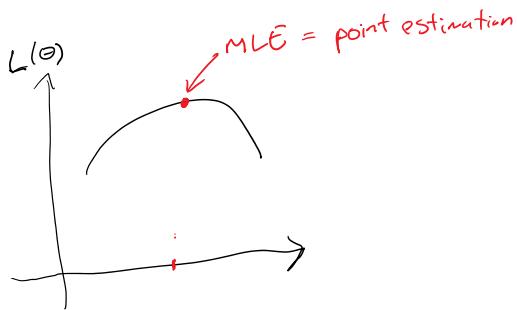
(?) ↑

write this down on exam if you want marks

$$\bar{x} \quad \gamma = 0.9$$







$$\begin{aligned} P[\text{likelihood} \geq k] &\geq \gamma \\ P[\log \text{likelihood} \geq k^*] &\geq \gamma \end{aligned}$$

(likelihood interval for  
any distribution)

looking for a range of  $\theta$  where there is a gamma chance it will be above some threshold.

$\Rightarrow \sim N(\mu, \frac{\sigma^2}{n})$ , is  $-\left(\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}}\right)^2$  the log likelihood? yes.

$$b) P\left[-\left(\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}}\right)^2 \geq k^*\right] \geq \gamma$$

$$P\left[\left(\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}}\right)^2 \leq k^{**}\right] \geq \gamma$$

constant with neg

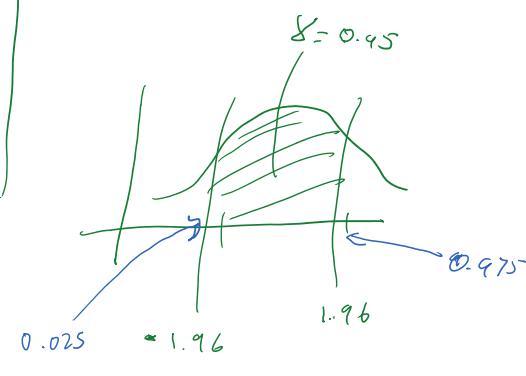
$$P\left[-\sqrt{k^{**}} \leq \frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} \leq \sqrt{k^{**}}\right] \geq \gamma$$

Only cause of U/T

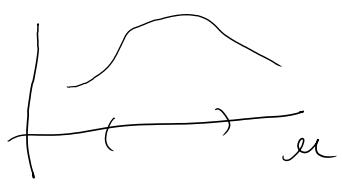
$$\Rightarrow P\left[-\frac{z_{1-\gamma}}{2} \leq \frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} \leq \frac{z_{1-\gamma}}{2}\right] \geq \gamma$$

Central limit theor

$$\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$



$$q_{\text{norm}}(0.025) = -1.96$$



vn

↓

interval for  $\mu$  since  
 $\mu$  is unknown

$$q_{\text{norm}}(0.025) = -1.96$$

$$q_{\text{norm}}(0.975) = 1.96$$

if  $\delta = 0.9$ .

$$\left( \frac{\bar{x} + \delta}{2} \right) = \frac{1.9}{2} = 0.95$$

# Transform the equation to have  $\mu$  on one side.

$$P\left[\frac{\bar{x} - z_{1+\alpha}}{\frac{\delta}{\sqrt{n}}} \leq \mu \leq \frac{\bar{x} + z_{1+\alpha}}{\frac{\delta}{\sqrt{n}}}\right] \Rightarrow \left(\bar{x} \pm z_{1+\alpha} \frac{\delta}{\sqrt{n}}\right)$$

Ex  $\bar{y} \sim N(\mu, \sigma^2=6)$ ,  $n=9$   $\bar{x}=5$ , calc 95% confidence interval

$\sigma^2$  is known ✓

what gamma? 0.95

what quantile should I look at?  $\frac{1+\alpha}{2} = 0.975 = \Phi^{-1}(0.95)$

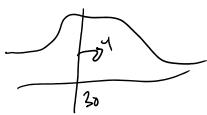
$$\Rightarrow 5 \pm 1.96 \cdot \frac{4}{\sqrt{9}} \Rightarrow \begin{cases} 7.61 \\ 2.38 \end{cases} \quad \begin{matrix} \rightarrow \text{interpretation: lower and upper bound} \\ \text{contain the true } \mu \end{matrix}$$

$$P[2.38 \leq \mu \leq 7.61] = 0.95 \quad \times$$

Likelihood is not a pdf, just a function.

$\bar{x}$  is a variable, because

lets say

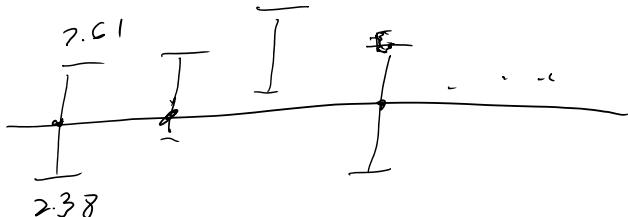


Sample 1	
$x_1 - x_3$	32 - 30
30	31 29   36

Sample 2	
$x_2 - x_3$	32 - 28
32	28 25   $\mu_2$

it changes every sample so its a variable  
confidence interval is a variable.

$$P[2.38 \leq \mu \leq 7.61] = 0.95 \quad \times$$



it looks at the whole thing not just one case

correct

$$P[lb \leq \mu \leq ub] = 0.95 \quad \checkmark$$

$$\underline{Ex} \quad N(\mu, \sigma^2 \rightarrow \text{unknown})$$

$$P\left[-k^{**} \leq \frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} \leq k^{**}\right] \geq \gamma \quad , \text{ we can't use C/T}$$

since  $\sigma^2$  is unknown

$$\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

$\rightarrow \frac{(n-1)\sigma^2}{\sigma^2} = \chi^2_{(n-1)}$  (chi-square distribution) ← don't worry you will learn this later from  
 parameter degrees of freedom ( $\delta f$ )  
 no parameter in standard normal

$$N(0, 1) \rightarrow \text{parameter free} \quad \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$\frac{\text{std Normal}}{\sqrt{\frac{\bar{x}^2}{df}}} \sim t \quad \text{will not be tested}$$

$$\frac{\bar{x}-\mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\bar{x}-\mu}{\frac{s}{\sqrt{n}}} \sim t - \text{distribution}$$

*magic*

$$\sqrt{\frac{(n-1)s^2}{\sigma^2}} = \sqrt{\frac{(n-1)}{(n-1)}}$$

when you know sigma use standard normal

when you don't know use  $t$ -distribution.

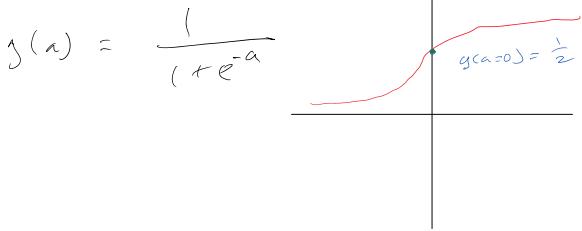


$$\begin{aligned}
 p(c_1|x) &= \frac{p(x|c_1)p(c_1)}{p(x)} \\
 &= \frac{p(x|c_1)p(c_1)}{p(x|c_1)p(c_1) + p(x|c_2)p(c_2)}, \quad \text{divide top and bottom by the top} \\
 &= \frac{1}{1 + \frac{p(x|c_2)p(c_2)}{p(x|c_1)p(c_1)}} \quad \left. \begin{array}{l} \text{important form} \\ \text{=} p(c_1|x) \text{ posterior prob} \\ \text{distribution of } c_1 \text{ given } x \end{array} \right\}
 \end{aligned}$$

$\Rightarrow$  This form tells us that the relative prob matters

$$\text{Recall} \\
 d(x) = \log \left( \frac{p(x|c_1)p(c_1)}{p(x|c_2)p(c_2)} \right) \Rightarrow \text{decision boundary } a(x) = 0$$

$$\begin{aligned}
 p(c_1|x) &= \frac{1}{1 + \frac{p(x|c_2)p(c_2)}{p(x|c_1)p(c_1)}} \\
 &= \left( \frac{1}{1 + e^{-a(x)}} \right) \equiv g(a(x)) \quad \text{sigmoid function}
 \end{aligned}$$



## Logistic Regression

Classification technique not a classification.

Suppose  $a(x) = \bar{x}^T \bar{w} + b \leftarrow$  suppose linear  
 $= \bar{x}^T \bar{w} \leftarrow$  ignore bias by absorbing  
 your collapse in your weight vector.

$$p(c_1|x) = \frac{1}{1 + e^{-\bar{w}^T \bar{x}}} = g(\bar{w}^T \bar{x})$$

## Learning

Optimization on likelihood  
 data  $\{(x_i, y_i)\}_{i=1}^N$   $y \in \{0, 1\}$

$$x = \{\bar{x}_i\} \quad y = \{y_i\}$$

$$\begin{aligned}
 p(x, y|\bar{w}) &= p(y|\bar{w}^T \bar{x})p(x|\bar{w}) \\
 &\quad \cancel{p(x)} \text{ because data doesn't depend on } w, y \text{ depends on } w \text{ though} \\
 &\Rightarrow x \perp w
 \end{aligned}$$

$$\propto p(y|\bar{w}^T \bar{x})p(x)$$

$$p(y|x, \bar{w})$$

$$p_i \equiv p(y_i=c_1 | x_i, \bar{w})$$

$$1-p_i \equiv p(y_i=c_2 | x_i, \bar{w})$$

$$p(y|x, \bar{w}) = \prod_{i:y_i=c_1} p_i \prod_{i:y_i=c_2} (1-p_i)$$

$$\begin{aligned} &= \prod_{i=1}^N p_i^{y_i} (1-p_i)^{1-y_i} \leftarrow \text{cross entropy loss function} \\ &= \prod_{i=1}^N g(\bar{x}_i^\top \bar{w})^{y_i} (1-g(\bar{x}_i^\top \bar{w}))^{1-y_i} \end{aligned}$$

$$E(w) = -\log p(y|x, \bar{w}) = \sum_{i=1}^N y_i \log(g(\bar{x}_i^\top \bar{w})) + (1-y_i) \log(1-g(\bar{x}_i^\top \bar{w})) \leftarrow \text{energy function}$$

$$p_i \equiv p(y_i=c_1 | x_i, \bar{w}) = g(\bar{x}_i^\top \bar{w})$$

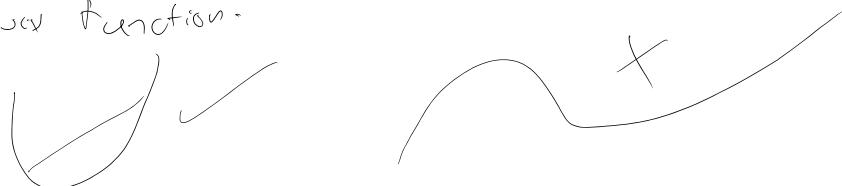
$$1-p_i \equiv p(y_i=c_2 | x_i, \bar{w}) = 1-g(\bar{x}_i^\top \bar{w})$$

iterative optimization technique, take gradient with respect to  $w$ , stop when slope of the loss - Gradient descent.

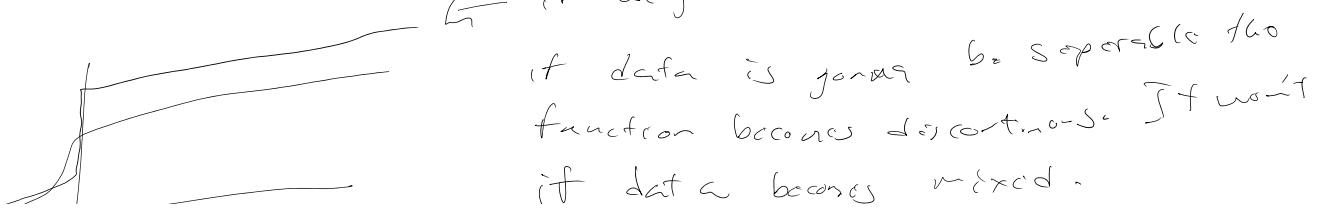
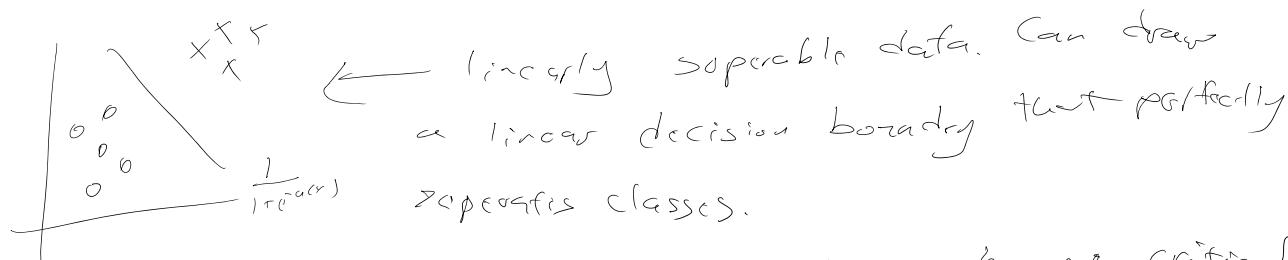
A2

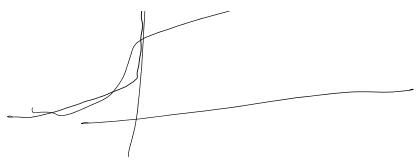
get two gradient descent code is provided  
# find derivative of the energy function.

convex function:



Unique minimum:





function becomes discontinuous  
if data becomes mixed.

∴ add a regularizer

$$E(\omega) = -\log p(x|y|\bar{\omega}) = \sum_{i=1}^n y_i \log(\bar{x}^T \bar{\omega}) + (1-y_i) \log(1 - g(\bar{x}_i^T \bar{\omega})) + \frac{1}{2} \bar{x}^T \bar{\omega}^T \bar{\omega}$$

Regularizer

## Review

~~Bayesian statistics~~ = parameters of joint distribution

data  $\leftarrow x_1, x_2, \dots, x_n \sim \text{Bern}(\theta)$

$\theta \sim \text{Beta}(\alpha, \beta)$   
prior

$$\pi(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$\pi(\theta | \text{data}) \propto \text{likelihood} * \text{prior}$

$$\underset{\text{sample}}{\downarrow} \propto \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\propto \theta^{\sum x_i + \alpha - 1} (1-\theta)^{n - \sum x_i + \beta - 1}$$

$$\theta | \text{data} \sim \text{Beta}(\sum x_i + \alpha, n - \sum x_i + \beta)$$

Ex if  $n=3$ ,  $\sum x_i = 2$ ,  $\alpha = 1$ ,  $\beta = 3$  and  $\theta \sim \text{Beta}(\alpha=2, \beta=2)$   
then what is  $\theta | \text{data} \sim \text{Beta}(5, 3)$ ?

$$\begin{aligned} & \text{Beta}(5, 3) \\ & = \text{Beta}(2+3, 3+2) \\ & = \text{Beta}(5, 5) \leftarrow \text{posterior, you know the parameters} \end{aligned}$$

mean of  $\text{Beta}(5, 3)$  is

$$E(\theta) = \frac{\alpha}{\alpha + \beta} = \frac{5}{5+3} = \frac{5}{8}$$

Calculating mode for beta ~~possible mid-term question~~

posterior is a combination prior and theta.

$$\theta | \text{data} \sim \text{Beta}(\sum x_i + \alpha, n - \sum x_i + \beta)$$

$$f_{\theta | \text{data}}(\theta) = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(\sum x_i + \alpha) \Gamma(n-\sum x_i + \beta)} \cdot \theta^{\sum x_i + \alpha - 1} (1-\theta)^{n - \sum x_i + \beta - 1}$$

$$\ln f_{\theta | \text{data}}(\theta) \propto (\sum x_i + \alpha - 1) \ln \theta + (n - \sum x_i + \beta - 1) \ln(1-\theta)$$

$$\frac{\partial \ln f_{\theta | \text{data}}(\theta)}{\partial \theta} \propto \frac{\sum x_i + \alpha - 1}{\theta} - \frac{(n - \sum x_i + \beta - 1)}{1-\theta}$$

$$\Rightarrow \frac{\sum x_i + \alpha - 1}{\theta} - \frac{n - \sum x_i + \beta - 1}{1-\theta} = 0$$

$$\Rightarrow \frac{\sum x_i + \alpha - 1}{\theta} = \frac{n - \sum x_i + \beta - 1}{1-\theta}$$

$$\frac{1-\theta}{\theta} = \frac{n - \sum x_i + \beta - 1}{\sum x_i + \alpha - 1}$$

$$\Rightarrow \frac{n + \alpha + \beta - 2}{\sum x_i + \alpha - 1}$$

For this example  
 $x_i \sim \text{Bern}(\theta)$

$\pi(\theta) = \text{Beta} \leftarrow \text{start with some distribution}$

post = Beta  $\leftarrow$  same distribution  
 $\Rightarrow$  "conjugate"

if prior distribution is post distribution then

it is a conjugate.

Ex  $x \sim \text{Pois}(x) / x > 0$   
 $\lambda \sim \text{Gamma}(\alpha, \beta)$   
 $x | \lambda \sim \text{Gamma}(\alpha^*, \beta^*)$

$\lambda \sim \text{Uniform}[0, \infty]$   
is a non informative prior because  
 $\lambda > 0$ . It's a good prior because information does not = 1.

if you have conjugate priors, we don't have to do  
 $m(s)$ , makes calculations easier.

$$\pi(\theta) \approx \text{Uniform}[0, 1] \leftarrow \text{non-informative prior}$$

if prior is not a valid pdf, it is an improper prior

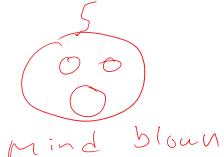
$$\text{Ex/ } \pi(\theta) \sim \text{Beta}(\alpha=1, \beta=1)$$

$\text{Beta}(\alpha, \beta)$  then

$$\text{pdf } \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\begin{aligned} &\Rightarrow \frac{\Gamma(2)}{\Gamma(1)\Gamma(1)} x^{1-1} (1-x)^{1-1} \\ &= \Gamma(2) \\ &= ((-1)!) \\ &= 1 \end{aligned} \quad \left. \begin{array}{l} \Gamma(3.5) = 2.5 \times 1.5 \times 0.5 \\ \Gamma(x+1) = x! \\ \Gamma(0.5) = \sqrt{\pi} \end{array} \right\}$$

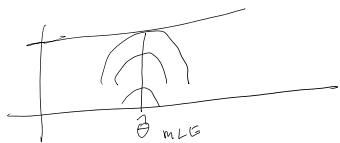
Note:  $\text{Beta}(1, 1) = \text{Uniform}[0, 1]$



$$\text{model: } \frac{\sum x_i}{n} = \bar{x} \quad \left| \begin{array}{l} \text{To know this} \\ \text{MLE} \Leftarrow \text{Bm}(x) \Leftarrow \hat{\theta} = \bar{x} \end{array} \right.$$

post  $\propto$  Likelihood

$$\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$



if non informative prior, whatever MLE is the mode



only if  
non informative  
prior!

Both curves look the same but will be shifted up and down. But the maximum is the same.  
This only the case when the prior is uninformative.

$\sum x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$  known

$$\mu \sim N(\mu_0, \sigma_0^2)$$

$$(\mu | s) \sim N(\mu^*, \sigma^*)$$

posterior

$$\exp \left[ -\frac{1}{2} \left( \frac{1}{\sigma^2} (\mu - \mu^*)^2 \right) \right]$$

$$\mu^* = \frac{\frac{1}{\sigma^2} \mu_0 + \frac{n}{\sigma^2} \bar{x}}{\frac{1}{\sigma^2} + \frac{n}{\sigma^2}} = \frac{\frac{1}{\sigma_0^2} \mu_0}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \mu_0 + \frac{\frac{n}{\sigma^2} \bar{x}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}}$$

$$\Rightarrow w_1 + w_2 = 1$$

$$= w_1 x + w_2 y$$

↳ weighted average

$$\text{avg} = \frac{3+5}{2} = \frac{1}{2}(3) + \frac{1}{2}(5)$$

$$w_1 = \frac{3}{3+5}$$

$$w_2 = \frac{5}{3+5}$$

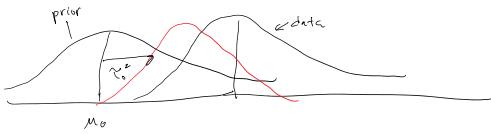
$$\text{avg} = \frac{3+5}{3+5} = \frac{1}{2}(3) + \frac{1}{2}(5)$$

$$= \frac{1}{3}(3) + \frac{1}{3}(5)$$

↓ in this case 3 is pulling the mean by  $\frac{2}{3}$  and 5 is pulling by  $\frac{1}{3}$



posterior mean  $\rightarrow$  weighted average of the prior mean ( $\mu_0$ ) and sample mean ( $\bar{x}$ )



• ← red is posterior  
right between prior  
and sample mean.  
end result.

If  $\sigma_0^2 = \sigma^2 = 1$ , then

$$\frac{\frac{1}{\sigma_0^2}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \mu_0 + \frac{\frac{n}{\sigma^2}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \bar{x} = \frac{1}{1+n} \mu_0 + \frac{n}{1+n} \bar{x}$$

$$(x_1, x_2, \dots, x_n) \text{ mean} = \frac{\sum x_i}{n}$$

$$(x_1, x_2, \dots, x_n, \mu_0) \text{ mean} = \frac{\sum x_i + n\mu_0}{n+1}$$

if  $\sigma_0^2 \rightarrow \infty$

$$\begin{aligned} \frac{\frac{1}{\sigma_0^2}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \mu_0 + \frac{\frac{n}{\sigma^2}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \bar{x} &= \frac{0}{0 + n} + \frac{n}{0 + n} \bar{x} \\ &= \bar{x} \end{aligned}$$

⇒ if variance of Normal is really big then if it is not a informative prior → whatever data is posterior is posterior will be the same one.

(case) you should do: something  $\mu$  is known but  $\sigma^2$  is unknown

$$x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$$

↑ unknown  
↓ known

prior /  $\frac{1}{\sigma^2} \sim \text{Gamma}(\alpha, \beta)$

$(\frac{1}{\sigma^2}) \leftarrow \text{post}$

$$\frac{B^\alpha}{\Gamma(\alpha)} \left( \frac{1}{\sigma^2} \right)^{\alpha-1} e^{-\beta \cdot \frac{1}{\sigma^2}}$$

normal part or is always on the bottom makes calculation easier

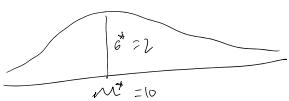
$\frac{B^\alpha}{\Gamma(\alpha)} \left( \frac{1}{\sigma^2} \right)^{\alpha-1} e^{-\beta \cdot \frac{1}{\sigma^2}}$

Recall

$$G(\alpha, \beta) = \frac{B^\alpha}{\Gamma(\alpha)} \times \alpha^{-1} e^{-\beta \alpha}$$

$$N \sim \left( \frac{1}{2\sigma^2} \right)^\frac{1}{2} \exp \left[ -\frac{1}{2\sigma^2} (y - \mu)^2 \right]$$

what is post? ← exercise. Pls do this f'ham.



normal distribution mean = mode = median

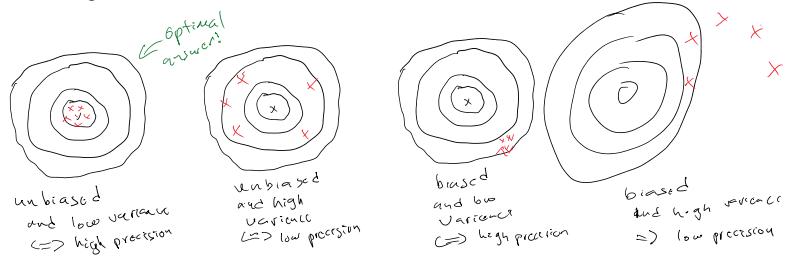


→ "credible region" instead of a confidence interval.

$$\text{MSE} = \underline{\text{Bias}^2} + \underline{\text{Var}}$$

← knows this f'ham  
its on the midterm

imagine archery



$$\frac{1}{\text{var}} = \text{precision}$$

Bias = shoot repeatedly, you are averaging the target.