

# STAT 241

# Personal Notes

---

\* These notes are strictly my own interpretation  
of the course materials.

Marcus Chan

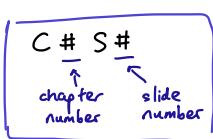
---

Taught by Michael Wallace  
UW CS '25

---



# Chapter 1: Introduction to Statistical Science



💡 Statistical science is the science of "empirical studies".

## EMPIRICAL STUDY (CIS24)

- 💡 An "empirical study" is one where we learn by observation and/or experimentation.
- 💡 Note these involve uncertainty - repeated experiments generate different results.
- 💡 But we model these uncertainties using probability models.

## UNIT (CIS25)

💡 A "unit" is an individual which we can take measurement(s).

## POPULATION (CIS26)

- 💡 A "population" is a collection of units.
  - eg - all current UW undergrad students
  - all donuts in Tim Hortons right now
- \* note: we need to be precise when defining populations or any other terms!
  - eg if we said "all UW students" this is ambiguous, since it might include grads, alumni, etc

## PROCESS (CIS27)

- 💡 A "process" is a system by which units are produced.
  - eg - hits on a particular website are units in a process
  - claims made by insurance policy holders are units in a process
- 💡 Note that although populations & processes are collections of units:
  - ① Populations are "static" (defined at one point in time), but
  - ② Processes usually occur over time.

## VARIATES (CIS32)

💡 "Variates" are characteristics of the units.

\* we usually represent these by letters  $x, y$  &  $z$ .

## CONTINUOUS VARIATES (CIS33)

- 💡 "Continuous variates" are those that can be measured (at least theoretically) to an infinite degree of accuracy.
  - eg height, weight, lifetime of a fuse, etc

## DISCRETE VARIATES (CIS33)

- 💡 "Discrete variates" are those that can only take finitely or countably many values.

eg # of car accidents on a certain stretch of highway / yr, etc.

- 💡 Note that depending on how we measure a continuous variate, it may become discrete.
  - eg if we measure weight w/ a scale that only goes to 2dp, the resulting variate is discrete!

- 💡 Ultimately the distinction affects
  - ① our assumptions of the data; and
  - ② the probability models we use
    - for discrete variates, we usually use discrete prob models (eg Poisson)
    - for cts variates, we usually use cts prob models (eg Gaussian)
    - but there are exceptions. (CIS43)

## CATEGORICAL VARIATES (CIS35)

- 💡 "Categorical variates" are those where the units fall into non-numeric categories, without any implied order.
  - eg hair color, university program

## ORDINAL VARIATES (CIS35)

- 💡 "Ordinal variates" are those where an ordering is implied, but not necessarily from a numeric measure.
  - eg strongly disagree, ..., strongly agree;
  - small, medium, large;
  - etc

## COMPLEX VARIATES (CIS37)

- 💡 "Complex variates" are those that are more unusual, and don't fall neatly into the other variate types.

- 💡 eg open-ended responses to a survey question  
We usually need processing to convert these into one of the other types.

eg text processing to convert a tweet's content into "positive", "negative" or "neutral"

## ATTRIBUTES [OF A POPULATION/PROCESS] (CISY8)

"Attributes" of a population/process are functions of a variate which is defined for all units in said population/process.

- eg (STAT 231 asmts) - mean # of completed asmts  
- prop. of asmts subbed in last 24 hrs  
(Kw Humane Society) - prop. of dogs that arrive in good health  
- mean # of owners of dogs in their care

## TYPES OF EMPIRICAL STUDIES (CIS50)

### SAMPLE SURVEY (CIS52)

A "sample survey" is where information is obtained about a finite population by

- ① selecting a "representative" sample of units from the population; and
- ② determining the variates of interest for each unit in the sample.

- eg - poll to predict who will win an election  
- survey of potential consumers to compare products & state their preference (eg Coke vs Pepsi)

### OBSERVATIONAL STUDY (CIS53)

An "observational study" is where information about a population/process is collected without any change to the sampled units' variates.

- eg a study of blood alcohol levels for students at a 8:30am Mon lecture

Usually, the following are true:

Observational	Survey
① Pop" of interest is infinite/conceptual	Pop" is finite/real
② Data collected <u>routinely</u> over time	Data collected <u>once</u>
③ More passive (sit and see)	More "aggressive" (specific questions asked)

\*but these are just guidelines - there are exceptions. (CIS55)

### EXPERIMENTAL STUDY (CIS54)

An "experimental study" is one where the experimenter intervenes and modifies some of the variates for the units in a study.

- eg same example as above, but some students are warned beforehand, whereas some are not.

# DATA SUMMARIES (CIS56)

- These are used for
- the estimation of attributes; and
  - checking fit for a model.

## MEASURES OF CENTRAL TENDENCY / LOCATION (CIS58)

We usually represent our data using the notation  $\{y_1, \dots, y_n\}$ , where each  $y_i \in \mathbb{R}$  and  $n$  is called the "sample size".

We also use lower-case for constants, and upper-case for random variables.

### ORDERED SAMPLE / ORDER STATISTICS (CIS59)

We call the "ordered sample" or "order statistics" of the data to be

$$y_{(1)}, \dots, y_{(n)}$$

where  $y_{(1)} \leq \dots \leq y_{(n)}$ ,  $y_{(1)} = \min\{y_1, \dots, y_n\}$  &  $y_{(n)} = \max\{y_1, \dots, y_n\}$ .

### SAMPLE MEAN/AVERAGE: $\bar{y}$ (CIS58)

The "sample mean", denoted by " $\bar{y}$ ", is equal to

$$\bar{y} := \frac{y_1 + \dots + y_n}{n} = \frac{1}{n} \sum_{i=1}^n y_i$$

\* the keyword "sample" is important!

### SAMPLE MEDIAN: $\hat{m}$ (CIS59)

The "sample median", denoted as " $\hat{m}$ ", is defined by

$$\hat{m} := \begin{cases} y_{(\frac{n+1}{2})}, & n \text{ is odd} \\ \frac{1}{2}(y_{(\frac{n}{2})} + y_{(\frac{n}{2}+1)}), & n \text{ is even} \end{cases}$$

Note that

- In symmetrical distributions,  $\bar{y} \approx \hat{m}$ ;
- In skewed distributions,  $\bar{y} \neq \hat{m}$  (there may be a significant gap between them). (CIS66)

### SAMPLE MODE (CIS61)

The "sample mode" is just the most common value(s) in a set of data.

In this case, the "sample modal class" is the group/class with the highest frequency.

## MEASURES OF VARIABILITY /

### DISPERSION (CIS67)

"Measures of variability" convey how "spread out" the data is.

### ROBUST [MEASURE] (CIS80)

We say a measure is "robust" if it is not significantly affected by extreme values.

e.g. IQR is robust, range is not

### SAMPLE VARIANCE & STANDARD DEVIATION:

#### $s^2, s$ (CIS69)

We define the "sample variance", denoted " $s^2$ ", of the data  $\{y_1, \dots, y_n\}$  to be

$$s^2 := \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n y_i^2 - n\bar{y}^2 \right]$$

The "sample standard deviation", denoted " $s$ ", is just the square root of the sample variance.

### "68-95" RULE FOR GAUSSIAN ESTIMATION

#### (CIS70)

Suppose the data  $\{y_1, \dots, y_n\}$  is from a Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$ . \*  $\mathcal{N}(\mu, \sigma^2) = N(\mu, \sigma^2)$

Then necessarily

① 68% of the sample lies in  $[\bar{y} - s, \bar{y} + s]$ ;

and

② 95% of the sample lies in  $[\bar{y} - 2s, \bar{y} + 2s]$ .

\* this can be verified in R using the code

```
> pnorm(1) - pnorm(-1)
> pnorm(2) - pnorm(-2)
```

### RANGE (CIS73)

The "range" is defined as

$$\text{range} = y_{(n)} - y_{(1)}$$

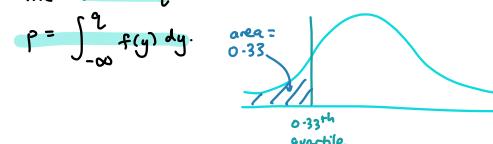
\* the range is very susceptible to outliers!

### QUANTILES & PERCENTILES (CIS74)

The " $p$ th quartile", also called the "(100p)th percentile", is the value such that a fraction  $p$  of the data fall at or below said value.

\* the median is the 50th quartile / 50th percentile.

In other words, the  $p$ th quartile of a distribution is the value  $q$ , such that



We can calculate quartiles in R using the code

```
> quantile(c(y1, ..., yn), p)
```

### QUARTILES: $q(0.25), \hat{m}, q(0.75)$ (CIS79)

The "lower quartile", or "first quartile", denoted by  $q(0.25)$ , is the 25th percentile.

The "upper quartile", or "third quartile", denoted by  $q(0.75)$ , is the 75th percentile.

The "second quartile" is just the median  $\hat{m}$ .

### INTERQUARTILE RANGE / IQR (CIS80)

The "interquartile range" is defined as

$$\text{IQR} = q(0.75) - q(0.25)$$

\* IQR is robust — it is not affected by extreme values.

\* if considering discrete data, the interpretation of IQRs can vary depending whether we consider the "interval" from  $q(0.25)$  to  $q(0.75)$  to be open, semi-open or closed.

# MEASURES OF SHAPE (CIS84)

## SAMPLE SKEWNESS (CIS88)

$\exists_1$  "Sample skewness" measures the asymmetry of the data, and is equal to

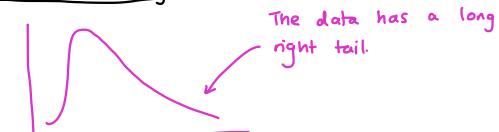
$$\text{sample skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^3}{\left[ \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{3/2}}$$

$\exists_2$  Interpretation of sample skewness's value:

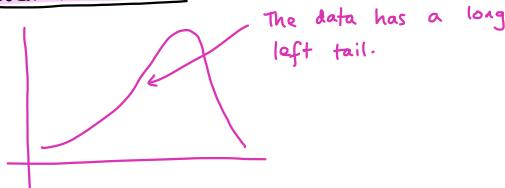
- ① If  $ss = 0 \Rightarrow$  distribution is symmetric; eg Gaussian, uniform



- ② If  $ss > 0 \Rightarrow$  distribution is positively skewed / skewed to the right;



- ③ If  $ss < 0 \Rightarrow$  distribution is negatively skewed / skewed to the left.



## SAMPLE KURTOSIS (CIS96)

$\exists_1$  "Sample kurtosis" measures whether data is concentrated in the central "peak" or in the tails, and is calculated by

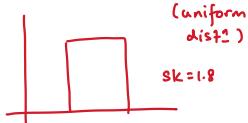
$$\text{sample kurtosis} = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^4}{\left[ \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^2}$$

$\exists_2$  Interpretation of sample kurtosis' value:

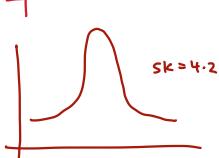
- ①  $sk = 3 \Rightarrow$  distribution looks "Gaussian" (bell-shaped);



- ②  $sk < 3 \Rightarrow$  distribution has shorter tails (more concentrated in the peak)



- ③  $sk > 3 \Rightarrow$  distribution has longer tails (less concentrated in the peak)



## ASSUMING A MODEL IS GAUSSIAN (CIS102)

$\exists_1$  To assume data can be reasonably modelled by a Gaussian distribution, we must have the following:

- ① The sample mean & median should be approximately equal;
- ② The sample skewness should be close to 0;
- ③ The sample kurtosis should be close to 3; and
- ④ ~95% of the observations should lie in the interval  $[\bar{y} - 2s, \bar{y} + 2s]$ .

## IN STATISTICS, WE DON'T PROVE THINGS! (CIS103)

$\exists_1$  In statistics, we don't prove assumptions are true, but instead find evidence against an assumption.

- ① If there is sufficient evidence against the assumption, then we say the data is "not consistent" with said assumption.
- ② Otherwise, we say the data is "consistent" with the assumption.

## FIVE NUMBER SUMMARY (CIS108)

$\exists_1$  The "five number summary" for a set of data is

- ① The minimum value  $y_{(1)}$ ;
- ②  $q_{(0.25)}$ ;
- ③  $q_{(0.5)}$ ;
- ④  $q_{(0.75)}$ ; &
- ⑤ The maximum value  $y_{(n)}$ .

$\exists_2$  In R, we can find the five number summary via the code

> summary(...)

# GRAPHICAL SUMMARIES (CIS112)

When displaying graphs, note that

- ① All graphs should be displayed at an appropriate size;
- ② Graphics should have clear titles which are fairly self-explanatory;
- ③ Axes should be labelled & units given where appropriate;
- ④ Choice of scales should be made with care; and
- ⑤ Graphics should not be used without thought, especially if there are better ways of displaying the information.

## HISTOGRAMS (CIS116)

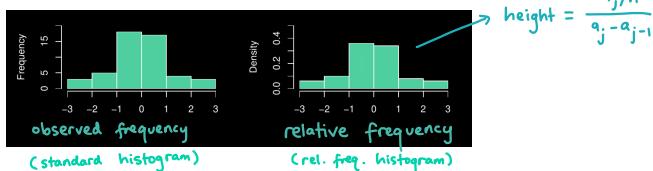
Essentially, histograms create a graphical summary of our data that we can use to compare with a pdf for crvs, or a pmf for a drv.

Let our data be  $y_1, \dots, y_n$ . Partition the range of the  $y$ 's into  $k$  non-overlapping intervals

$$I_j = [a_{j-1}, a_j], \quad j=1, 2, \dots, k.$$

Let  $f_j = \# \text{ of values from } \{y_1, \dots, y_n\} \text{ in } I_j$ . The  $f_j$ 's are called the "observed frequencies".

Then, draw a rectangle above each of the intervals with height proportional to the observed/relative frequency.



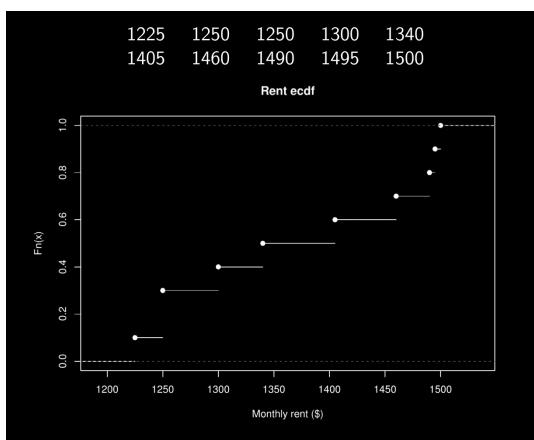
## EMPIRICAL CDF (CIS124)

An "empirical cdf" lets us compare the distribution of a dataset with a cdf of a random variable.

Mathematically, the empirical cdf is defined

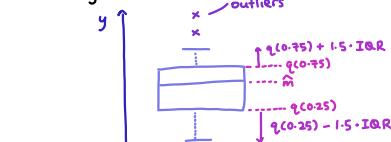
by

$$\hat{F}(y) = \frac{\#\text{ of values in } \{y_1, \dots, y_n\} \text{ which are } \leq y}{n} \quad \forall y \in \mathbb{R}.$$

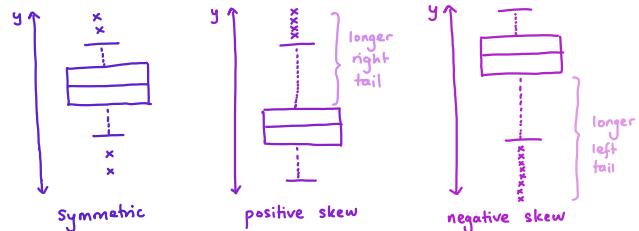


## BOX-PLOT (CIS139)

"Box-plots" give a graphical summary of the shape of a dataset's distribution in a similar way to the five number summary.

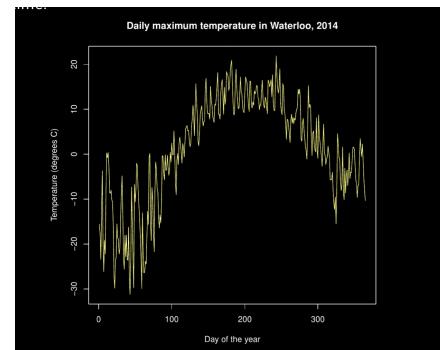


Box-plots can also show the skewness of a distribution:



## RUN CHART (CIS154)

A "run-chart" gives a graphical summary of data which are varying over time.



# SCATTERPLOTS (CIS157)

## BIVARIATE VS UNIVARIATE DATA (CIS157)

- Q1:** "Bivariate data" is of the form  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , where  $x_i, y_i \in \mathbb{R}$ . In contrast, "univariate data" is of the form  $\{y_1, \dots, y_n\}$  for  $y_i \in \mathbb{R}$ .

### SCATTER-PLOT (CIS158)

A "scatter-plot" for bivariate data is simply a plot of the  $(x_i, y_i)$ 's.



### SAMPLE CORRELATION: $r$ (CIS162)

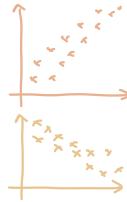
The "sample correlation", denoted " $r$ ", gives us a numerical summary of a bivariate dataset.

For data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ ,

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

In particular,  $r \in [-1, 1]$ , and measures the linear relationship between  $x$  &  $y$ .

- ① If  $r \approx -1$ , we say there is a strong negative linear relationship between the two variates.
- ② If  $r \approx +1$ , we say there is a strong positive linear relationship between the two variates.
- \*  $|r| \approx 1$  does not imply a causal relationship (correlation does not imply causation!)
- ③ If  $r \approx 0$ , we say there is no linear relationship between the two variates.
- \*  $r \approx 0$  does not imply  $x$  &  $y$  are unrelated — it just implies they are not linearly correlated.



Here  $r=0$  but obviously the data is related quadratically.

### RESPONSE & EXPLANATORY VARIATES (CIS171)

In an experiment, the "explanatory variate" is the variate that attempts to explain/determine the distribution of the "response variate".

\* explanatory variate = "independent" variable  
response variate = "dependent" variable.

# BIVARIATE CATEGORICAL DATA (CIS172)

We use the following survey as motivation:

- ① Hometown in Canada, like hockey
- ② Hometown not in Canada, like hockey
- ③ Hometown in Canada, dislike hockey
- ④ Hometown not in Canada, dislike hockey

Sample results from that survey:

	Canada ✓	Canada X	$\Sigma$
Hockey ✓	33	9	42
Hockey X	22	43	65
$\Sigma$	55	52	107

### RELATIVE RISK (CIS176)

Let  $A \subseteq \mathbb{X}$  &  $B \subseteq \mathbb{Y}$  be events in bivariate data " $X \times Y$ ".

Then the "relative risk" of "A with B" is equal to

$$\text{relative risk} = \frac{P(AB|B)}{P(A \cap B^c|B)}$$

e.g. in the survey above,

$$\begin{aligned} \text{relative risk of liking hockey} \\ \text{among those w/ a Canadian hometown} &= \frac{\text{prop. of Canada hometown who like hockey}}{\text{prop. of non-Canada hometown who dislike hockey}} \\ &= \frac{(33/55)}{(9/52)} \\ &= 3.467 \end{aligned}$$

# DATA ANALYSIS & STATISTICAL INFERENCE (CIS182)

## DESCRIPTIVE STATISTICS (CIS183)

"Descriptive statistics" is the portrayal of data (or parts of it) in numerical & graphical ways.  
\* all our previous work falls under this category!

## STATISTICAL INFERENCE (CIS184)

"Statistical inference" is the process of drawing general conclusions for a population/process based off of data obtained in a study about said population/process.

eg "based off my sample, I expect 90% of asmts this term to be submitted within the final 24 hrs of the deadline"

## INDUCTIVE VS DEDUCTIVE REASONING (CIS185)

- 1: "Inductive reasoning" occurs when we reason from the "specific" (observed data about a sample) to the "general" (the target population/process).
- 2: In contrast, "deductive reasoning" occurs when we use general results to prove theorems.  
\* proof by induction = deductive reasoning!

## ESTIMATION PROBLEMS (CIS187)

In "estimation problems," we are concerned about estimating one or more attributes of a population/process.

eg - estimate the prop. of STAT 231 students who like poutine  
- "fitting" a probability distribution for a process.

## HYPOTHESIS TESTING PROBLEMS (CIS188)

In a "hypothesis testing problem", we use the data to assess the truth of some question/hypothesis.

eg is it true a higher proportion of math majors than CS majors like poutine?

## PREDICTION PROBLEMS (CIS189)

In a "prediction problem", we use the data to predict a future value of a variate for a unit to be selected from the population/process.

eg given the past performance of a stock/other data, predict the value of the stock at some point in the future.

# Chapter 2: Statistical Models and Maximum Likelihood Estimation

## STATISTICAL MODELS (C2S191)

💡 A "statistical model" is a mathematical model that incorporates probability.

💡 These are useful since they can describe many different processes.

- eg - the daily closing value of CAD  
- when catastrophic events occur (eg pandemics)  
- the effect of drinking alcohol on your health

💡 We use random variables to represent a variate/characteristic of a randomly selected unit from the population/process.

eg let  $Y$  = how long I need to wait for the next game on an online video game.

💡 Statistical models can also be used to quantify any uncertainties obtained when drawing conclusions from data.

eg how the observed mean/variance of data differs from the actual mean/variance of data (eg goals scored in hockey)

💡 In particular, we can formulate questions of interest as parameters of the statistical model.

eg In the last example, say  $X = \#$  of hockey goals in a particular game

and suppose

$$X \sim Po(\theta).$$

We can then estimate  $\theta$  (ie the mean # of goals scored).

💡 We can then make decisions based on the results of our models, and use computers to simulate the processes.

# CHOOSING A PROBABILITY MODEL (C2S198)

When choosing a probability model, we use some or all of the following:

- ① Background knowledge / assumptions about the population/process that lead to certain distributions;
- ② Past experience with data sets from the population/process which show certain distributions are suitable;
- ③ Mathematical convenience (ie the tradeoff between complexity & accuracy), or
- ④ A current data set which the model can be assessed.

## "ALL MODELS ARE WRONG, BUT SOME ARE USEFUL" (C2S199)

Note that no statistical model is ever perfect, but that does not mean we cannot learn anything from imperfect ones.

(Quote from John Box)

## FAMILIES OF PROBABILITY DISTRIBUTIONS (C2S200)

Recall the following probability distributions:

- ① Poisson( $\theta$ )
- ② Exponential( $\theta$ )  
\* " $\theta$ " = mean of the distribution  
(not  $\frac{1}{\text{mean}}$ ).
- ③ Binomial( $n, \theta$ )
- ④ Gaussian( $\theta$ ) = Gaussian( $\mu, \sigma$ )
- ⑤ Multinomial( $n, \theta_1, \dots, \theta_n$ )
- ⑥ Geometric( $\theta$ )

## Y IS PARAMETERIZED BY $\theta \cdot f(y; \theta)$ (C2S205)

In particular, for each "family" of distributions, we get a different model for each value of  $\theta$ .

Thus, we say the random variable is "parameterized" by  $\theta$ .

If the r.v. is  $Y$ , we write the pf/pdf of  $Y$  as  $f(y; \theta)$  for  $y \in A = \text{range}(Y)$  to emphasize the dependence of the model on  $\theta$ .

## ESTIMATION OF UNKNOWN PARAMETERS (C2S206)

To determine how well the model fits the data, we need a value of  $\theta$  obtained from the data.

We usually denote this value  $\hat{\theta}$ .

- \* don't confuse  $\theta$  &  $\hat{\theta}$ !
  - $\theta$  = the underlying "true" value
  - $\hat{\theta}$  = our own estimate

This process is referred to as "estimating" the value of  $\theta$ .

## STEPS IN CHOOSING A MODEL (C2S208)

Suppose we have an experiment which involves collecting data to increase knowledge about a certain phenomena or to answer questions about a phenomena that has been carefully designed.

To choose a model for this experiment, we use the following steps:

- ① Collect/examine the data;  
\* more about this in Chap 3.
- ② Propose a model;  
eg  $G(\mu, \sigma)$
- ③ Fit the model;  
eg find  $\hat{\mu}, \hat{\sigma}$
- ④ Check the model;
- ⑤ If required, propose a revised model and return to ③,
- ⑥ Lastly, draw conclusions using the chosen model & the observed data.

# MAXIMUM LIKELIHOOD ESTIMATION (C2S210)

POINT ESTIMATE [OF A PARAMETER]:  $\hat{\theta}$  (C1S215)

A "point estimate" of a parameter, say  $\theta$ , is the value of a function of the observed data  $y$  and the other known quantities (eg the sample size  $n$ ).

We denote this estimate by " $\hat{\theta}$ ", where  $\hat{\theta} = \hat{\theta}(y)$ .

\* note  $\hat{\theta}$  is a function of  $y$ , and so  $\hat{\theta}$  depends on the value of  $y$  (the observed data).

For example:

①  $\mathcal{N}(\mu, \sigma)$ : estimate  $\mu$  by  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$  (the sample mean)

②  $\text{Bin}(n, \theta)$ : estimate  $\theta$  by  $\hat{\theta} = \frac{y}{n}$  (the sample proportion)

PARAMETER SPACE:  $\Omega$

The "parameter space"  $\Omega$  of a parameter  $\theta$  is the set of all possible values  $\theta$  can take.

LIKELIHOOD FUNCTION [FOR DRV]

$L(\theta)$  (C2S224)

Let  $y$  be potential data that will be used to estimate  $\theta$ , and let  $y$  be the actual observed data.

Suppose  $y$  is a drv.

Then, the "likelihood function for  $\theta$ " is defined to be

$$L(\theta) = L(\theta; y) = P(Y=y; \theta) \text{ for } \theta \in \Omega,$$

where  $\Omega$  is the parameter space of  $\theta$ .

\*  $L$  is technically a function of  $\theta$  &  $y$ , but for brevity we usually just write  $L(\theta)$ .

MAXIMUM LIKELIHOOD ESTIMATE / m.l. ESTIMATE:  $\hat{\theta}$  (C2S225)

The "maximum likelihood (ie m.l.) estimate" for given data  $y$  is the value of  $\theta$  which maximizes  $L(\theta)$ , and we denote it by  $\hat{\theta}$ .

In particular, generally  $\hat{\theta}$  satisfies

$$\frac{dL(\theta)}{d\theta} \Big|_{\theta=\hat{\theta}} = 0.$$

Why? - most distributions look like with a single "max" peak  
- so the only place the derivative will be 0 is at the peak, which we want.

RELATIVE LIKELIHOOD FUNCTION:  $R(\theta)$  (C2S234)

Let  $\hat{\theta}$  be the MLE of  $L(\theta)$ . Then, the "relative likelihood function" is

$$R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} \text{ for } \theta \in \Omega.$$

Note that

$$\textcircled{1} \quad 0 \leq R(\theta) \leq 1;$$

\textcircled{2}  $L(\hat{\theta})$  is a constant; and

\textcircled{3}  $R(\hat{\theta}) = 1$ , and so  $R$  is maximized at  $\theta = \hat{\theta}$ .

RELATIVE LIKELIHOOD FOR BINOMIAL DATA:

$$R(\theta) = \frac{\theta^y (1-\theta)^{n-y}}{\hat{\theta}^y (1-\hat{\theta})^{n-y}}, \quad \hat{\theta} = \frac{y}{n} \quad (\text{C2S235})$$

For binomial data, necessarily

$$R(\theta) = \frac{\theta^y (1-\theta)^{n-y}}{\hat{\theta}^y (1-\hat{\theta})^{n-y}}, \quad \hat{\theta} = \frac{y}{n}.$$

why?  $\rightarrow L(\theta) = (\hat{\theta})^\theta (1-\hat{\theta})^{n-y}$

$$= \theta^y (1-\theta)^{n-y}.$$

$$\text{Then } L(\hat{\theta}) = \hat{\theta}^y (1-\hat{\theta})^{n-y}.$$

$$(\hat{\theta} = \frac{y}{n} \text{ from earlier})$$

$$\Rightarrow R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} = \frac{\theta^y (1-\theta)^{n-y}}{\hat{\theta}^y (1-\hat{\theta})^{n-y}}.$$

\* when computing relative likelihoods, we can ignore any constants wrt  $\theta$  as they will cancel out in the computation of  $R(\theta)$ .

LOG LIKELIHOOD FUNCTION:  $\ell(\theta)$

(C2S237)

The "log likelihood function" is defined to be

$$\ell(\theta) = \log L(\theta) \quad \forall \theta \in \Omega.$$

\* log = ln for this course!

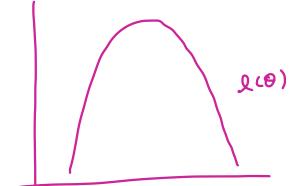
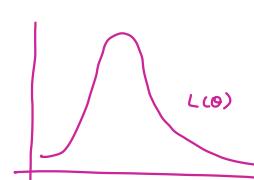
Note that  $\ell(\theta)$  is maximized for the same value of  $\theta$  as the regular likelihood function.

\* ie  $\ell'(\hat{\theta}) = 0 \Leftrightarrow \ell'(\hat{\theta}) = 0$ .

$\ell(\theta)$  is also preferred over  $L(\theta)$  because it is usually easier to take derivatives of  $\ell$  (which typically involves sums) over  $L$  (which typically involves products).

However, note  $\ell(\theta)$  has a different "shape" than  $L(\theta)$  (it looks more "quadratic").

eg  $L(\theta) = \theta^y (1-\theta)^{n-y}$



# LIKELIHOOD FUNCTION FOR INDEPENDENT EXPERIMENTS (C2S244)

Suppose we observe data  $Y = (Y_1, \dots, Y_n)$  that are iid each with p.f.  $P(Y_i = y_i; \theta)$ . Then the (combined) likelihood function for  $\theta$  based on the data  $(y_1, \dots, y_n)$  is

$$L(\theta) = \prod_{i=1}^n L_i(\theta) = \prod_{i=1}^n P(Y_i = y_i; \theta) \quad \forall \theta \in \Omega.$$

## RELATIVE LIKELIHOOD FOR POISSON DATA:

$$R(\theta) = \frac{\theta^n e^{-n\theta}}{\hat{\theta}^n e^{-n\hat{\theta}}}, \quad \hat{\theta} = \bar{y} \quad (\text{C2S254})$$

For Poisson data, necessarily

$$R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} = \frac{\theta^n e^{-n\theta}}{\hat{\theta}^n e^{-n\hat{\theta}}}, \quad \hat{\theta} = \bar{y}$$

Proof First, see that

$$P(Y_i = y_i; \theta) = \frac{\theta^{y_i} e^{-\theta}}{y_i!}.$$

Therefore

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n P(Y_i = y_i; \theta) = \prod_{i=1}^n \frac{\theta^{y_i} e^{-\theta}}{y_i!} \\ &= \prod_{i=1}^n \frac{1}{y_i!} \prod_{i=1}^n \theta^{y_i} \prod_{i=1}^n e^{-\theta} \\ &= \frac{\prod_{i=1}^n y_i!}{\theta^n} e^{-n\theta} \quad (\text{we ditch the constant}) \\ &= \theta^n e^{-n\theta}, \quad (\because \bar{y} = \frac{1}{n} \sum y_i) \end{aligned}$$

and so

$$L(\theta) = \log L(\theta) = n\bar{y} \log(\theta) - n\theta.$$

Thus

$$L'(\theta) = \frac{n\bar{y}}{\theta} - n \quad (= 0)$$

and so  $L$  (and thus  $L$ ) is maximized when  $\theta = \bar{y}$  ( $= \hat{\theta}$ ).

Therefore

$$R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} = \frac{\theta^n e^{-n\theta}}{\hat{\theta}^n e^{-n\hat{\theta}}}, \quad \hat{\theta} = \bar{y}. \quad \blacksquare$$

## RANDOM SAMPLE: $Y_1, \dots, Y_n$ (C2S256)

Suppose  $Y_1, \dots, Y_n$  are iid with p.f.  $P(Y_i = y_i; \theta) = f(y_i; \theta)$ . We call  $Y_1, \dots, Y_n$  a "random sample".

## LIKELIHOOD FUNCTION FOR A RANDOM SAMPLE (C2S257)

Let  $Y_1, \dots, Y_n$  be a random sample, with p.f.  $P(Y_i = y_i; \theta) = f(y_i; \theta)$ .

Let  $y_1, \dots, y_n$  be a realization of (ie the observed data from) the random sample.

Then the likelihood function for  $\theta$  based on the observed sample is

$$L(\theta) = \prod_{i=1}^n P(Y_i = y_i; \theta) \quad \forall \theta \in \Omega.$$

Proof:  $L(\theta) = P(\text{observing the data } y_1, \dots, y_n \text{ given } \theta)$   
 $= P(Y_1 = y_1, \dots, Y_n = y_n; \theta)$   
 $= P(Y_1 = y_1; \theta) \dots P(Y_n = y_n; \theta) \quad (\text{by independence})$   
 $= \prod_{i=1}^n P(Y_i = y_i; \theta). \quad \blacksquare$

# LIKELIHOOD FOR CONTINUOUS RANDOM VARIABLES

## (C2S258)

### LIKELIHOOD FUNCTION FOR CRV

#### (C2S262)

Let  $\mathbf{y} = (y_1, \dots, y_n)$  be a random sample from a continuous distribution with pdf  $f(y; \theta)$  for  $\theta \in \Omega$ .

Let  $y = (y_1, \dots, y_n)$  be a realization of  $\mathbf{Y}$ .

Then, the likelihood function for  $\theta$  based on the observed data  $y = (y_1, \dots, y_n)$  is defined to be

$$L(\theta) = L(\theta; y) = \prod_{i=1}^n f(y_i; \theta) \quad \forall \theta \in \Omega.$$

### MLE FOR $\text{Exp}(\theta)$ : $\hat{\theta} = \bar{y}$ (C2S266)

Let  $\mathbf{Y} \sim \text{Exp}(\theta)$ , and let  $(y_1, \dots, y_n)$  be the observed data from a sample of size  $n$ . Then the maximum likelihood estimate is necessarily  $\hat{\theta} = \bar{y}$ .

Proof. See that

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n \frac{1}{\theta} e^{-\frac{y_i}{\theta}} \\ &= \theta^{-n} e^{-\bar{y}/\theta}, \end{aligned}$$

and so

$$l(\theta) = \log L(\theta) = -n \log \theta - \frac{n\bar{y}}{\theta} \quad (=0).$$

Hence

$$l'(\theta) = -\frac{n}{\theta} + \frac{n\bar{y}}{\theta^2} \quad (=0)$$

and it follows  $l$  (and so  $L$ ) is maximized when  $\theta = \bar{y}$  ( $= \hat{\theta}$ ).  $\blacksquare$

### LIKELIHOOD FUNCTION FOR $\mathcal{N}(\mu, \sigma^2)$ :

$$L(\theta) = (2\pi)^{-n/2} \sigma^{-n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right] \quad (\text{C2S267})$$

Let  $y_1, \dots, y_n$  be observations from  $\mathbf{Y} \sim \mathcal{N}(\mu, \sigma^2)$ .

Then necessarily

$$L(\theta) = (2\pi)^{-n/2} \sigma^{-n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right].$$

$$\begin{aligned} \text{Proof. } L(\theta) &= \prod_{i=1}^n f(y_i; \mu, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2} (y_i - \mu)^2\right] \\ &= (2\pi)^{-\frac{n}{2}} \sigma^{-n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right]. \quad \blacksquare \end{aligned}$$

In particular, the MLE is

$$\hat{\mu} = \bar{y}, \quad \hat{\sigma} = \left[ \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{\frac{1}{2}}.$$

$$\text{Proof. First, see that } l(\theta) = \log(L(\theta)) = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2.$$

$$\text{Then } \frac{\partial l}{\partial \mu} = \frac{n}{\sigma^2} (\bar{y} - \mu) \quad \& \quad \frac{\partial l}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (y_i - \mu)^2.$$

Thus

$$\frac{\partial l}{\partial \mu} = 0 \Rightarrow \hat{\mu} = \bar{y} \quad \&$$

$$\frac{\partial l}{\partial \sigma} = 0, \quad \hat{\mu} = \bar{y} \Rightarrow \hat{\sigma} = \left[ \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{\frac{1}{2}}. \quad \blacksquare$$

### INVARIANCE PROPERTY OF MLEs

#### (C2S272)

Let  $\hat{\theta}$  be the MLE of a parameter  $\theta$ .

Then  $g(\hat{\theta})$  is necessarily the MLE of  $g(\theta)$ .

e.g. suppose  $\mathbf{Y} \sim \text{Poi}(\theta)$ ,  $\hat{\theta} = \bar{y}$ .

$$\text{Then } P(Y \geq 3) = 1 - P(Y \leq 2) = 1 - \sum_{y=0}^2 \frac{\theta^y e^{-\theta}}{y!}.$$

But this is a function of  $\theta$ , so the MLE of  $P(Y \geq 3)$  is

$$1 - \sum_{y=0}^2 \frac{\hat{\theta}^y e^{-\hat{\theta}}}{y!}.$$

\* in R, we calculate this via "1 - ppois(2, 3)"

We should always clarify when/where we use the invariance property.

# CHECKING MODEL FIT (C2S276)

## COMPARING OBSERVED VS EXPECTED FREQUENCIES [FOR DRV] (C2S277)

To check whether a model fits a given set of data, we can compare the observed frequencies & the expected frequencies using a table.

eg Suppose a hockey team scored the following # of goals in these # of games:

Goals	0	1	2	3	4	5	6	7
Games	2	17	21	18	15	7	1	1

Let's say we assume the data can be modelled by a Poisson distn, say  $Y \sim \text{Poi}(\theta)$ . Then, we estimate  $\theta$  using the MLE of  $\theta$ , aka  $\hat{\theta}$ :

$$\hat{\theta} = \bar{y} = \frac{1}{82} (2(0) + 17(1) + \dots + 1(7)) = 2.695.$$

Next, we calculate the expected frequencies.

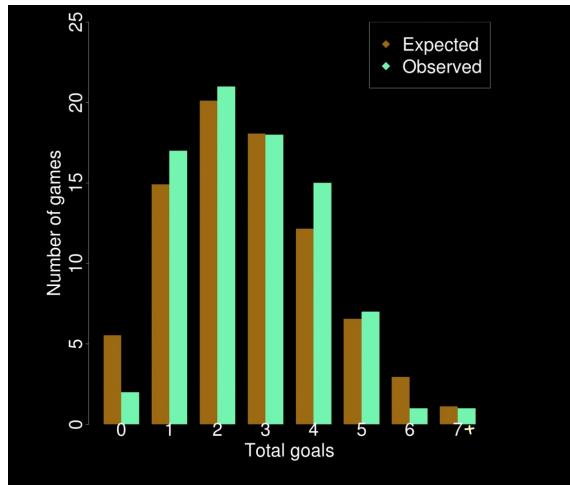
Since the range of Poi is technically  $0, 1, 2, \dots$ , we need to account for the "right tail" by grouping all the values  $\geq 7$  into "one" value:

Goals	0	1	2	3	4	5	6	7
obs.	2	17	21	18	15	7	1	1
Exp.	5.54	14.93	20.11	18.07	12.17	6.56	2.95	1.67

where

$$\text{exp value for } i = nP(Y=i) = 82 \frac{e^{-2.695}}{i!} (2.695)^i.$$

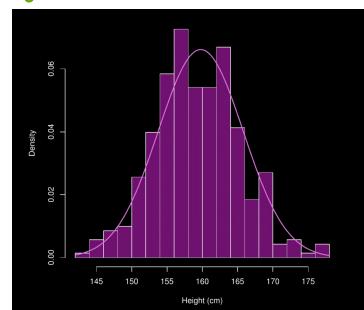
We may also plot the expected/observed frequencies via a bar plot.



## COMPARING OBSERVED VS EXPECTED FREQUENCIES [FOR CRV] (C2S300)

We can do something similar for continuous random variables.

eg Consider the dataset



We have

$$\bar{y} = 159.77,$$

$$s^2 = 36.36,$$

$$s = 6.03.$$

Let  $Y$  be the data.

How reasonable is it to model the data via a Gaussian distribution?

Suppose it is; ie  $Y \sim \mathcal{N}(\mu, \sigma^2)$ .

We estimate

$\mu \approx$  the sample mean (MLE); &

$\sigma \approx$  the sample sd (not the MLE),

so  $Y \sim \mathcal{N}(159.77, 6.03)$ .

We then can estimate the exp. probabilities  $Y$  falls into one of the intervals of the histogram outlined above; eg

$$P(160 \leq Y \leq 162) = P\left(\frac{160 - 159.77}{6.03} < Z < \frac{162 - 159.77}{6.03}\right) \\ = 0.129$$

\*in R, we can calculate this via

> pnorm(0.370) - pnorm(0.038)

or

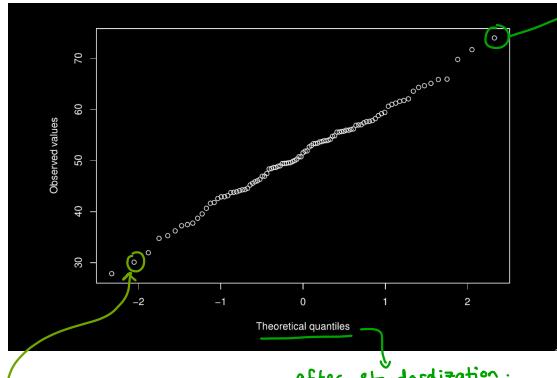
> pnorm(162, 159.77, 6.03) - pnorm(160, 159.77, 6.03).

and thus calculate the exp. # of values to fall within the given interval; eg

$e_j = 351p_j$ , where  $I_j$  is the  $j^{\text{th}}$  interval, and compare this with the observed values.

## QQ / QUANTITY - QUANTITY PLOTS (C2S311)

- B1** A "QQ plot" plots the observed values / quantiles from the sample data on the y-axis over the theoretical values obtained by fitting a model to said data on the x-axis.
- B2** In particular, we may standardize the theoretical values and plot that on the x-axis instead.



each dot corresponds to a quantile; ie the  $q^{\text{th}}$  quartile.

- ① The y-value corresponds to the value  $y$  such that " $q$ " of the data is  $\leq y$ .
- ② The x-value corresponds to the value  $x$  such that if we fit the rv  $Y$  to a model, and standardize said model to be  $Z$ , then  $P(Z \leq x) = q$ .

**B3** If we model  $Y \sim N(\mu, \sigma)$ , then the QQ-plot of the points

$$(\Phi^{-1}\left(\frac{i}{n+1}\right), y_{(i)}) \quad \text{for } 1 \leq i \leq n,$$

where  $y_{(1)}, \dots, y_{(n)}$  is the observed data, should be approximately a straight line if the normal distribution is a good fit for said data.

## USING QQ-PLOTS TO INFERENCE SHAPE OF DISTRIBUTION (C2S341)

We can use QQ-plots to infer the underlying shape of a distribution:

- ① If the points are along a straight line, then this indicates the data is normal.



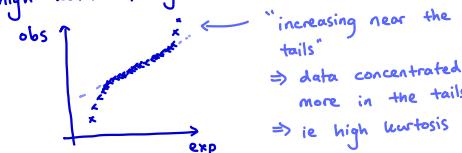
- ② If the data is S-shaped, this indicates symmetry (ie low skewness).

→ then, the relative abundance of points in the "center" vs. tails implies the magnitude of the kurtosis.

low kurtosis, symmetric



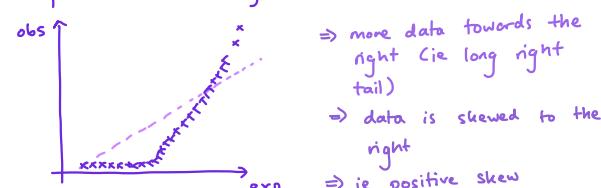
high kurtosis, symmetric



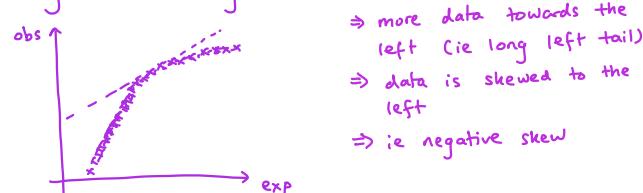
- ③ If the data is U-shaped, this indicates asymmetry.

→ then, the relative abundance of points in the left vs. right tails implies the magnitude and sign of the skewness.

positive skew, asymmetric



negative skew, asymmetric



## NORMALITY CHECKING SUMMARY (C2S344)

To assume data is a good fit for a Gaussian model, we need to check:

- ① The sample mean & median are approximately equal;
- ② The sample skewness is close to 0;
- ③ The sample kurtosis is close to 3;
- ④ Approximately 95% of the observations lie in  $[\bar{y} - 2s, \bar{y} + 2s]$ ;
- ⑤ Histograms & ecdfs should show agreement between the data & theoretical distribution;
- ⑥ The QQ-plot should roughly be a straight line.

## UNBIASED ESTIMATOR: $S^2$ (C2S352)

The "unbiased estimator" for data is defined to be

$$S^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

## $k^{\text{th}}$ POPULATION MOMENT: $\mu_k$ (C2S353)

The " $k^{\text{th}}$  population moment" of  $Y$  is defined to be

$$\mu_k = E[Y^k].$$

In particular,

- ①  $\mu_1 = E[Y]$ ;
- ②  $\mu_2 = \text{Var}(Y) + E[Y]^2$ .

## $k^{\text{th}}$ SAMPLE MOMENT: $m_k$ (C2S355)

Let  $y_1, \dots, y_n$  be a sample.

Then, the " $k^{\text{th}}$  sample moment" is defined to be

$$m_k = \frac{1}{n} \sum_{i=1}^n y_i^k.$$

## METHOD OF MOMENTS FOR ESTIMATION (C2S358)

The "method of moments" allows us to estimate parameters for a model, based off the data we are using the model for.

Steps:

- ① Compute the first  $p$  sample moments, where  $p = \# \text{ of parameters}$ .
- ② Relate the population moments to the true parameter values.
- ③ Use the sample moments to solve the resulting system of equations to estimate the parameters.

## EXAMPLE 1: $G(\mu, \sigma)$ (C2S356)

Problem:

"Suppose  $Y \sim G(\mu, \sigma)$ . Use the sample  $y_1, \dots, y_n$  to estimate  $\mu$  and  $\sigma$ ".

Sol<sup>n</sup>. Since  $Y \sim G(\mu, \sigma)$ , and

$$\mu = E(Y) = \mu_1.$$

$$\sigma^2 = E(Y^2) - E(Y)^2 = \mu_2 - \mu_1^2.$$

we can estimate the values of  $\mu$  &  $\sigma$  by

$$\hat{\mu} = m_1, \quad \hat{\sigma}^2 = m_2 - m_1^2.$$

Hence

$$\hat{\mu} = m_1 = \frac{1}{n} \sum_{i=1}^n y_i$$

&

$$\begin{aligned} \hat{\sigma}^2 &= m_2 - m_1^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n y_i^2 \right) - \bar{y}^2 \\ &= \frac{1}{n} \left( \sum_{i=1}^n (y_i - \bar{y})^2 \right). \end{aligned}$$

In Note: we use the " $\hat{\cdot}$ " notation for both MLE and method of moments!

## EXAMPLE 2: $\text{Unif}(a, b)$ (C2S361)

Problem:

"Suppose  $y_1, \dots, y_n$  are independently sampled from a continuous uniform distribution on  $(a, b)$ . What are the method of moments estimates on  $(a, b)$ ?"

Sol<sup>n</sup>. We need to estimate  $2$  parameters, so we require

$$\mu_1 = E(Y), \quad \mu_2 = E(Y^2)$$

and hence we need to use

$$m_1 = \frac{1}{n} \sum_{i=1}^n y_i, \quad m_2 = \frac{1}{n} \sum_{i=1}^n y_i^2.$$

\*remember  $m_1$  &  $m_2$  are both numbers (since they are based off the sample!)

Then, using LOTUS,

$$\begin{aligned} \mu_1 &= \int_a^b \frac{y}{b-a} dy = \frac{1}{b-a} \left[ \frac{y^2}{2} \right]_a^b \\ &= \frac{1}{2} \left( \frac{1}{b-a} \right) (b^2 - a^2) \\ &\quad & \end{aligned}$$

$$\mu_2 = \int_a^b \frac{y^2}{b-a} dy = \dots = \frac{1}{3} (a^2 + b^2 + ab).$$

We then estimate  $\mu_1 \approx m_1$  &  $\mu_2 \approx m_2$ , so that

$$m_1 = \frac{1}{2} (\hat{a} + \hat{b}) \Rightarrow \hat{a} = 2m_1 - \hat{b},$$

$$\& \quad m_2 = \frac{1}{3} (a^2 + b^2 + ab) \Rightarrow (\hat{b} - m_1)^2 = 3(m_2 - m_1^2)$$

using the appropriate subs<sup>ts</sup>s.

Solving for  $\hat{a}$  &  $\hat{b}$  yields the estimates

$$\hat{b} = m_1 + \sqrt{3(m_2 - m_1^2)}$$

$$\& \quad \hat{a} = m_1 - \sqrt{3(m_2 - m_1^2)},$$

and by evaluating  $m_1$  &  $m_2$  we could then compute  $\hat{a}$  &  $\hat{b}$ .

Moreover, note that

$$m_2 - m_1^2 = \frac{n-1}{n} s^2,$$

and so we could also write the above as

$$\hat{a} = m_1 - \sqrt{\frac{(n-1)}{n} s^2}, \quad \hat{b} = m_1 + \sqrt{\frac{(n-1)}{n} s^2}.$$

## EXAMPLE 3: CONTEST & PRIZES (C2S372)

Problem:

"A contest awards prizes as follows:

- $P(\text{win \$1}) = a$ ;
- $P(\text{win \$10}) = b$ ;
- $P(\text{lose}) = 1 - a - b$ .

You buy five tickets and win three times, including one \\$10 win.

Use MM to estimate  $a$  &  $b$ ."

Sol<sup>n</sup>. Again, we have two parameters, so we need

$$\mu_1 = E(Y), \quad \mu_2 = E(Y^2)$$

and use the sample moments

$$M_1 = \frac{1}{n} \sum y_i, \quad M_2 = \frac{1}{n} \sum y_i^2.$$

Since  $Y$  is a drv, we use

$$E(Y^k) = \sum_{y \in A} y^k f(y),$$

and for this example

$$A = \{0, 1, 10\}, \quad f(0) = 1 - a - b, \quad f(1) = a, \quad f(10) = b.$$

Hence

$$\mu_1 = E(Y) = 0(1-a-b) + 1(a) + 10(b) = a + 10b;$$

$$\& \quad \mu_2 = E(Y^2) = 0^2(1-a-b) + 1^2(a) + 10^2(b) = a + 100b.$$

Then, we estimate  $\mu_i$  with  $m_i$  to get

$$M_1 = \hat{a} + 10\hat{b}, \quad M_2 = \hat{a} + 100\hat{b}.$$

Solving for  $\hat{a}$  &  $\hat{b}$  yields that

$$\hat{b} = \frac{M_2 - M_1}{90}, \quad \hat{a} = M_1 - \frac{M_2 - M_1}{90}.$$

Finally, for our sample we observed  $\{0, 0, 1, 1, 10\}$  and so

$$M_1 = 2.4, \quad M_2 = 20.4$$

and so

$$\hat{a} = 0.4, \quad \hat{b} = 0.2. \quad \square$$

# Chapter 3: Planning and Conducting Empirical Studies

B1 Recall "empirical studies" are those where data collected can be used to learn about a population/process.

\* we use this "Pfizer vs. Moderna" study for examples:

[www.nejm.org/doi/full/10.1056/NEJMoa2115463](http://www.nejm.org/doi/full/10.1056/NEJMoa2115463)

so it might be helpful to have the study open whilst reading this chapter.

## PPDAC (C3S384)

B1 We can design an empirical study using "PPDAC".

B2 In particular, this stands for

- ① Problem — a clear statement of the study's objectives;
- ② Plan — the procedures in the study, how the data is collected
- ③ Data — the physical collection of the data
- ④ Analysis — analysis of said data
- ⑤ Conclusion — conclusions drawn from said analysis, and their limitations

## PROBLEM (C3S393)

- The "problem" addresses
- ① what group of things/people do we want the conclusions to apply?
  - ② what variates can we define?
  - ③ what are the questions we are trying to answer?
  - ④ what conclusions are we trying to draw?

## TARGET POPULATION/PROCESS (C3S394)

- The "target population/process" is the collection of units to which the experimenters conducting the empirical study wish the conclusions to apply.
- In the problem, the units & target population/process must be defined.

eg in the vaccine study, possible target pop<sup>n</sup>s/processes:

- ① people in VA health-care system now and in future
- ② unvaccinated people in the VA health care system now in the future
- ③ ① & ② but limiting the time period to the duration of the COVID-19 pandemic.

## VARIATES [IN EMPIRICAL STUDIES] (C3S398)

- A "variate" is a characteristic of a unit.

- To determine the variates, look at what is measured or recorded on each unit.

eg for the vaccine study, the variates include

- which vaccine each participant took;  
(ie Pfizer / Moderna)
- outcome indicators such as COVID-19 infection, symptoms, hospitalization, and death;
- age, sex, race, residence, geographic location;
- etc

## ATTRIBUTES [IN EMPIRICAL STUDIES] (C3S402)

- "Attributes" are functions of variates over a population.

- In the problem step, the questions of interest are specified in terms of the attributes of the target population.

eg in the vaccine study, possible attributes include

- ① the proportion of people in the target pop<sup>n</sup> who would contract COVID-19 after receiving the Pfizer vaccine within 24 weeks;
- ② the proportion of people in the target pop<sup>n</sup> who would contract COVID-19 after receiving the Moderna vaccine within 24 weeks;
- ③ the difference in the preceding two numbers.

## TYPES OF PROBLEMS (C3S405)

- Types of problems an empirical study can solve:

- ① "Descriptive" — determine a particular attribute of the population.
    - eg - the national unemployment rate
    - estimating the relative efficacy of the two vaccines among all those who received it at the time of the study
  - ② "Causative" — determine the existence (or lack of) of a causal relationship between two variates.
    - eg - does a new hockey helmet reduce the risk of concussion
    - whether giving someone the Moderna vaccine instead of the Pfizer vaccine reduces their risk of COVID-19
  - ③ "Predictive" — predict the response for a given unit.
    - eg - predict e-cig weekly sales if sales tax on them is doubled
    - estimating relative efficacy of Pfizer & Moderna
- Note that we usually cannot answer causative problems from observational studies.
- Causative & descriptive problems are also hard to distinguish.

# PLAN (C3S411)

- In the plan, we
  - decide what units are available for study;
  - what units will be examined; &
  - what variates will be collected and how.

## STUDY POPULATION / PROCESS (C3S411)

The "study population/process" is the collection of units available to be included in the study.

Note the study population is a strict subset of the target population.

eg - veterans of age  $\geq 18$  years, no previous COVID infection, etc (in vaccine study)

## STUDY ERROR (C3S423)

"Study error" occurs when the attributes in the study population differ from said attributes in the target population.

eg say  $Y_T \sim \text{Bin}(n, \theta_T)$  represents the # of people in a random sample of size  $n$  from the target pop who support Brexit.

Say  $Y_S \sim \text{Bin}(n, \theta_S)$  represents the # of people in a random sample of size  $n$  from the study pop who support Brexit.

We might be concerned  $\theta_T \neq \theta_S$ . (C3S422)

Note that just the study/target populations being different is not study error — the difference must be in their attributes.

Moreover, note study error concerns populations; we do not care about the study/target samples.

Hence, we must be careful when thinking about the attributes of interest in a study.

In particular, as the values of the target or study populations' attributes are unknown, the study error cannot be quantified.

Instead, we generally rely on expertise from other sources to determine whether conclusions derived from the study population may apply to the target population.

eg whether studies on mice apply to humans.  
(study) (target)

## SAMPLING PROTOCOL (C3S430)

The "sampling protocol" is the procedure used to select a sample of units from the study population.

In practice, obtaining a (truly) random sample is difficult/impossible/expensive, so less rigorous sampling methods are usually used.

eg "matching" in the vaccine study

## SAMPLE SIZE (C3S430)

The "sample size" is the number of units sampled from the sampling protocol.

## SAMPLE ERROR (C3S435)

"Sample error" occurs when the attributes in the sample differ from the attributes in the study population.

\* again, it must be a difference in the attributes, not just because the two groups differ!

Note sample error does not care about the target population & sample!

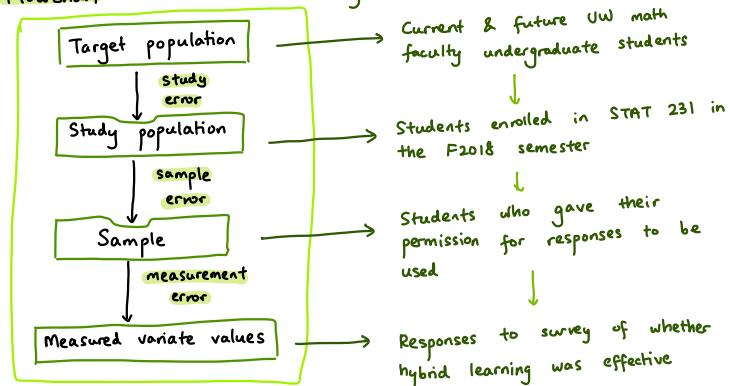
## MEASUREMENT ERROR (C3S442)

"Measurement error" occurs if the measured and true values of a variate are not identical.

eg - measuring blood pressure  
- patients more stressed in doctor's office  
- so reading is higher  
- "white coat hypertension"

## STEPS IN THE PLAN (C3S445)

Flowchart:



## DATA (C3S454)

💡 1 The "data" step concerns collecting data according to the plan.

💡 2 To do this, the  
① variates must be clearly defined; &  
② satisfactory methods of measuring them must be used.

## RECORDING DATA (C3S455)

💡 1 Note mistakes can occur in recording data into a DB, and so for more complex investigations it is useful to put checks in place to avoid these mistakes & detect those that are made.

💡 2 Moreover, when lots of data is used, database design and management is important.

💡 3 Also, if data is recorded longitudinally (ie over a period of time), departures from the plan might occur; these must be recorded.

eg persons might drop out of a long-term medical study because of adverse reactions to a treatment.

💡 4 Such departures will affect the Analysis & Conclusion steps.

## ANALYSIS (C3S456)

💡 1 In the "analysis" step, we analyze the data collected.

💡 2 This includes  
① numerical & graphical summaries of the data;  
② selecting an appropriate model; &  
③ checking if said model is a good fit.

💡 3 We usually formulate these questions in terms of the model parameters.

eg "if  $Y \sim \text{Bin}(n, \theta)$  &  $\theta = P(\text{new drug cures a disease})$ , what is  $\theta$ ?"

💡 4 Departures from the plan that affect the analysis must also be noted.

## CONCLUSION (C3S458)

💡 1 In the "conclusion" step, the questions posed in the problem are answered to the extent permitted by the data.

💡 2 In other words, the conclusion is directed by the problem.

💡 3 The conclusion must also feature  
① a discussion and/or quantification of potential study, sample & measurement errors;  
② departures from the plan that affect the analysis;  
③ the limitations of the study.

# Chapter 4: Estimation

## STATISTICAL MODELS & ESTIMATION (C4S463)

💡 In choosing a model for the analysis

step of PPDAC we need to consider:

① Model A: a model for variation in the population/process being studied which includes the attributes which are to be estimated; and

② Model B: a model which takes into account how the data were collected & which is constructed in conjunction with model A.

e.g. (See C4S466 for more details)

$y = \# \text{ of } X \text{ in a randomly chosen sample from the target pop}^n \text{ who have had COVID-19.}$

For model A, we may assume

$$y \sim \text{Bin}(n, \theta_T),$$

where  $\theta_T = \text{proportion of target pop}^n \text{ who have had COVID-19}$

(not "probability person has COVID")

For model B, we take into account target & study populations are not the same.

We assume

$$Y \sim \text{Bin}(n, \theta)$$

where  $\theta = \text{proportion of study pop}^n \text{ who have had COVID-19.}$

\* If  $\theta_T \neq \theta$ , this represents study error.

💡 For this course, we assume

① data arises from a random sample from the study population; &

② variates are measured without error.

💡 This only means we are able to estimate attributes of interest about the study population, not the target population.

\* if we make any inferences about the target pop<sup>n</sup>, we have to state our assumptions.

# ESTIMATORS & SAMPLING DISTRIBUTIONS (C4S474)

Note that sampling is an inherently random process.

## RANDOM VARIABLE ASSOCIATED WITH $\bar{Y}$ : $\bar{Y}$ (C4S491)

Let  $Y_1, \dots, Y_n$  be iid. Then we define

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

In particular, if  $Y_i \sim \mathcal{N}(\mu, \sigma^2)$ , then

$$\bar{Y} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right).$$

## ESTIMATOR (C4S498)

An "estimator" is a rule that tells us how to process the data to obtain an estimate of an unknown parameter  $\theta$ .

## POINT ESTIMATORS: $\hat{\theta}$ (C4S496)

Let  $Y_1, \dots, Y_n$  be potential observations in a random sample.

Consider the point estimate

$$\hat{\theta} = g(Y_1, \dots, Y_n).$$

Then, we can associate  $\hat{\theta}$  with a random variable

$$\star \quad \hat{\theta} = g(Y_1, \dots, Y_n).$$

e.g. the random variable associated w/  $\hat{\theta} = \bar{y} = \frac{1}{n} \sum_i y_i$   
is  $\hat{\theta} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ .

\*  $\hat{\theta}$  = estimate (single value); &  
 $\hat{\theta}$  = random variable!

## SAMPLING DISTRIBUTION [OF AN ESTIMATOR] (C4S500)

The "sampling distribution" of an estimator  $\hat{\theta}$  is the distribution of  $\hat{\theta}$ .

## GAUSSIAN SAMPLING DISTRIBUTION (C4S511)

Let  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ , so  $\bar{Y} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

(this is the sampling distribution of the sample mean.)

Vary parameters, and how does it affect the sampling distribution:

	↑ sample size, $n$	↑ std dev, $\sigma$	↑ mean, $\mu$
location	does not change	does not change	moves to the right
spread	decreases	increases	does not change
Shape	does not change	does not change	does not change

Thus, the probability we draw a sample that yields an estimate  $\hat{\mu}$  close to  $\mu$

- ① increases as  $n$  increases;
- ② decreases as  $\sigma$  increases; &
- ③ does not change with  $\mu$ .

\* in particular, because

$$P(|\hat{\mu} - \mu| \leq \epsilon) = P(\mu - \epsilon \leq \bar{Y} \leq \mu + \epsilon) \\ = P\left(\frac{-\epsilon\sqrt{n}}{\sigma} \leq Z \leq \frac{\epsilon\sqrt{n}}{\sigma}\right) \text{ (as } \bar{Y} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)).$$

Moreover, see that  $sd(\bar{Y}) \approx \frac{\sigma}{\sqrt{n}}$ , and so

- ①  $sd(\bar{Y})$  decreases as  $n$  increases, and so more of our sample estimates will be closer to  $\mu$ ;
- ②  $sd(\bar{Y})$  increases as  $\sigma$  increases, and so less of our sample estimates will be closer to  $\mu$ ;
- ③  $sd(\bar{Y})$  does not change with  $\mu$ . (C4S540)

## NORMAL APPROXIMATIONS (C4S525)

If  $Y_1, \dots, Y_n$  are iid with mean  $\mu$  & variance  $\sigma^2$ , then by the CLT for large enough samples we have

$$\frac{\bar{Y} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} = Z_n \rightarrow \mathcal{N}(0, 1).$$

Particular examples:

- ① Binomial — If  $Y \sim \text{Bin}(n, \theta)$ , then

$$\frac{\frac{Y}{n} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \sim \mathcal{N}(0, 1).$$

- ② Exponential — If  $Y_i \sim \text{Exp}(\theta)$ , then for large  $n$

$$\frac{\bar{Y} - \theta}{\frac{\theta}{\sqrt{n}}} \sim \mathcal{N}(0, 1).$$

- ③ Poisson — If  $\theta \geq 5$  then if  $Y \sim \text{Poi}(\theta)$  then

$$Y \approx \mathcal{N}(\theta, \sqrt{\theta}).$$

If  $Y_i \sim \text{Poi}(\theta)$ , then for large  $n$

$$\frac{\bar{Y} - \theta}{\sqrt{\theta}} \sim \mathcal{N}(0, 1).$$

# COMPARING ESTIMATORS (C4S551)

BIAS [OF AN ESTIMATOR]:

Bias ( $\hat{\theta}$ ) (C4S553)

The "bias" of an estimator  $\hat{\theta}$  is given by

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta.$$

If the bias is zero, we say the estimator is "unbiased".

EXAMPLE:  $Y \sim \text{Bin}(n, \theta)$ ,  $\hat{\theta} = \frac{Y}{n}$  (C4S554)

Problem:

"Suppose  $Y \sim \text{Bin}(n, \theta)$ . Show  $\hat{\theta} = \frac{Y}{n}$  is unbiased."

$$\begin{aligned} \text{Soln. } \text{Bias}(\hat{\theta}) &= E(\hat{\theta}) - \theta \\ &= E\left(\frac{Y}{n}\right) - \theta \\ &= \frac{1}{n}E(Y) - \theta \\ &= \frac{1}{n}(n\theta) - \theta = 0. \end{aligned}$$

EXAMPLE:  $\hat{\sigma}^2$  IN  $\mathcal{G}(\mu, \sigma)$  (C4S555)

Problem:

"Consider  $Y_1, \dots, Y_n$  iid  $\mathcal{G}(\mu, \sigma)$ . What is the bias of the ML estimator  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$ ?"

Soln. Note  $\text{Bias}(\hat{\sigma}^2) = E(\hat{\sigma}^2) - \sigma^2$ .

$$\begin{aligned} \text{Then } E[\hat{\sigma}^2] &= E\left(\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n Y_i^2 - 2\bar{Y}Y + \bar{Y}^2\right) \end{aligned}$$

Since  $\bar{Y} = \frac{1}{n} \sum Y_i$ , thus  $\sum Y_i = n\bar{Y}$ . So

$$\begin{aligned} E[\hat{\sigma}^2] &= \frac{1}{n} E\left[\sum_{i=1}^n (Y_i^2) - 2n\bar{Y}\bar{Y} + n\bar{Y}^2\right] \\ &= \frac{1}{n} E\left[\sum_{i=1}^n (Y_i^2) - n\bar{Y}^2\right] \\ &= \frac{1}{n} \left( \sum_{i=1}^n E[Y_i^2] - nE(\bar{Y}^2) \right) \end{aligned}$$

Then, note  $\text{Var}(Y) = E(Y^2) - E(Y)^2$ .

For  $Y \sim \mathcal{G}(\mu, \sigma)$ , thus

$$\therefore \sigma^2 = E(Y^2) - \mu^2, \text{ & so } E(Y^2) = \mu^2 + \sigma^2.$$

Since  $\bar{Y} \sim \mathcal{G}(\mu, \frac{\sigma^2}{n})$ , we can show

$$E(\bar{Y}^2) = \frac{\sigma^2}{n} + \mu^2.$$

Thus

$$\begin{aligned} E[\hat{\sigma}^2] &= \frac{1}{n} \left( \sum_{i=1}^n E[Y_i^2] - nE(\bar{Y}^2) \right) \\ &= \frac{1}{n} \left[ \sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right) \right] \\ &= \frac{1}{n} (n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2) \\ &= \frac{n-1}{n} \sigma^2. \end{aligned}$$

So the bias is

$$\text{Bias}(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}$$

which is not zero!

\* the MLE & mom estimator of the variance slightly underestimates the true variance.

\* note this bias decreases as  $n$  increases.

MEAN SQUARED ERROR / MSE [OF AN ESTIMATOR]

(C4S562)

The "mean squared error" of an estimator is

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

We prefer estimators with a smaller MSE.

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

<< BIAS-VARIANCE DECOMPOSITION OF THE MSE >> (C4S570)

Notes:

- ① Large variance & small bias are both undesirable.
- ② If the estimator is unbiased, then the MSE is just the variance.

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

Then

$$\begin{aligned} E((\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)) &= (E(\hat{\theta} - E(\hat{\theta}))) (E(\hat{\theta}) - \theta) \\ &\quad \text{constant} \\ &= 0 \times (E(\hat{\theta}) - \theta) \\ &= 0. \end{aligned}$$

Proof follows.  $\square$

# EFFICIENCY (C4S575)

SCORE [OF A PARAMETER]:  $U(\theta; Y)$  (C4S578)

The "score" of an unknown parameter  $\theta$  is the gradient of the log-likelihood function; ie

$$U(\theta; Y) = \frac{\partial}{\partial \theta} L(\theta; Y).$$

Notes:

①  $U(\theta; Y)$  is a random variable;

$$\text{② } U(\theta; Y) = \frac{\partial}{\partial \theta} \log L(\theta; Y) = \frac{1}{L(\theta; Y)} \frac{\partial}{\partial \theta} L(\theta; Y);$$

$$\text{③ } E[U(\theta; Y); \theta] = 0.$$

## EXAMPLE: $E(U(\theta))$ OF $\text{Exp}(\theta)$ (C4S580)

Problem:

"Suppose  $Y_1, \dots, Y_n$  are iid  $\text{Exp}(\theta)$  r.v. Show expected value of the score is 0."

Sol1. First, see that

$$L(\theta; Y) = \prod_{i=1}^n \frac{1}{\theta} e^{-\frac{y_i}{\theta}} = \theta^{-n} e^{-\sum_{i=1}^n \frac{y_i}{\theta}}.$$

Thus

$$L(\theta) = -n \log(\theta) - \frac{1}{\theta} \sum_{i=1}^n y_i.$$

Hence

$$U(\theta; Y) = \frac{\partial}{\partial \theta} \left[ -n \log(\theta) - \frac{1}{\theta} \sum_{i=1}^n y_i \right] \\ = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n y_i.$$

So

$$E(U(\theta)) = E\left(-\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n y_i\right) \\ = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n E(Y_i) \\ = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n \theta \\ = -\frac{n}{\theta} + \frac{n\theta}{\theta^2} = 0,$$

as expected.  $\square$

## FISHER INFORMATION [OF A PARAMETER]:

$I(\theta)$  (C4S584)

The "Fisher Information" of an parameter  $\theta$  is the variance of its score; ie

$$I(\theta) = E\left(\left[\frac{\partial}{\partial \theta} \log L(\theta; Y)\right]^2 | \theta\right).$$

We can also write

$$I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \log L(\theta; Y) | \theta\right].$$

Proof. Note that

$$\frac{\partial^2}{\partial \theta^2} \log L(\theta; Y) = \frac{\frac{\partial^2}{\partial \theta^2} L(\theta; Y)}{L(\theta; Y)} - \left(\frac{\partial}{\partial \theta} \log L(\theta; Y)\right)^2$$

Then  $E\left(\frac{\partial^2}{\partial \theta^2} L(\theta; Y) / L(\theta; Y)\right) = 0$ , and so taking expectations of both sides, yields that

$$I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \log L(\theta; Y) | \theta\right]$$

as needed.

Hence, the information tells us about the shape of the log-likelihood.

In turn, the shape of  $L(\theta)$  near the maximum likelihood tells us how many values of  $\theta$  lead to similar values of the log-likelihood itself.

If we have iid rv  $Y_1, \dots, Y_n$ , then if

$$X_1(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \log L(\theta; Y_i) | \theta\right]$$

then

$$X(\theta) = n X_1(\theta).$$

## EXAMPLE: $X(\theta)$ OF $\text{Exp}(\theta)$ (C4S589)

Problem:

"Let  $Y_1, \dots, Y_n$  be iid  $\text{Exp}(\theta)$  rv. Find the Fisher Information."

Sol1. Earlier we showed that

$$\frac{\partial}{\partial \theta} \log L(\theta; Y) = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n Y_i.$$

Thus

$$\frac{\partial^2}{\partial \theta^2} \log L(\theta; Y) = \frac{n}{\theta^2} - \frac{2}{\theta^3} \sum_{i=1}^n Y_i.$$

Since  $E[Y_i] = \theta$ , thus

$$X(\theta) = -E\left[\frac{n}{\theta^2} - \frac{2}{\theta^3} \sum_{i=1}^n Y_i | \theta\right] \\ = -E\left[\frac{n}{\theta^2} - \frac{2n\theta}{\theta^3}\right]$$

$$\therefore X(\theta) = \frac{n}{\theta^2}.$$

$$\text{Var}(\tilde{\theta}) \geq \frac{1}{X(\theta)}$$

<< CRAMER-RAO LOWER BOUND >>

(C4S594)

let  $\tilde{\theta}$  be an unbiased estimator. Then necessarily

$$\text{Var}(\tilde{\theta}) \geq \frac{1}{X(\theta)}.$$

MINIMUM-VARIANCE UNBIASED ESTIMATOR / MVUE (C4S594)

A "minimum-variance unbiased estimator" is  $\tilde{\theta}$  such that

$$\text{Var}(\tilde{\theta}) = \frac{1}{X(\theta)}.$$

EFFICIENCY [OF AN UNBIASED ESTIMATOR]:  $e(\theta)$  (C4S600)

The "efficiency" of an unbiased estimator  $\tilde{\theta}$  of a parameter  $\theta$  is

$$e(\tilde{\theta}) = \frac{1}{\text{Var}(\tilde{\theta})}.$$

If  $\text{Var}(\tilde{\theta}) = \frac{1}{X(\theta)}$ , then  $e(\tilde{\theta}) = 1$ , and in this case we say the estimator is "efficient".

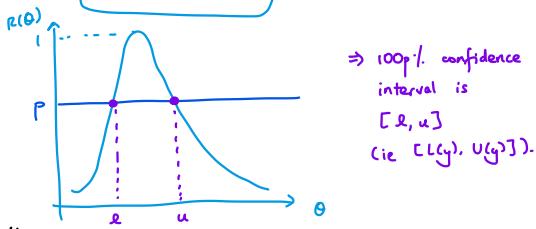
# INTERVAL ESTIMATION (C4S602)

## LIKELIHOOD INTERVAL (C4S614)

Q<sub>1</sub>: A "100% likelihood interval" for the parameter

$\theta$  is the set

$$\{\theta \mid R(\theta) \geq p\}.$$



$\Rightarrow$  100% confidence interval is  $[l, u]$  (ie  $[L_g, U_g]$ ).

Q<sub>2</sub>: Interpretations: if the data is in a

- ① 50% likelihood interval  $\rightarrow$  very plausible
- ② 10% likelihood interval  $\rightarrow$  plausible
- ③ 5% likelihood interval  $\rightarrow$  implausible
- ④ 1% likelihood interval  $\rightarrow$  very implausible

Q<sub>3</sub>: In particular, increasing the sample size  $n$  decreases the width of the likelihood intervals.  
 $\rightarrow$  dist<sup>n</sup> becomes narrower as  $n \uparrow$

## LOG RELATIVE LIKELIHOOD FUNCTION:

### $r(\theta)$ (C4S627)

Q<sub>1</sub>: The "log relative likelihood function" of  $\theta$  is

$$r(\theta) = \log R(\theta) = l(\theta) - l(\hat{\theta}).$$

Q<sub>2</sub>: To obtain a 100% likelihood interval, we plot  $r(\theta)$  and draw a line at

$$r(\theta) = \log(p).$$

# CONFIDENCE INTERVALS & PIVOTAL QUANTITIES (C4S632)

## COVERAGE PROBABILITIES [OF INTERVAL ESTIMATORS] (C4S633)

$\Theta_1$ : let  $Y = (Y_1, \dots, Y_n)$  be the potential data to be collected.

$\Theta_2$ : let  $[L(Y), U(Y)]$  be an "interval estimator" which can be used to construct the possible values  $\theta$  can take.

$\Theta_3$ : Then, the "coverage probability" for the interval estimator  $[L(Y), U(Y)]$  is equal to

$$P(\theta \in [L(Y), U(Y)]) = P(L(Y) \leq \theta \leq U(Y)).$$

$\Theta_4$ : We choose  $L(Y), U(Y)$  such that

- ① The coverage probability is large; &
  - ② The interval is as narrow as possible.
- } but these conflict!

$\Theta_5$ : Usually, we fix the coverage probability and try to find the narrowest interval.

## CONFIDENCE INTERVAL & COEFFICIENT (C4S640)

$\Theta_1$ : A "100p% confidence interval" for a parameter  $\theta$  is an interval estimate  $[L(\bar{y}), U(\bar{y})]$  such that

$$P(\theta \in [L(\bar{y}), U(\bar{y})]) = P(L(\bar{y}) \leq \theta \leq U(\bar{y})) = p.$$

$\Theta_2$ :  $p$  is called the "confidence coefficient" of the interval.

$\Theta_3$ : Note that

①  $\theta$  is an unknown constant of the population, not a random variable.

★ ② So, we cannot say "the probability  $\theta$  lies between  $L(\theta)$  &  $U(\theta)$  is  $p$ ".

$\Theta_4$ : But, we can say we are "100p% confident" that the interval contains the true (and unknown) value of  $\theta$ .

$\Theta_5$ : Note greater confidence corresponds to a wider confidence interval!

## PIVOTAL QUANTITY (C4S652)

$\Theta_1$ : A "pivotal quantity"  $Q = Q(Y; \theta)$  is a function of the data  $Y$  & the unknown parameter  $\theta$  such that the distribution of the random variable  $Q$  is completely known.

$$\text{eg } \frac{\bar{Y} - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

$\Theta_2$ : To use a pivotal quantity to construct a confidence interval:

- ① Determine numbers  $a, b$  such that  $P(a \leq Q(Y; \theta) \leq b) = p$
- ② Re-express  $a \leq Q(Y; \theta) \leq b$  in the form  $L(Y) \leq \theta \leq U(Y)$ ;
- ③ Then

$$p = P(L(Y) \leq \theta \leq U(Y)) = P(a \leq Q(Y; \theta) \leq b).$$

- ④ For observed data  $y$ , the interval  $[L(y), U(y)]$  is a 100p% confidence interval for  $\theta$ .

## GAUSSIAN DATA (C4S673)

$\Theta_1$ : Let  $Z \sim \mathcal{N}(0, 1)$ . Then, a 100p% confidence interval for a sample size of  $n$  is

$$(\bar{y} - \frac{\sigma}{\sqrt{n}}, \bar{y} + \frac{\sigma}{\sqrt{n}}),$$

where

$$P(Z \leq a) = \frac{1+p}{2}.$$

in R:  
 $qnorm((1+p)/2)$

## TWO-SIDED, EQUAL-TAILED CIS FOR $\mu$ (C4S674)

$\Theta_1$ : A 100p% confidence interval for  $\mu$  is of the form

$$\text{point estimate} \pm \text{distribution quantile} \times \text{sd(estimate)}.$$

\* note not all CIs are symmetric in general!

## ASYMPTOTIC / APPROXIMATE PIVOTAL QUANTITY (C4S682)

$\Theta_1$ : An "asymptotic/approximate pivotal quantity" is a set of random variables  $Q_n = Q_n(Y_1, \dots, Y_n; \theta)$  such that as  $n \rightarrow \infty$ , the distribution of  $Q_n$  ceases to depend on  $\theta$  or other unknown information.

$\Theta_2$ : These can be used to construct approximate CIs for  $\theta$ .

EXAMPLE:  $\text{Bin}(n, \theta)$ ,  $\hat{\theta} = \frac{\bar{Y}}{n}$  (C4S686)

Problem:

"Let  $Y \sim \text{Bin}(n, \theta)$ ,  $\hat{\theta} = \frac{\bar{Y}}{n}$ . Find an approximate 95% CI for  $\theta$ ."

Sol2. By CLT,

$$\frac{\hat{\theta} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \sim \mathcal{N}(0, 1) \text{ approximately.}$$

Moreover,

$$Q_n = \frac{\hat{\theta} - \theta}{\sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}} \sim \mathcal{N}(0, 1) \text{ approximately.}$$

Then since

$$0.95 = P(-1.96 \leq Z \leq 1.96)$$

Hence

$$0.95 = P(-1.96 \leq \frac{\hat{\theta} - \theta}{\sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}} \leq 1.96)$$

thus

$$0.95 \approx P(\hat{\theta} - 1.96 \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}} \leq \theta \leq \hat{\theta} + 1.96 \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}})$$

and so an approximate 95% CI is

$$\hat{\theta} \pm 1.96 \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}.$$

## INTERVAL ESTIMATION (C4S687)

$\Theta_1$ : Ways of finding interval estimates for an unknown parameter:

- ① Use a 100p% likelihood interval-
- ② Use a 100p% confidence interval if an exact pivotal quantity exists; or
- ③ Use a 100p% approximate confidence interval based on an approximate pivotal quantity (usually using CLT).

## SAMPLE SIZE CALCULATION (C4S695)

$\Theta_1$  Suppose we want to estimate  $\theta$ , the proportion of units in a large population who have a specific characteristic, and we plan to select  $n$  units at random.

$\Theta_2$  Suppose we use the 100% CI.

$$\hat{\theta} \pm a \frac{s}{\sqrt{n}}$$

$\Theta_3$  We can specify we want a CI of width  $\leq 2l$ ; ie

$$a \frac{s}{\sqrt{n}} \leq l,$$

or

$$n \geq \left(\frac{a}{l}\right)^2 s^2,$$

which tells us the minimum value of  $n$  needed for the CI to be of width at most  $2l$ .

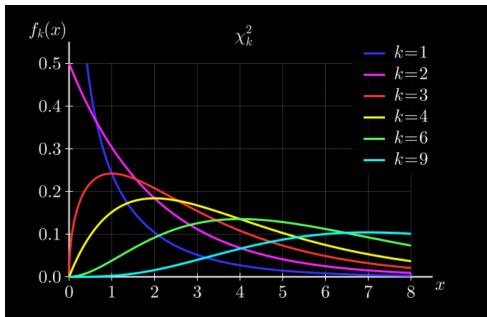
## CHI-SQUARED DISTRIBUTION: $\chi_k^2$ (C4S701)

$\Theta_1$  The "chi-squared distribution" is parameterized by its degrees of freedom  $k$ .

$\Theta_2$   $k$  affects the shape of the resulting pdf:

$$\text{pdf} = \frac{1}{2^{k/2} \Gamma(k/2)} y^{\frac{k}{2}-1} e^{-\frac{y}{2}}.$$

\* not needed to know.



$\Theta_3$  Properties:

① If  $W_1, W_2, \dots, W_n$  are iid with  $W_i \sim \chi_{k_i}^2$ , then

$$S = \sum_{i=1}^n W_i \sim \chi_{\sum k_i}^2.$$

② If  $Z \sim \mathcal{N}(0,1)$ , then

$$Z^2 \sim W \sim \chi_1^2.$$

③ If  $Z_1, \dots, Z_n \sim \mathcal{N}(0,1)$ , then

$$S = \sum_{i=1}^n Z_i^2 \sim \chi_n^2.$$

④ Also,

$$W \sim \chi_2^2 = \text{Exp}(2).$$

## LIKELIHOOD RATIO STATISTIC: $\Lambda(\theta)$ (C4S716)

$\Theta_1$  The "likelihood ratio statistic" is defined to be

$$\Lambda(\theta) = -2 \log \left( \frac{L(\theta)}{L(\hat{\theta})} \right) = -2 \log \left( \frac{L(\theta; Y)}{L(\hat{\theta}; Y)} \right).$$

\*  $\Lambda$  is a random variable!

$\Theta_2$  For large enough  $n$ , we can show

$$\Lambda \sim \chi^2.$$

## LIKELIHOOD BASED CONFIDENCE INTERVAL (C4S724)

$\Theta_1$  Note a 100% likelihood interval is an approximate

100% confidence interval, where

$$q = P(W \leq -2 \log(p)), \quad W \sim \chi_1^2.$$

See slides for proof.

$\Theta_2$  In particular, since  $\Lambda(\theta) \sim \chi^2$  for large  $n$ , the likelihood interval can be written like

$$\{\theta : R(\theta) \geq p\} = \{\theta : -2 \log \left[ \frac{L(\theta)}{L(\hat{\theta})} \right] \leq -2 \log(p)\}.$$

$\Theta_3$  Hence, the confidence coefficient is

$$\begin{aligned} P(\Lambda(\theta) \leq -2 \log(p)) &\approx P(W \leq -2 \log(p)) \\ &= P(Z \leq \sqrt{-2 \log(p)}) \\ &= 2P(Z < \sqrt{-2 \log(p)}) + 1. \end{aligned}$$

$$* P(W \leq c) = 2P(Z \leq \sqrt{c}) - 1.$$

# GAUSSIAN DATA: UNKNOWN $\mu$ & $\sigma$ (C4S74)

$\text{Q}_1$  Let  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ , where  $\mu$  &  $\sigma$  are unknown.

$\text{Q}_2$  We can use the MLE estimator for  $\mu$ :

$$\hat{\mu} = \bar{Y}.$$

$\text{Q}_3$  We use the point estimator for  $\sigma^2$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

- we prefer  $S^2$  since  $E(S^2) = \sigma^2$ .

## t-DISTRIBUTION: $t_k$ (C4S74)

$\text{Q}_1$  A rv  $T$  is said to have a "Student's t" distribution if its pdf is

$$f(t; k) = c_k \left(1 + \frac{t^2}{k}\right)^{-\frac{k+1}{2}},$$

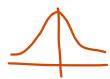
where

$$c_k = \frac{\Gamma(\frac{k+1}{2})}{\sqrt{k\pi} \Gamma(\frac{k}{2})}.$$

$\text{Q}_2$  The parameter  $k$  is called the "degrees of freedom", and we write  $T \sim t_k$  or  $T \sim t(k)$ .

$\text{Q}_3$  Notes:

- ① The t distribution is unimodal and symmetric about 0;
- ② For large  $k$ ,  $t_k \approx \mathcal{N}(0, 1)$ .



$$Z \sim \mathcal{N}(0, 1), U \sim \chi_k^2 \Rightarrow \frac{Z}{\sqrt{U/k}} \sim t_k \quad (\text{C4S754})$$

$\text{Q}_1$  Let  $Z \sim \mathcal{N}(0, 1)$  &  $U \sim \chi_k^2$  be independent. Then

$$T \sim \frac{Z}{\sqrt{U/k}}$$

has a t-distribution with  $k$  degrees of freedom.

$$Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2), \mu, \sigma \text{ UNKNOWN} \Rightarrow$$

$$\frac{Y - \mu}{S/\sqrt{n}} \sim t_{n-1} \quad (\text{C4S757})$$

$\text{Q}_1$  First, see that if

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

then

$$U = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

$\text{Q}_2$  In particular, if  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$  where  $\mu$  &  $\sigma$  are unknown, then

$$\frac{Y - \mu}{S/\sqrt{n}} \sim t_{n-1}.$$

\*this is a pivotal quantity!

## CONFIDENCE INTERVAL FOR $\mu$ IF $\sigma$ IS UNKNOWN (C4S760)

$\text{Q}_1$  Let  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$  be iid where  $\sigma$  is unknown.

Then necessarily a 100% CI for  $\mu$  is

$$(\bar{Y} - a \frac{S}{\sqrt{n}}, \bar{Y} + a \frac{S}{\sqrt{n}})$$

\*don't forget the  $\underline{n-1}$ !!

where  $P(T \leq a) = \frac{1-p}{2}$ , where  $T \sim t_{n-1}$ .

\*in R, the command "pt(b, df)" returns  $P(T \leq b)$ , where  $T \sim t_{df}$ .

\*the command "qt(t, df)" returns  $t$  s.t.  $P(T \leq t) = q$ , where  $T \sim t_{df}$ .

## HOW PARAMETERS AFFECT WIDTH OF CI (C4S788)

$\text{Q}_1$  The width of the CI is  $2a \frac{S}{\sqrt{n}}$ , where  $P(T \leq a) = \frac{1-p}{2}$ ,  $T \sim t_{n-1}$ .

$\text{Q}_2$  Note that

- ① ↑ confidence level  $\Rightarrow$  new CI is wider
- ② ↑ sample size  $\Rightarrow$  new CI is narrower  
as  $k \uparrow$ ,  $t_k$  becomes less concentrated at peak
- ③ ↑ sample std dev  $\Rightarrow$  new CI is wider
- ④ ↑ (or ↓) sample mean  $\Rightarrow$  new CI's width is unchanged

## SAMPLE SIZE CALCULATION (C4S789)

$\text{Q}_1$  In these, we assume  $\sigma$  is known.

- since 's' depends on  $n$ .

$\text{Q}_2$  Our CI is thus

$$\bar{Y} \pm a \frac{\sigma}{\sqrt{n}}, \quad P(Z \leq a) = \frac{1-p}{2}.$$

\*we assume population std dev = sample std dev.

If we want this to have width  $2R$ , then we choose  $n$  such that

$$n \approx \left(\frac{a\sigma}{R}\right)^2.$$

## CI FOR $\sigma^2$ (C4S796)

$\text{Q}_1$  Let  $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$  be iid, where  $\mu$  &  $\sigma$  are unknown. Then we have shown

$$W = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

$\text{Q}_2$  We pick  $a, b$  such that

$$P(W \leq a) = \frac{1-p}{2}, \quad P(W \geq b) = \frac{1-p}{2}$$

or in other words

$$P(a \leq W \leq b) = p.$$

\*since  $\chi^2$  is not symmetric.

$\text{Q}_3$  Thus, our coverage is

$$P\left(a \leq \frac{(n-1)S^2}{\sigma^2} \leq b\right) = p.$$

which can be rearranged to

$$P\left(\frac{(n-1)S^2}{b} \leq \sigma^2 \leq \frac{(n-1)S^2}{a}\right) = p.$$

$\text{Q}_4$  Thus, a 100% CI for  $\sigma^2$  is

$$\left(\frac{(n-1)S^2}{b}, \frac{(n-1)S^2}{a}\right).$$

\*this is not symmetric!

\*in R, we use `pchisq(w, df)` & `qchisq(w, df)`.

$\text{Q}_5$  A 100% CI for  $\sigma$  is

$$\left(\sqrt{\frac{(n-1)S^2}{b}}, \sqrt{\frac{(n-1)S^2}{a}}\right).$$