

STAT 241

Personal Notes

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



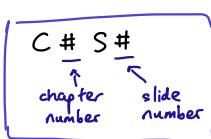
Marcus Chan

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
 ① the estimation of attributes; and
 ② 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

- ① In symmetrical distributions, $\bar{y} \approx \hat{m}$;
- but
- ② 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

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

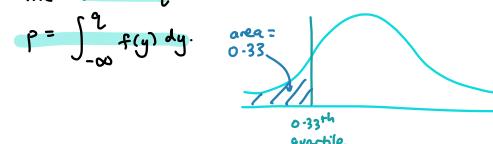
* the range is very susceptible to outliers!

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

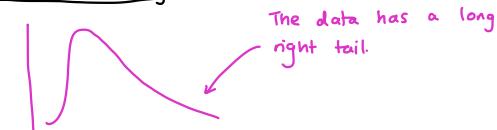
$$\text{sample skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^3}{\left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{3/2}}$$

\exists_2 Interpretation of sample skewness's value:

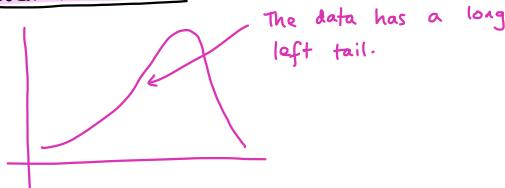
- ① If $ss = 0 \Rightarrow$ distribution is symmetric; eg Gaussian, uniform



- 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

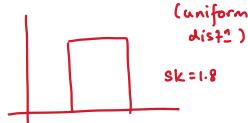
$$\text{sample kurtosis} = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^4}{\left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^2}$$

\exists_2 Interpretation of sample kurtosis' value:

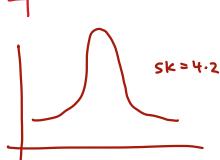
- ① $sk = 3 \Rightarrow$ distribution looks "Gaussian" (bell-shaped);



- 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 $\approx 95\%$ 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)

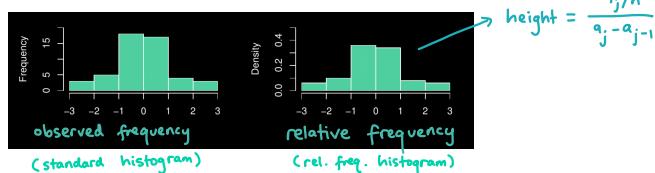
Essentially, histograms create a graphical summary of our data that we can use to compare with a pdf for crvs, or a pmf for a drv.

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.



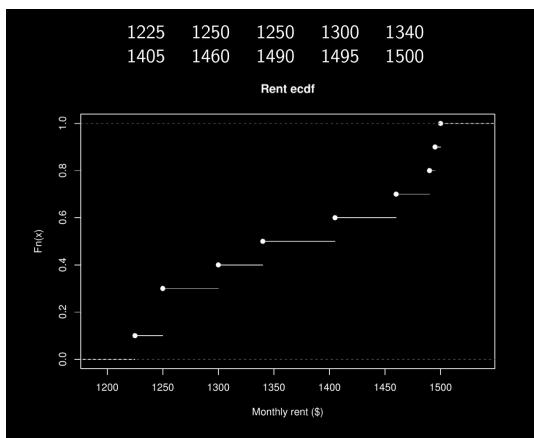
EMPIRICAL CDF (CIS124)

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

Mathematically, the empirical cdf is defined

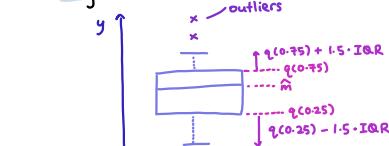
by

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

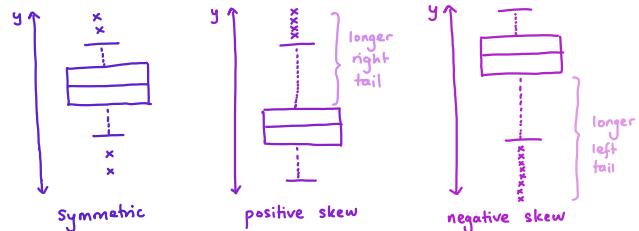


BOX-PLOT (CIS139)

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

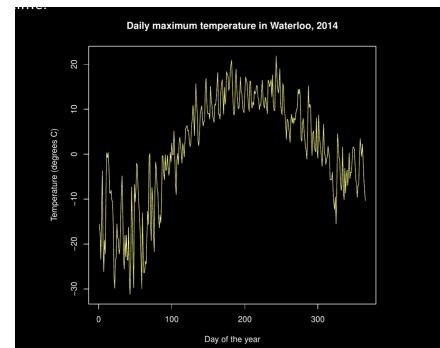


Box-plots can also show the skewness of a distribution:



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)

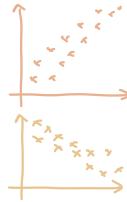
The "sample correlation", denoted " r ", gives us a numerical summary of a bivariate dataset.

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	Σ
Hockey ✓	33	9	42
Hockey X	22	43	65
Σ	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 θ (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(θ)
- ② Exponential(θ)
* " θ " = mean of the distribution
(not $\frac{1}{\text{mean}}$).
- ③ Binomial(n, θ)
- ④ Gaussian(θ) = Gaussian(μ, σ)
- ⑤ Multinomial($n, \theta_1, \dots, \theta_n$)
- ⑥ Geometric(θ)

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 θ .

Thus, we say the random variable is "parameterized" by θ .

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 θ .

ESTIMATION OF UNKNOWN PARAMETERS (C2S206)

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

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

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

This process is referred to as "estimating" the value of θ .

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

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

PARAMETER SPACE: Ω

The "parameter space" Ω of a parameter θ is the set of all possible values θ can take.

LIKELIHOOD FUNCTION [FOR DRV]

$L(\theta)$ (C2S224)

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

Suppose y is a drv.

Then, the "likelihood function for θ " is defined to be

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

where Ω is the parameter space of θ .

* L is technically a function of θ & 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 θ 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 θ 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!

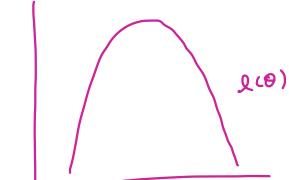
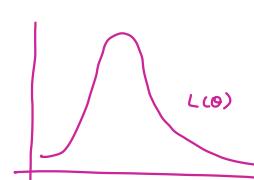
Note that $\ell(\theta)$ is maximized for the same value of θ as the regular likelihood function.

* 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 ℓ (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 θ 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 θ 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 θ 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

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

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

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 θ .

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 θ , 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 θ using the MLE of θ , 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.

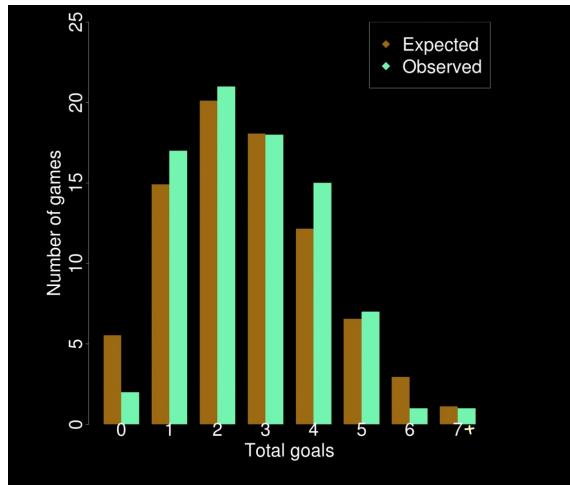
Since the range of Poi is technically $0, 1, 2, \dots$, we need to account for the "right tail" by grouping all the values ≥ 7 into "one" value:

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

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

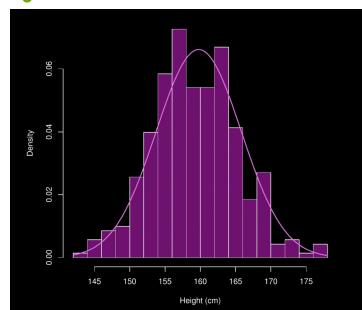
We may also plot the expected/observed frequencies via a bar plot.



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

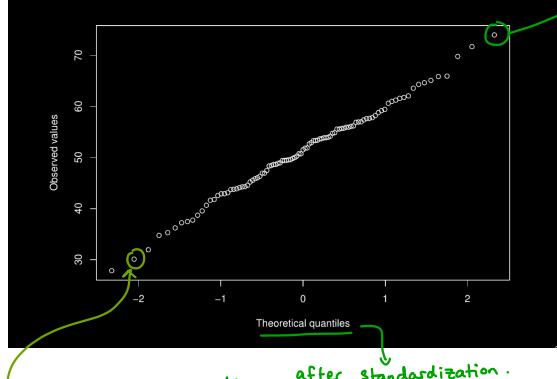
> pnorm(162, 159.77, 6.03) - pnorm(160, 159.77, 6.03).

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

$e_j = 351p_j$, where I_j is the j^{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^{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

"leveling out"
⇒ very little data in the tails
⇒ data concentrated in the peak
⇒ ie low kurtosis

high kurtosis, symmetric

"increasing near the tails"
⇒ data concentrated more in the tails
⇒ ie high kurtosis

- 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

⇒ more data towards the right (ie long right tail)
⇒ data is skewed to the right
⇒ ie positive skew

negative skew, asymmetric

⇒ more data towards the left (ie long left tail)
⇒ data is skewed to the left
⇒ ie negative skew

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^{th} POPULATION MOMENT: μ_k (C2S353)

The " k^{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^{th} SAMPLE MOMENT: m_k (C2S355)

Let y_1, \dots, y_n be a sample.

Then, the " k^{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 μ and σ ".

Solⁿ. 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 μ & σ 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) ?"

Solⁿ. 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^{ts}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 ."

Solⁿ. 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

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

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

Then, we estimate μ_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

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.

e.g. in the vaccine study, possible target popⁿs/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.

e.g. 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.

e.g. in the vaccine study, possible attributes include

- ① the proportion of people in the target popⁿ who would contract COVID-19 after receiving the Pfizer vaccine within 24 weeks;
- ② the proportion of people in the target popⁿ 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 ≥ 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!

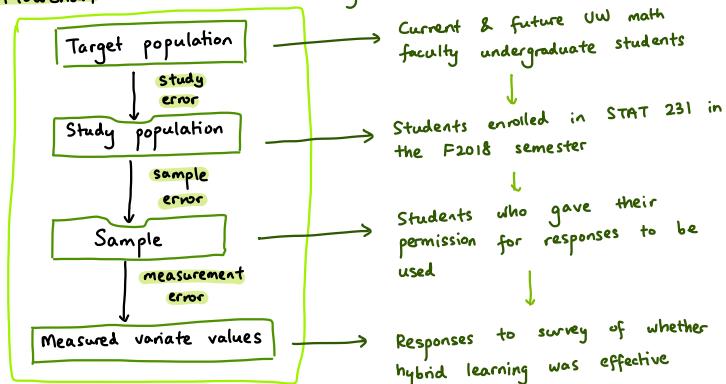
MEASUREMENT ERROR (C3S442)

"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:



eg STAT 231 example (C3S447)

Current & future UW math faculty undergraduate students

↓
Students enrolled in STAT 231 in the F2018 semester

↓
Students who gave their permission for responses to be used

↓
Responses to survey of whether hybrid learning was effective

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 θ ?"

💡 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)
 $y = \# \text{ of } X \text{ in a randomly chosen sample from}$

$\text{the target pop}^n \text{ who have had COVID-19.}$

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 popⁿ, 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 θ .

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, σ	↑ mean, μ
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 μ

- ① increases as n increases;
- ② decreases as σ increases; &
- ③ does not change with μ .

* 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 μ ;
- ② $sd(\bar{Y})$ increases as σ increases, and so less of our sample estimates will be closer to μ ;
- ③ $sd(\bar{Y})$ does not change with μ . (C4S540)

NORMAL APPROXIMATIONS (C4S525)

If Y_1, \dots, Y_n are iid with mean μ & variance σ^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 θ 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 θ 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 θ 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 θ 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".

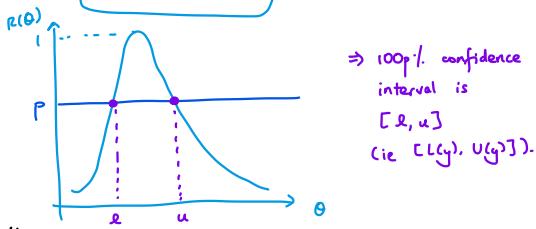
INTERVAL ESTIMATION (C4S602)

LIKELIHOOD INTERVAL (C4S614)

Q₁: A "100% likelihood interval" for the parameter

θ is the set

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



\Rightarrow 100% confidence interval is $[l, u]$ (ie $[L(y), U(y)]$).

Q₂: 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

Q₃: In particular, increasing the sample size n decreases the width of the likelihood intervals.
 \rightarrow distⁿ becomes narrower as $n \uparrow$

LOG RELATIVE LIKELIHOOD FUNCTION:

$r(\theta)$ (C4S627)

Q₁: The "log relative likelihood function" of θ is

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

Q₂: 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)

Θ_1 : let $Y = (Y_1, \dots, Y_n)$ be the potential data to be collected.

Θ_2 : let $[L(Y), U(Y)]$ be an "interval estimator" which can be used to construct the possible values θ can take.

Θ_3 : Then, the "coverage probability" for the interval estimator $[L(Y), U(Y)]$ is equal to

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

Θ_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!

Θ_5 : Usually, we fix the coverage probability and try to find the narrowest interval.

CONFIDENCE INTERVAL & COEFFICIENT (C4S640)

Θ_1 : A "100p% confidence interval" for a parameter θ is an interval estimate $[L(\bar{y}), U(\bar{y})]$ such that

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

Θ_2 : p is called the "confidence coefficient" of the interval.

Θ_3 : Note that

① θ is an unknown constant of the population, not a random variable.

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

Θ_4 : But, we can say we are "100p% confident" that the interval contains the true (and unknown) value of θ .

Θ_5 : Note greater confidence corresponds to a wider confidence interval!

PIVOTAL QUANTITY (C4S652)

Θ_1 : A "pivotal quantity" $Q = Q(Y; \theta)$ is a function of the data Y & the unknown parameter θ such that the distribution of the random variable Q is completely known.

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

Θ_2 : To use a pivotal quantity to construct a confidence interval:

① Determine numbers a, b such that $P(a \leq Q(Y; \theta) \leq b) = p$

* not necessarily symmetric
But try to find narrowest interval.

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 θ .

GAUSSIAN DATA (C4S673)

Θ_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 μ (C4S674)

Θ_1 : A 100p% confidence interval for μ 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)

Θ_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 θ or other unknown information.

Θ_2 : These can be used to construct approximate CIs for θ .

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 θ ."

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)

Θ_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)

Θ_1 Suppose we want to estimate θ , the proportion of units in a large population who have a specific characteristic, and we plan to select n units at random.

Θ_2 Suppose we use the 100% CI.

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

Θ_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 $2l$.

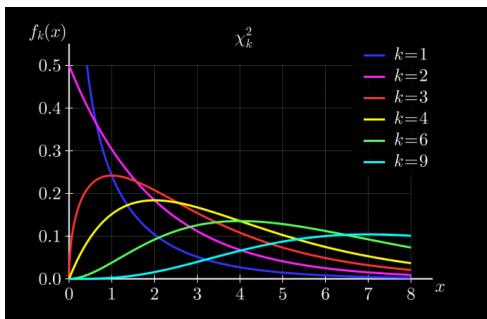
CHI-SQUARED DISTRIBUTION: χ_k^2 (C4S701)

Θ_1 The "chi-squared distribution" is parameterized by its degrees of freedom k .

Θ_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.



Θ_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)

Θ_1 The "likelihood ratio statistic" is defined to be

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

* Λ is a random variable!

Θ_2 For large enough n , we can show

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

LIKELIHOOD BASED CONFIDENCE INTERVAL (C4S724)

Θ_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.

Θ_2 In particular, since $\Lambda(\theta) \sim \chi^2$ for large n , the likelihood interval can be written like

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

Θ_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 μ & σ (C4S74)

Q1 Let $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$, where μ & σ are unknown.

Q2 We can use the MLE estimator for μ :

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

Q3 We use the point estimator for σ^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)

Q1 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

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

Q2 The parameter k is called the "degrees of freedom", and we write $T \sim t_k$ or $T \sim t(k)$.

Q3 Notes:

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



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

Q1 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.

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

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

Q1 First, see that if

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

then

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

Q2 In particular, if $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ where μ & σ are unknown, then

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

*this is a pivotal quantity!

CONFIDENCE INTERVAL FOR μ IF σ IS UNKNOWN (C4S760)

Q1 Let $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ be iid where σ is unknown.

Then necessarily a 100% CI for μ 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)

Q1 The width of the CI is $2a\frac{s}{\sqrt{n}}$, where

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

Q2 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)

Q1 In these, we assume σ is known.

- since 's' depends on n .

Q2 Our CI is thus

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

*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 σ^2 (C4S796)

Q1 Let $Y_1, \dots, Y_n \sim \mathcal{N}(\mu, \sigma^2)$ be iid, where μ & σ are unknown. Then we have shown

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

Q2 We pick a, b such that

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

or in other words

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

*since χ^2 is not symmetric.

Q3 Thus, our coverage is

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

which can be rearranged to

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

Q4 Thus, a 100% CI for σ^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)`.

Q5 A 100% CI for σ is

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

Chapter 5: Tests of Hypothesis

NULL HYPOTHESIS: H_0 (CSS825)

- 💡 The "null hypothesis" is a single "default" hypothesis.
eg "the defendant is innocent"
- 💡 Hypothesis testing is based on collecting data, and based on said data determining how plausible H_0 is.

ALTERNATIVE HYPOTHESIS: H_A/H_1 (CSS826)

- 💡 The "alternative hypothesis" is the alternative to the null.
Often, H_A is just the negation of H_0 .
eg "the defendant is guilty"

TEST STATISTIC / DISCREPANCY MEASURE (CSS833)

- 💡 A "test statistic" is a function of the data $D = g(Y)$ that is constructed to measure the "agreement" between the data Y & H_0 .

Note:

- ① D is a random variable.
- ② If we observe $Y=y$, we use $d=g(y)$ to denote the observed value of D .
- 💡 We usually define D so $d=0$ represents the best possible agreement between Y & H_0 .
- 💡 Note "large" values of d indicate poor agreement between Y & H_0 .

p-VALUE (CSS845)

- 💡 Suppose we use the test statistic $D=D(Y)$ to test the hypothesis H_0 . Let $d=g(y)$ be the observed value of D . Then, the "p-value" of H_0 using D is $P(D \geq d; H_0)$.

The p-value is the probability of observing a value of the test statistic greater than or equal to the observed value of the test statistic assuming H_0 is true.

- 💡 In particular, a small p-value tells us that if H_0 were true, it would be unlikely to have observed data at least as surprising as the data we actually observed.

STEPS OF A HYPOTHESIS TEST (CSS849)

- 💡 Steps:
 - Specify H_0 to be tested using data Y .
 - Define a test statistic $D(Y)$, where large values of D imply the data is less "consistent" with H_0 .
 - Let $d=D(y)$; ie the 'observed value' of D .
 - Calculate the p-value $P(D \geq d; H_0)$.
 - Draw a conclusion based on the p-value.

INTERPRETING THE p-VALUE (CSS852)

- 💡 If $d=D(y)$ is large, and thus the p-value $P(D \geq d; H_0)$ is small, then either
 - ① H_0 is true, but by chance we observed an event that is very unlikely when H_0 is true; or
 - ② H_0 is false.
- 💡 What does a "small" p-value mean?

p-value There is — evidence against H_0 based on the data.
0.1 < p no
0.05 < p ≤ 0.1 weak
0.01 < p ≤ 0.05 some
0.001 < p ≤ 0.01 strong
p ≤ 0.001 very strong

* these are only guidelines!

- 💡 Depending on the p-value, we may state we "reject", or fail to "reject", the null hypothesis.

* we never accept H_0 (or H_1)!

TYPE I & II ERRORS: α , β (CSS861)

- 💡 A "type I error" is the probability we reject H_0 when it is actually true.
(ie false positive)
- 💡 A "type II error" is the probability we fail to reject H_0 when it is actually false.
(ie false negative)

POWER [OF A TEST] (CSS864)

- 💡 The "power" of a test is $1-\beta$, where β is the corresponding type 2 error.
- 💡 A more powerful test is more desirable.
- 💡 In particular,

$$\text{power} = P(\text{reject } H_0 \mid H_0 \text{ is false}).$$

STEPS ON COMPUTING THE POWER OF A TEST (CSS871)

- 💡 Steps to finding a test's power of $H_0: \theta = \theta_0$ against an alternative of $\theta \neq \theta_0$ at significance level α :
 - Identify the "rejection region"; ie the test statistics that would lead us to reject H_0 .
 - For a specified value of $\theta \neq \theta_0$, compute the probability a sample would yield a test statistic in the rejection region.

p-HACKING (CSS878)

"p-hacking" is repeating experiments or being selective with one's results to falsely engineer a "significant" result.

ONE-SIDED TEST (CSS881)

Q₁: A "one-sided test" is a hypothesis test where

$$\begin{aligned} H_0: \theta &= \theta' \\ H_1: \theta &> \theta' \quad (\text{or } \theta < \theta'). \end{aligned}$$

Q₂: We may use the test statistic

$$D = \max \{ Y - \theta', 0 \}.$$

*symmetric for case where $\theta < \theta'$.

Our p-value is thus $P(D \geq d)$, where
 $d = y - \theta'$.

HYPOTHESIS TESTING FOR $G(\mu, \sigma)$ PARAMETERS

TESTING $H_0: \mu = \mu_0$, σ UNKNOWN

(CSS895)

\exists_1 Let $Y_1, \dots, Y_n \sim G(\mu, \sigma)$ be a random sample, where σ is unknown.

Recall

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

\exists_2 To test $H_0: \mu = \mu_0$, we use the test statistic

$$D = \frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}}.$$

Why? Notice $E[\bar{Y}] = \mu_0$ if H_0 is true.

Our question is "is $D=d$ surprisingly large?" ie our p-value is

$$\begin{aligned} P(D \geq d; H_0) &= P\left(\frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}} \geq d; H_0\right) \\ &= P(|T| \geq d), \quad T \sim t_{n-1} \\ &= 1 - P(-d \leq T \leq d). \end{aligned}$$

\exists_3 Our p-value is thus

$$p = 1 - P(-d \leq T \leq d)$$

where $T \sim t_{n-1}$.

\exists_4 In R, we may use

t.test(y, mu=1)

TESTS OF HYPOTHESIS & CIs (CSS914)

\exists_1 The p-value for testing $H_0: \theta = \theta_0$ is $\geq p$ iff the value $\theta = \theta_0$ is inside a $100(1-p)\%$ CI (using the same pivotal quantity).

e.g. $H_0: \mu = \mu_0$ for $G(\mu, \sigma)$ data.

$$\begin{aligned} \text{Then } p\text{-value} \geq 0.05 &\Leftrightarrow P\left(\frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}} \geq \frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}}\right) \geq 0.05 \\ &\Leftrightarrow P\left(|T| \geq \frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}}\right) \geq 0.05, \quad T \sim t_{n-1} \\ &\Leftrightarrow P\left(|T| \leq \frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}}\right) \leq 0.95 \\ &\Leftrightarrow \frac{|\bar{Y} - \mu_0|}{S/\sqrt{n}} \leq a, \quad 0.95 = P(|T| \leq a) \\ &\Leftrightarrow \mu_0 \in [\bar{Y} - a \frac{s}{\sqrt{n}}, \bar{Y} + a \frac{s}{\sqrt{n}}] \end{aligned}$$

which is a 95% CI for μ .

POWERING A STUDY (CSS919)

\exists_1 Usually, we fix σ to a certain value and ask what sample size we require to attain a specified level of power.

\exists_2 For $\alpha = 0.05$, the power is

$$P(\bar{Y} < \mu_0 - 1.96 \frac{s}{\sqrt{n}}) + P(\bar{Y} > \mu_0 + 1.96 \frac{s}{\sqrt{n}}).$$

So we should find n s.t. $1 - \beta = P(\bar{Y} < \mu_0 - 1.96 \frac{s}{\sqrt{n}}) + P(\bar{Y} > \mu_0 + 1.96 \frac{s}{\sqrt{n}})$.

\exists_3 More generally, to obtain power $1 - \beta$ we seek n s.t.

$$P(\bar{Y} < \mu_0 - z_{\alpha/2} \frac{s}{\sqrt{n}}) = 1 - \beta,$$

where $\bar{Y} \sim G(\mu, \frac{\sigma}{\sqrt{n}})$, μ & σ are known, and $z_{\alpha/2}$ is such that $P(z < a) = \frac{a}{2}$ for $z \sim N(0, 1)$.

\exists_4 This is equivalent to

$$P\left(z = \frac{\bar{Y} - \mu}{s/\sqrt{n}} < \frac{\mu_0 - z_{\alpha/2} \frac{s}{\sqrt{n}} - \mu}{s/\sqrt{n}}\right) = 1 - \beta$$

which can be rearranged to

$$n \geq \left(\frac{\sigma(z_{\alpha/2} + z_{1-\beta})}{\mu_0 - \mu}\right)^2, \quad P(z < a) = 1 - \beta.$$

$$\begin{aligned} \text{Proof: } P\left(z < \frac{\mu_0 - z_{\alpha/2} \frac{s}{\sqrt{n}} - \mu}{s/\sqrt{n}}\right) &= 1 - \beta \Rightarrow z_{1-\beta} < \frac{\mu_0 - \mu}{s/\sqrt{n}} - z_{\alpha/2} \\ &\Rightarrow \frac{\sqrt{n}}{\sigma} < \frac{z_{\alpha/2} + z_{1-\beta}}{\mu_0 - \mu} \\ &\Rightarrow n \geq \left[\frac{\sigma(z_{\alpha/2} + z_{1-\beta})}{\mu_0 - \mu}\right]^2. \end{aligned}$$

\exists_5 Note the result holds for $\mu_0 < \mu$ or $\mu_0 > \mu$.

TESTING $H_0: \sigma^2 = \sigma_0^2$, μ UNKNOWN

(CSS932)

\exists_1 Recall

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}.$$

\exists_2 We can use the pivotal quantity

$$U = \frac{(n-1)s^2}{\sigma_0^2}.$$

If $H_0: \sigma = \sigma_0$ is true, then it follows that

$$U \sim \chi^2_{n-1}.$$

\exists_3 Note large and small values of U provide evidence against H_0 .

\exists_4 The p-value is then approximately

$$P = \begin{cases} 2P(U \leq u), & P(U \leq u) < 0.5 \\ 2P(U \geq u), & P(U \geq u) \leq 0.5, \end{cases}$$

where

$$U \sim \chi^2_{n-1}, \quad u = \frac{(n-1)s^2}{\sigma_0^2}.$$

★ END of content for MT2.

LIKELIHOOD RATIO TEST STATISTIC

B₁: Recall the likelihood ratio statistic

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

where $\hat{\theta}$ = MLE of θ .

B₂: For large n , we showed $\Lambda \sim \chi^2_{n-1}$.

TESTING $H_0: \theta = \theta_0$ USING Λ (C5S958)

To test $H_0: \theta = \theta_0$ using Λ :

① Find $L(\theta)$ & $R(\theta) = \frac{L(\theta)}{L(\hat{\theta})}$ from the sample.

② Compute

$$\lambda(\theta_0) = -2\log(R(\theta_0)) = -2\log\left(\frac{L(\theta_0)}{L(\hat{\theta})}\right),$$

i.e. the observed value of Λ .

③ Find the p-value

$$p = P(W \geq \lambda(\theta_0)) = 2[1 - P(Z \leq \sqrt{-2\log(R(\theta_0))})]$$

where $Z \sim N(0,1)$ & $W \sim \chi^2_1$.

Chapter 6: Gaussian Response Models

Idea: we want to study the relationships between variates (in bivariate data).

One possible method: sample correlation.

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}, \quad S_{\alpha\beta} = \sum_{i=1}^n (\alpha_i - \bar{\alpha})(\beta_i - \bar{\beta}), \quad \alpha, \beta \in \{x, y\}$$

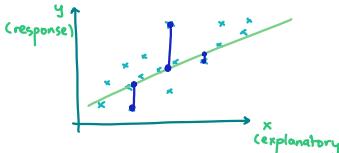
LEAST SQUARES ESTIMATES

How can we fit a straight line to bivariate data?



RESIDUALS (C6S978)

The "residuals" are the distances between the fitted line and the data.



LEAST SQUARES ESTIMATE (C6S979)

Usually, we find the fitted line $y = \alpha + \beta x$ that minimizes the sum of squares of the residuals.

Estimates of α & β , ie " $\hat{\alpha}$ " & " $\hat{\beta}$ ", are called the "least squares estimate".

We want to find $\hat{\alpha}$ & $\hat{\beta}$ that minimizes

$$g(\alpha, \beta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

which are given by

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}, \quad \hat{\beta} = \frac{S_{xy}}{S_{xx}}$$

Proof Sketch. We can get $\hat{\alpha}$ & $\hat{\beta}$ by solving the simultaneous eqns

$$\begin{cases} \frac{\partial g}{\partial \alpha} = \frac{\partial}{\partial \alpha} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i)(-1) = 0 \\ \frac{\partial g}{\partial \beta} = \frac{\partial}{\partial \beta} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i)(-x_i) = 0 \end{cases}$$

which resolves to $\alpha = \bar{y} - \hat{\beta} \bar{x}$ & $\sum_{i=1}^n (y_i - \alpha - \beta x_i)x_i = 0$. Set $\alpha = \hat{\alpha}$, $\beta = \hat{\beta}$. Algebra gives us the desired result (see slides for full details.)

In particular,

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} \cdot \sqrt{\frac{S_{xy}}{S_{xx}}} = r \sqrt{\frac{S_{xy}}{S_{xx}}}$$

Thus, the sign of $\hat{\beta}$ = sign of r , and $\hat{\beta}$ and r are linearly related.

In R, we can do this using

> lm(y ~ x)

or we can create a "model object"

> mod <- lm(y ~ x)
> summary(mod)

SIMPLE LINEAR REGRESSION (C6S998)

Idea:

① σ — the sd of the x 's (unknown)

② $\mu(x) = \alpha + \beta x$ — the mean y -value in the study population with x -value x .

* x = explanatory variate
 y = response variate

2 In particular,

- $\alpha = \mu(0)$ = mean y -value amongst data s.t. $x=0$
(not really useful)

- β represents the "increase" in the mean y -value in the study population for a one 'unit' increase in the x -value.
(this is the same regardless of x)

3 We assume $y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma)$ for $i=1, \dots, n$, and so

σ represents the variability in the response variate y in the study population for each value of the explanatory variate x .

LIKELIHOOD FUNCTION FOR α & β (C6S1004)

1 Our likelihood function for α & β is

$$L(\alpha, \beta) = \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2\right)$$

(since we assume $y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma)$)

2 So, to maximize $L(\alpha, \beta)$, we minimize

$$\sum_{i=1}^n (y_i - \alpha - \beta x_i)$$

but this is just the least squares problem!

3 Therefore, the MLEs of α & β are

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}, \quad \hat{\beta} = \frac{S_{xy}}{S_{xx}}$$

REGRESSION PARAMETERS: $\hat{\alpha}$, $\hat{\beta}$

We call the values of $\hat{\alpha}$ and $\hat{\beta}$ above the "regression parameters".

DISTRIBUTION OF $\hat{\beta}$ (C6S1011)

$\hat{\beta}$ A corresponding estimator for β is

$$\hat{\beta} = \frac{1}{S_{xx}} \sum_{i=1}^n (x_i - \bar{x}) y_i.$$

$\hat{\beta}$ If each $y_i \sim G(\alpha + \beta x_i, \sigma)$ independently, then it follows that

$$\hat{\beta} \sim G\left(\beta, \frac{\sigma}{\sqrt{S_{xx}}}\right).$$

Proof sketch. $\hat{\beta}$ is a linear combination of Gaussian rv.

$\Rightarrow \hat{\beta}$ is also a Gaussian rv. \square

PIVOTAL QUANTITY OF β IF σ IS KNOWN (C6S1013)

$\hat{\beta}$ Thus, if σ is known, a pivotal quantity for $\hat{\beta}$ is

$$\frac{\hat{\beta} - \beta}{\sigma / \sqrt{S_{xx}}} \sim G(0, 1).$$

ESTIMATING σ^2 IN SIMPLE LINEAR REGRESSION & THE MEAN SQUARED ERROR: S_e^2 (C6S1014)

$\hat{\beta}$ If σ^2 is unknown, we estimate it via

$$S_e^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i - \hat{\beta} \bar{x})^2 = \frac{1}{n-2} (S_{yy} - \hat{\beta} S_{xy}).$$

the "mean squared error" the "sum of squared errors"

PIVOTAL QUANTITY OF β IF σ IS UNKNOWN (C6S1020)

$\hat{\beta}$ We can show that

$$\frac{(n-2) S_e^2}{\sigma^2} \sim \chi_{n-2}^2. \quad * \text{ note } n-2.$$

$\hat{\beta}$ Hence, recall that

$$T = \frac{z}{\sqrt{U/k}} \sim t_k, \quad U \sim \chi_k^2, \quad z \sim G(0, 1)$$

and so

$$\frac{\hat{\beta} - \beta}{S_e / \sqrt{S_{xx}}} \sim t_{n-2}.$$

CI FOR β (C6S1021)

$\hat{\beta}$ Thus, if $P(-a \leq T \leq a) = p$ for $T \sim t_{n-2}$, a 100p% confidence interval for β is

$$P\left(\hat{\beta} - a \frac{S_e}{\sqrt{S_{xx}}} \leq \beta \leq \hat{\beta} + a \frac{S_e}{\sqrt{S_{xx}}}\right).$$

(or $P(T \leq a) = \frac{1-p}{2}$).

$\hat{\beta}$ For testing $H_0: \beta = \beta_0$, the p-value is

$$p = 2 \left[1 - P\left(T \leq \frac{|\hat{\beta} - \beta_0|}{S_e / \sqrt{S_{xx}}}\right) \right]$$

where $T \sim t_{n-2}$.

HYPOTHESIS OF NO (LINEAR) RELATIONSHIP (C6S1023)

$\hat{\beta}$ A discrepancy measure for testing $H_0: \beta = \beta_0$ is

$$\frac{|\hat{\beta} - \beta_0|}{S_e / \sqrt{S_{xx}}} \sim t_{n-2}$$

which is larger if the data are surprising if H_0 is true.

$\hat{\beta}$ Since $\mu(x) = \alpha + \beta x$, a test of $H_0: \beta = 0$ is a test of the hypothesis that $\mu(x)$ does not depend on x .

CI FOR $\mu(x) = \alpha + \beta x$ (C6S1034)

$\hat{\beta}$ A point estimate for $\mu(x)$ is

$$\hat{\mu}(x) = \hat{\alpha} + \hat{\beta} x = \bar{y} + \hat{\beta}(x - \bar{x})$$

and so the corresponding estimator is

$$\tilde{\mu}(x) = \bar{y} + \hat{\beta}(x - \bar{x}).$$

$\hat{\beta}$ We can show that

$$\tilde{\mu}(x) = \sum_{i=1}^n \left(\frac{1}{n} + (x - \bar{x}) \frac{(x_i - \bar{x})}{S_{xx}} \right) y_i$$

where $y_i \sim G(\alpha + \beta x_i, \sigma)$ for each i .

$\hat{\beta}$ This has distribution

$$\tilde{\mu}(x) \sim G\left(\mu(x), \sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}\right)$$

where $\tilde{\mu}(x) = \hat{\alpha} + \hat{\beta} x$ & $\mu(x) = \alpha + \beta x$.

$\hat{\beta}$ Equivalently

$$\frac{\tilde{\mu}(x) - \mu(x)}{\sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim G(0, 1).$$

Since σ is unknown, we use the pivotal quantity

$$\frac{\tilde{\mu}(x) - \mu(x)}{S_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim t_{n-2}$$

$\hat{\beta}$ A 100p% CI for $\mu(x) = \alpha + \beta x$ is

$$\hat{\alpha} + \hat{\beta} x \pm \sigma S_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

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

CI FOR α (C6S1038)

$\hat{\beta}$ Since $\mu(0) = \alpha + \beta(0)$, a 100p% CI for α is given by

$$\hat{\alpha} \pm a S_e \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

CI FOR AN INDIVIDUAL RESPONSE Y AT X

Q1 Question:

"What is the CI for y such that $x?$ "

PREDICTION INTERVAL [FOR A FUTURE RESPONSE Y] (C6S1049)

Q1 A "100p% prediction interval" for a future response

y is

$$\hat{y} + \hat{\beta}x \pm s_e \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}$$

* note the $+1$
(which is absent
in CIs).

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

How do we get this?

Let y = potential observation for given value of x .

We then have

$$y = \mu(x) + R, \quad R \sim N(0, \sigma)$$

independent of Y_1, \dots, Y_n . We established

$$y \sim N(\mu(x) + R, \sigma^2) \quad & \hat{\mu}(x) \sim N(\mu(x), \sigma^2 \sqrt{1 + \frac{(x - \bar{x})^2}{s_{xx}}}).$$

Then

$$y - \hat{\mu}(x) = y - \mu(x) + \mu(x) - \hat{\mu}(x) \\ = R + [\mu(x) - \hat{\mu}(x)].$$

Note R is independent of $\hat{\mu}(x)$ since it is not connected to the existing sample.

Thus the equation is a linear combination of ind, normally dist rv, and so is also normally dist.

$$\begin{aligned} E(y - \hat{\mu}(x)) &= E(R + [\mu(x) - \hat{\mu}(x)]) \\ &= E(R) + E(\mu(x)) - E(\hat{\mu}(x)) \\ &= 0 + \mu(x) - \mu(x) = 0. \end{aligned}$$

$$\begin{aligned} \text{Var}(y - \hat{\mu}(x)) &= \text{Var}(y) + \text{Var}(\hat{\mu}(x)) \\ &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}} \right] \\ &= \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}} \right]. \end{aligned}$$

Thus

$$y - \hat{\mu}(x) \sim N(0, \sigma^2 \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}).$$

Thus

$$\frac{y - \hat{\mu}(x)}{\sigma \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}} \sim N(0, 1)$$

and since σ is unknown we use

$$\frac{y - \hat{\mu}(x)}{s_e \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}} \sim t_{n-2}.$$

The corresponding CI is the one above.

Q2 In R, we can use

> predict(data, data.frame(x=75), interval='prediction')

GAUSSIAN RESPONSE MODELS

B1 The general form of a Gaussian response model is

$$Y_i \sim \mathcal{N}(\mu(x_i), \sigma), \quad i=1, \dots, n$$

independently where x_i are assumed to be known constants (possibly vectors).

B2 We can also write this as

$$Y_i = \mu(x_i) + R_i,$$

where $R_i \sim \mathcal{N}(0, \sigma)$ independently.

B3 In particular:

- ① $\mu(x_i)$ is a "deterministic" component;
- ② R_i is a "random" component.

LINEAR REGRESSION MODELS (C6S1056)

B1 In "linear regression models", the deterministic component takes the form

$$E(Y_i) = \mu(x_i) = \beta_0 + \sum_{j=1}^k \beta_j x_{ij}$$

so that $E(Y_i)$ is a linear function of a vector of explanatory variables for unit i ($x_i = (x_{i1}, \dots, x_{ik})$) & unknown parameters β_0, \dots, β_k .

B2 In particular,

- ① The β_j are called the "regression coefficients"; &
- ② The x_{ij} are called the "covariates".

MULTIPLE LINEAR REGRESSION (C6S1058)

B1 If we wish to fit the model

$$E(Y_i) = \mu(x_i) = \beta_0 + \sum_{j=1}^k \beta_j x_{ij}$$

then we seek parameters β_0, \dots, β_k that minimize

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{ij})^2$$

B2 In R, we use

```
> mod <- lm(y ~ x1 + x2)
> summary(mod)
```

Coefficients:	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.01375	5.01527	-0.202	0.84133
x1	0.73142	0.07664	9.544	3.83e-10 ***
x2	0.28225	0.09850	2.866	0.00797 **

Then $\hat{\beta}_0 = -0.01375$
 $\hat{\beta}_1 = 0.73142$
 $\hat{\beta}_2 = 0.28225$

INTERPRETING $\hat{\beta}_j$ (C6S1060)

B1 $\hat{\beta}_j$ can be interpreted as the amount of increase in response y when x_j increases by one unit when the other predictors $x_1, x_2, \dots, x_{j-1}, x_{j+1}, \dots, x_k$ are held fixed.

HYPOTHESIS TEST OF $H_0: \beta_j = 0$ (C6S1061)

B1 To test $H_0: \beta_j = 0$, $H_1: \beta_j \neq 0$, we use the test statistic

$$t_j = \frac{\hat{\beta}_j}{\text{std. error}} = \frac{\text{estimate}}{\text{std. error}}$$

B2 If H_0 is true, then

$$t_j \sim t_{n-k-1}$$

where $k = \# \text{ of parameters}$.

SUM OF SQUARED ERRORS / SSE (C6S1063)

B1 The "sum of squared errors" is

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{where } \hat{y}_i = \hat{\beta}_0 + \sum_{j=1}^k \hat{\beta}_j x_{ij}.$$

B2 The smaller this is, the less 'error' in our model fit.

R² STATISTIC (C6S1064)

B1 The "R² statistic" is

$$R^2 = 1 - \frac{\text{SSE}}{S_{yy}}.$$

B2 In particular,

$$R^2 = \frac{\text{variation explained by regression}}{\text{total variation}}$$

ADJUSTED R²-STATISTIC (C6S1065)

B1 The "adjusted R²" is

$$\text{adj. } R^2 = 1 - \frac{\text{SSE} / (n-k-1)}{S_{yy} / (n-1)}$$

where $k = \# \text{ of explanatory variables}$.

B2 This tries to "compensate" that adding more & more variables will (potentially artificially) increase R².

ASSUMPTIONS (C6S1067)

B1 Assumptions we make for Gaussian linear response models:

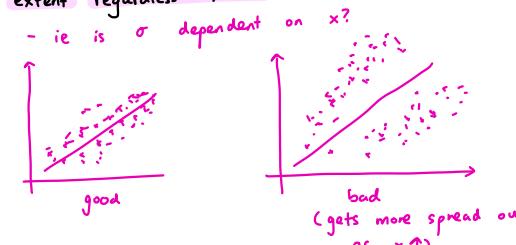
- ① Y_i (given covariates x_i) has a Gaussian distribution with std dev σ which does not depend on the covariates; &
- ② $E(Y_i) = \mu(x_i)$ is a linear combination of known covariates $x_i = (x_{i1}, \dots, x_{in})$ and the unknown regression coefficients β_0, \dots, β_k .

B2 We must check these are suitable!

GRAPHICAL METHOD TO CHECK MODEL (C6S1069)

B1 We can use graphical methods to do this: in particular,

- ① Make a scatterplot of y against x :
- ② Do the points seem to fit reasonably along a straight line?
 - ie is it linear
- ③ Are the points generally "spread out" to the same extent regardless of x ?
 - ie is σ dependent on x ?



bad
(gets more spread out as $x \uparrow$)

USING RESIDUALS TO CHECK MODEL (C6S1070)

\mathbb{B}_1 We let our "fitted response" to be

$$\hat{\mu}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}$$

- for simple LR, $\hat{\mu}_i = \hat{\beta} + \hat{\beta}' x$.

\mathbb{B}_2 The "residuals" are

$$\hat{r}_i = y_i - \hat{\mu}_i.$$

this represents what has been "left over" after the model has been fitted to the data.

\mathbb{B}_3 We assume $y_i \sim \mathcal{N}(\mu(x_i), \sigma)$, and in particular,

$$Y_i = \mu_i + R_i.$$

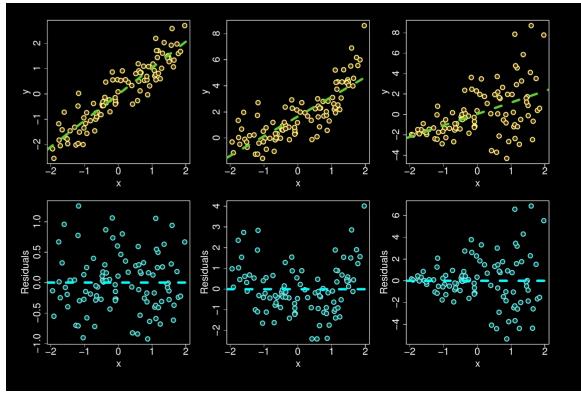
where $R_i \sim \mathcal{N}(0, \sigma^2)$.

RESIDUAL PLOTS (C6S1073)

\mathbb{B}_1 A "residual plot" is a plot of the points (x_i, \hat{r}_i) .

\mathbb{B}_2 If the model assumptions hold, the points

- ① should lie more or less within a horizontal band around the line $\hat{r}_i = 0$
- ② with no obvious pattern.



STANDARDIZED RESIDUALS: \hat{r}_i^* (C6S1077)

\mathbb{B}_1 We define the "standardized residuals" to be

$$\hat{r}_i^* = \frac{\hat{r}_i}{s_e} = \frac{y_i - \hat{\mu}_i}{s_e}.$$

\mathbb{B}_2 If we plot (x_i, \hat{r}_i^*) instead of (x_i, \hat{r}_i) :

- ① The plot looks the same, but be "scaled";
- ② The \hat{r}_i^* values lie in the range $(-3, 3)$ since $\hat{r}_i^* \sim \mathcal{N}(0, 1)$.

RESIDUAL PLOT TYPE 2 (C6S1080)

\mathbb{B}_1 If we have a more general linear model

$$E(Y_i) = \mu_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

we can plot $(\hat{\mu}_i, \hat{r}_i^*)$, where

$$\hat{\mu}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}.$$

\mathbb{B}_2 We can use this to check the assumption about the form of $\mu(x_i)$;

we check if the points appear randomly scattered around a horizontal line at 0.

QQ-PLOT OF RESIDUALS (C6S1082)

\mathbb{B}_1 Since $\hat{r}_i^* \approx \mathcal{N}(0, 1)$, a QQ-plot should give approximately a straight line if the model assumptions hold.

MULTICOLLINEARITY (C6S1094)

\mathbb{B}_1 "Multicollinearity" describes a situation when two (or more) of our explanatory variates are highly correlated.

\mathbb{B}_2 This can occur when we have collected data on several variates on the same subject.

\mathbb{B}_3 This can make us deduce incorrect conclusions.

PREDICTING BEYOND THE RANGE OF COVARIATES (C6S1097)

\mathbb{B}_1 We may be tempted to predict an outcome for a covariate value outside the range of those in our dataset.

\mathbb{B}_2 However,

- ① Our model assumptions may no longer hold, and we have no way of checking them.
- ② Our predictions might also not make sense.

BINARY OUTCOMES (C6S1099)

ODDS RATIO [OF AN EVENT]: odds(E)
(C6S1105)

💡 The "odds" of an event E is

$$\text{odds}(E) = \frac{P(E)}{1-P(E)}$$

💡 If $\text{odds}(E) = \frac{a}{b}$ ("the odds of E are a to b"), then

$$P(E) = \frac{a}{a+b}.$$

GENERALIZED LINEAR MODELS / GLMs (C6S1108)

💡 GLMs have the following properties:

- ① A probability distribution for the outcome variable;
- ② A linear model $\eta = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$; &
- ③ A "link function" relating the linear model to the parameters of the outcome distribution.

LOGISTIC REGRESSION (C6S1109)

💡 "Logistic regression" is a GLM for binary outcome data.

💡 Assumptions:

- ① Outcome can be modelled by a binomial rv;
- ② We want to model p , the probability of success.

💡 A common link function is "logit":

$$g(p) = \text{logit}(p) = \log\left(\frac{p}{1-p}\right), \quad p \in [0,1]$$

which maps from $[0,1] \rightarrow [-\infty, +\infty]$.

💡 This is the log odds of success.

USING logit IN LOGISTIC REGRESSION (C6S1110)

💡 Assume we had a single explanatory variate x_i , and let p_i be the probability unit i experiences the outcome.

💡 We can use logit to relate this probability to a linear model of our data:

$$g(p_i) = \text{logit}(p_i) = \log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_i$$

💡 This can be rewritten as

$$p_i = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x_i))}$$

💡 In R, we use

```
< mod1 <- glm(offer ~ grade, family = 'binomial')  
< summary(mod1)
```

to get the estimates for the linear model.

💡 We use

```
< predict(mod1, newdata = data.frame(grade=80), type='response')
```

to get the odds directly.

ODDS RATIO & LOG ODDS RATIO

(C6S1121)

💡 Let O_1 be the odds of E_1 , & O_2 be the odds of E_2 .

💡 The "odds ratio" of E_1 relative to E_2 is

$$\text{odds ratio} = \frac{O_1}{O_2}.$$

💡 The "log odds ratio" of E_1 relative to E_2 is

$$\text{log odds ratio} = \log\left(\frac{O_1}{O_2}\right).$$

💡 If $\gamma_i = \beta_0 + \beta_1 x_i$, then $\hat{\beta}_1$ is our estimate of the log odds ratio of a one unit increase in x .

ASSUMPTIONS FOR LOGISTIC REGRESSION (C6S1128)

💡 Assumptions for logistic regression:

- ① Events are independent; &
- ② A linear relationship exists between predictors and the log odds.

💡 One option: split the data into tertiles / quantiles etc.

COMPARING MEANS OF TWO POPULATIONS

TWO-SAMPLE GAUSSIAN PROBLEM (C6S1149)

Θ_1 A "two-sample Gaussian problem" involves

$$\begin{aligned} Y_{1i} &\sim \mathcal{N}(\mu_1, \sigma^2), \quad i=1, \dots, n_1 \text{ independently}; \\ Y_{2i} &\sim \mathcal{N}(\mu_2, \sigma^2), \quad i=1, \dots, n_2 \text{ independently}. \end{aligned}$$

Θ_2 This is a special case of the Gaussian response model.

HYPOTHESIS TEST THAT TWO MEANS ARE THE SAME (C6S1140)

Θ_1 To check if $\mu_1 = \mu_2$, we use

$$H_0: \mu_1 = \mu_2$$

or equivalently

$$H_0: \mu_1 - \mu_2 = 0.$$

POINT ESTIMATORS FOR μ_1 & μ_2 (C6S1141)

Θ_1 First, note the ML estimators for μ_1 & μ_2 are

$$\begin{aligned} \hat{\mu}_1 &= \bar{Y}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} Y_{1i} \\ \hat{\mu}_2 &= \bar{Y}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} Y_{2i} \end{aligned}$$

and so a point estimator for $\mu_1 - \mu_2$ is

$$\hat{\mu}_1 - \hat{\mu}_2 = \bar{Y}_1 - \bar{Y}_2.$$

Θ_2 Since $\hat{\mu}_1 = \bar{Y}_1 \sim \mathcal{N}(\mu_1, \frac{\sigma^2}{n_1})$, $\hat{\mu}_2 = \bar{Y}_2 \sim \mathcal{N}(\mu_2, \frac{\sigma^2}{n_2})$ independently, it follows that

$$\hat{\mu}_1 - \hat{\mu}_2 = \bar{Y}_1 - \bar{Y}_2 \sim \mathcal{N}(\mu_1 - \mu_2, \sigma^2 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}})$$

and so

$$\frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim \mathcal{N}(0, 1).$$

POINT ESTIMATOR FOR σ^2 ; THE POOLED ESTIMATOR FOR VARIANCE: S_p^2 (C6S1143)

Θ_1 First, define

$$\begin{aligned} S_1^2 &= \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_1)^2 \\ S_2^2 &= \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (Y_{2i} - \bar{Y}_2)^2 \end{aligned}$$

which are the point estimators of σ^2 based on only the Y_{1i} , & only the Y_{2i} , respectively.

Θ_2 Our point estimator of σ^2 is then

$$\begin{aligned} S_p^2 &= \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \\ &= \frac{1}{n_1 + n_2 - 2} \left[\sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_1)^2 + \sum_{i=1}^{n_2} (Y_{2i} - \bar{Y}_2)^2 \right] \end{aligned}$$

which is the "pooled estimator of variance".

Θ_3 Note $E(S_p^2) = \sigma^2$, so the estimator is unbiased. (It is NOT the ML estimator.)

Θ_4 Our pivotal quantity for S_p^2 is

$$\frac{(n_1 + n_2 - 2)S_p^2}{\sigma^2} \sim \chi^2_{n_1 + n_2 - 2}.$$

PIVOTAL QUANTITY FOR $\mu_1 - \mu_2$ (C6S1146)

Θ_1 From the previous results, thus

$$\frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

(since $Z \sim \mathcal{N}(0, 1)$, $U \sim \chi^2_{n-1}$ independently $\Rightarrow T = \frac{Z}{\sqrt{U/n}} \sim t_{n-1}$)

CI FOR $\mu_1 - \mu_2$ (C6S1147)

Θ_1 A 100 $\gamma\%$ CI for $\mu_1 - \mu_2$ is thus

$$CI = \bar{Y}_1 - \bar{Y}_2 \pm a s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

where $P(T \leq a) = \frac{1+\gamma}{2}$ & $T \sim t_{n_1 + n_2 - 2}$.

TEST STATISTIC FOR $H_0: \mu_1 - \mu_2 = 0$ (C6S1149)

Θ_1 The test statistic for $H_0: \mu_1 - \mu_2 = 0$ is

$$D = \frac{|\bar{Y}_1 - \bar{Y}_2 - 0|}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

with observed value

$$d = \frac{|\bar{Y}_1 - \bar{Y}_2 - 0|}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}.$$

Θ_2 In particular, large values of d would be surprising if H_0 is true.

p-VALUE FOR $H_0: \mu_1 - \mu_2 = 0$ (C6S1150)

Θ_1 The p-value is thus

$$p = 2[1 - P(T \leq d)]$$

where $T \sim t_{n_1 + n_2 - 2}$.

COMPARISON OF 2 MEANS WITH UNEQUAL VARIANCES (C6S1156)

APPROXIMATE PIVOTAL QUANTITY FOR $\mu_1 - \mu_2$ (C6S1157)

Θ_1 Suppose instead that

$$\begin{aligned} Y_{1i} &\sim \mathcal{N}(\mu_1, \sigma_1^2), \quad i=1, \dots, n_1 \text{ independently} \\ Y_{2i} &\sim \mathcal{N}(\mu_2, \sigma_2^2), \quad i=1, \dots, n_2 \text{ independently} \end{aligned}$$

where we don't assume $\sigma_1 = \sigma_2$.

Θ_2 If n_1, n_2 are large (≥ 30), we can use the approximate pivotal quantity

$$\frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \approx \mathcal{N}(0, 1).$$

Θ_3 Thus, an approximate CI for $\mu_1 - \mu_2$ is

$$CI = \mu_1 - \mu_2 \pm a \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

where $P(Z \leq a) = \frac{1+\gamma}{2}$, $Z \sim \mathcal{N}(0, 1)$.

PAIRED DATA

Q₁ "Paired data" considers scenarios where the y_{1i} 's are related to the y_{2i} 's.
eg y_{1i} = movies, y_{2i} = their sequels

Q₂ Suppose once again that

$$y_{1i} \sim G(\mu_1, \sigma_1^2), \quad i=1, \dots, n, \text{ independently}$$
$$y_{2i} \sim G(\mu_2, \sigma_2^2), \quad i=1, \dots, n, \text{ independently}$$

but the set of y_{1i}, y_{2i} is not independent with each other.

Q₃ Then

$$\text{Var}(\bar{Y}_1 - \bar{Y}_2) = \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n} - 2\text{Cov}(Y_{1i}, Y_{2i}).$$

which is smaller than for an unpaired experiment.

Q₄ To make inferences about $\mu_1 - \mu_2$, we analyze the within-pair differences

$$Y_i = Y_{1i} - Y_{2i} \quad \forall i=1, \dots, n$$

by assuming

$$Y_i = Y_{1i} - Y_{2i} \sim G(\mu_1 - \mu_2, \sigma) \quad \forall i=1, \dots, n$$

independently.

Q₅ To test $H_0: \mu_1 - \mu_2 = 0$, we use the test statistic

$$D = \frac{|\bar{Y} - 0|}{s/\sqrt{n}} \sim t_{n-1} \quad (\text{if } H_0 \text{ is true}).$$

and our p-value is

$$p = 2(1 - P(T \leq d))$$

where $T \sim t_{n-1}$.

Chapter 7: Multinomial Models & Goodness of Fit Tests

MULTINOMIAL LIKELIHOOD FUNCTION (C7S1184)

Θ_1 : We consider

$$(Y_1, \dots, Y_K) \sim \text{Multinomial}(n, \theta_1, \dots, \theta_k).$$

Θ_2 : The likelihood function based on y_1, \dots, y_k is

$$L(\theta_1, \dots, \theta_k) = P(Y_1 = \theta_1, \dots, Y_k = \theta_k; \theta_1, \dots, \theta_k) \\ = \frac{n!}{y_1! \dots y_k!} \theta_1^{y_1} \dots \theta_k^{y_k}.$$

Ditching the constant gets us that

$$L(\theta_1, \dots, \theta_k) = \theta_1^{y_1} \dots \theta_k^{y_k}.$$

TESTING $H_0: \theta_1 = \dots = \theta_k$ (C7S1184)

Θ_1 : Suppose we wish to test

$$H_0: \theta_1 = \dots = \theta_k \Leftrightarrow H_0: \theta = \theta_0 = (\frac{1}{k}, \dots, \frac{1}{k}).$$

Θ_2 : We note that the MLE of θ_j is

$$\hat{\theta}_j = \frac{y_j}{n}.$$

LIKELIHOOD RATIO TEST STATISTIC FOR $H_0: \theta = \theta_0$ (C7S1187)

Θ_1 : The likelihood ratio test statistic for H_0 is then

$$\Lambda(\theta_0) = -2 \log\left(\frac{L(\theta_0)}{L(\tilde{\theta})}\right)$$

where

$$\begin{aligned} \tilde{\theta} &= \left(\frac{y_1}{n}, \dots, \frac{y_k}{n} \right) \\ \theta_0 &= \left(\frac{1}{k}, \dots, \frac{1}{k} \right) \\ L(\theta_1, \dots, \theta_k) &= \theta_1^{y_1} \dots \theta_k^{y_k} = \prod_{j=1}^k \theta_j^{y_j}. \end{aligned}$$

Θ_2 : This evaluates to

$$\Lambda(\theta_0) = 2 \sum_{j=1}^k y_j \log\left(\frac{y_j}{E_j}\right), \quad E_j = \frac{n}{k}$$

where E_j is the "expected frequency" of y_j .

P-VALUE FOR $H_0: \theta = \theta_0$ (C7S1193)

Θ_1 : The observed value of $\Lambda(\theta_0)$ is

$$\lambda(\theta_0) = 2 \sum_{j=1}^k y_j \log\left(\frac{y_j}{e_j}\right)$$

$$\text{where } e_j = \frac{n}{k}.$$

Θ_2 : For a sufficiently large sample, we can use the p-value

$$\text{p-value} = P(W \geq \lambda(\theta_0)), \quad W \sim \chi^2_{k-1-p}$$

where $p = \# \text{ of parameters estimated in forming } H_0$.

* degrees of freedom = $k-1-p$ because these are the parameters that are "free to move".

(When we estimate parameters, they are "locked").

Θ_3 : Guideline: "sufficiently large" roughly implies $e_j \geq 5$ for all j .

- if this is not satisfied, we can "collapse" two or more adjacent categories with the smallest expected probabilities.

PEARSON'S χ^2 GOODNESS OF FIT STATISTIC (C7S1201)

\exists_1 The "Pearson's χ^2 test statistic" is

$$D = \sum_{j=1}^k \frac{(y_j - E_j)^2}{E_j} \sim \chi^2_{k-1-p}$$

\exists_2 For large n , D & Λ are asymptotically equivalent & have the same asymptotic χ^2 distribution.

\exists_3 Note we need to account for any parameters we estimated.

e.g. If we assume $\text{Poi}(\theta)$, we need to account for

- sample size, n
 - MLE of θ , $\hat{\theta}$
- } so degrees of freedom.
 $v = k-1-1 = k-2$.

\exists_4 Our p-value is thus

$$p = P(W \geq d), \quad d = \sum_{j=1}^k \frac{(y_j - E_j)^2}{E_j}, \quad W \sim \chi^2_{k-1-p}.$$

TWO-WAY TABLES & INDEPENDENCE TESTS (C7S1218)

\exists_1 "2-way tables" have the following form:

	B	\bar{B}	total
A	y_{11}	y_{12}	$r_1 = y_{11} + y_{12}$
\bar{A}	y_{21}	y_{22}	$n - r_1$
total	$c_1 = y_{11} + y_{21}$	$n - c_1$	n

\exists_2 We are concerned whether there is a relationship between A & B, and in particular, whether they are independent.

MODEL FOR TEST OF INDEPENDENCE (C7S1224)

\exists_1 We define the random variables

- $Y_{11} = \# \text{ of } A \cap B \text{ outcomes}$
- $Y_{12} = \# \text{ of } A \cap \bar{B} \text{ outcomes}$
- $Y_{21} = \# \text{ of } \bar{A} \cap B \text{ outcomes}$
- $Y_{22} = \# \text{ of } \bar{A} \cap \bar{B} \text{ outcomes.}$

\exists_2 Then our model is

$$(Y_{11}, Y_{12}, Y_{21}, Y_{22}) \sim \text{Multinomial}(n, \theta_{11}, \theta_{12}, \theta_{21}, \theta_{22})$$

where

$$\begin{aligned} \theta_{11} &= P(A \cap B) & \theta_{12} &= P(A \cap \bar{B}) \\ \theta_{21} &= P(\bar{A} \cap B) & \theta_{22} &= P(\bar{A} \cap \bar{B}) \end{aligned}$$

HYPOTHESIS OF INDEPENDENCE (C7S1227)

\exists_1 To test whether A & B are independent, we use the null hypothesis

$$H_0: P(A \cap B) = P(A)P(B)$$

\exists_2 This is equivalent to

$$H_0: \theta_{11} = \alpha\beta, \quad \alpha = P(A), \quad \beta = P(B)$$

LIKELIHOOD FUNCTION FOR $H_0: \theta_{11} = \alpha\beta$ (C7S1229)

\exists_1 The likelihood function (ignoring constants) is

$$L(\theta_{11}, \theta_{12}, \theta_{21}, \theta_{22}) = \theta_{11}^{y_{11}} \theta_{12}^{y_{12}} \theta_{21}^{y_{21}} \theta_{22}^{y_{22}}$$

\exists_2 The ML estimates are

$$\hat{\theta}_{ij} = \frac{y_{ij}}{n}, \quad i=1,2, j=1,2$$

with corresponding estimators

$$\tilde{\theta}_{ij} = \frac{Y_{ij}}{n}, \quad i=1,2, j=1,2$$

PARAMETER ESTIMATION UNDER

$H_0: \theta_{11} = \alpha\beta$ (C7S1230)

\exists_1 If $H_0: \theta_{11} = \alpha\beta$ is true, then the likelihood function is

$$\begin{aligned} L(\theta) &= \theta_{11}^{y_{11}} \theta_{12}^{y_{12}} \theta_{21}^{y_{21}} \theta_{22}^{y_{22}} \\ \Rightarrow L(\alpha\beta) &= (\alpha\beta)^{y_{11}} [\alpha(1-\beta)]^{y_{12}} [(1-\alpha)\beta]^{y_{21}} [(1-\alpha)(1-\beta)]^{y_{22}} \\ &= \alpha^{y_{11}+y_{21}} (1-\alpha)^{y_{12}+y_{22}} \beta^{y_{11}+y_{21}} (1-\beta)^{y_{12}+y_{22}}. \end{aligned}$$

\exists_2 The ML estimates under H_0 are

$$\hat{\alpha} = \frac{y_{11} + y_{12}}{n}, \quad \hat{\beta} = \frac{y_{11} + y_{21}}{n}$$

with corresponding ML estimators

$$\tilde{\alpha} = \frac{Y_{11} + Y_{12}}{n}, \quad \tilde{\beta} = \frac{Y_{11} + Y_{21}}{n}$$

LIKELIHOOD RATIO TEST STATISTIC FOR

$$H_0: \theta_{11} = \alpha\beta \quad (\text{C7S1233})$$

\mathbb{U}_1 The likelihood ratio statistic is

$$\begin{aligned}\Lambda(\theta) &= -2\log\left(\frac{L(\tilde{\alpha}, \tilde{\beta})}{L(\hat{\theta}_{11}, \hat{\theta}_{12}, \hat{\theta}_{21}, \hat{\theta}_{22})}\right) \\ &= -2\log\left(\frac{\tilde{\alpha}^{Y_{11}} \tilde{\beta}^{Y_{12}} (1-\tilde{\alpha})^{Y_{21}} (1-\tilde{\beta})^{Y_{22}}}{\hat{\theta}_{11}^{Y_{11}} \hat{\theta}_{12}^{Y_{12}} \hat{\theta}_{21}^{Y_{21}} \hat{\theta}_{22}^{Y_{22}}}\right) \\ &\approx \chi^2.\end{aligned}$$

where

$$\text{degrees of freedom} = 4 - 1 - 2 = 1.$$

\downarrow
sample size, since we are estimating α & β

\mathbb{U}_2 This is also equivalent to

$$\Lambda(\theta) = 2[Y_{11}\log(\frac{Y_{11}}{E_{11}}) + Y_{12}\log(\frac{Y_{12}}{E_{12}}) + Y_{21}\log(\frac{Y_{21}}{E_{21}}) + Y_{22}\log(\frac{Y_{22}}{E_{22}})]$$

where

$$\begin{aligned}E_{11} &= n\hat{\alpha}\hat{\beta}, & E_{12} &= n\hat{\alpha}(1-\hat{\beta}), \\ E_{21} &= n(1-\hat{\alpha})\hat{\beta}, & E_{22} &= n(1-\hat{\alpha})(1-\hat{\beta}).\end{aligned}$$

* using $\Lambda(\theta) = 2 \sum_{j=1}^k Y_j \log(\frac{Y_j}{E_j})$.

\mathbb{U}_3 Our p-value is thus approximately

$$P = P(W \geq \lambda) = 2[1 - P(Z \leq \sqrt{\lambda})],$$

$$W \sim \chi^2_1, Z \sim N(0,1)$$

where λ is our observed value of Λ .

LARGER TWO-WAY TABLES (C7S1254)

\mathbb{U}_1 Let Y_{ij} = number of individuals in category A_i & category B_j in a sample size of n .
 \mathbb{U}_2 Let

$$\theta_{ij} = P(A_i; B_j).$$

Then our model is

$$(Y_{11}, Y_{12}, \dots, Y_{ab}) \sim \text{Multinomial}(n; \theta_{11}, \theta_{12}, \dots, \theta_{ab}).$$

HYPOTHESIS FOR INDEPENDENCE (C7S1256)

\mathbb{U}_1 Let

$$\alpha_i = P(A_i), \quad \beta_j = P(B_j).$$

\mathbb{U}_2 To test if A & B are independent variates, we test

$$H_0: \theta_{ij} = \alpha_i \beta_j \quad \forall i=1, \dots, a, j=1, \dots, b.$$

\mathbb{U}_3 Under H_0 , we can show the expected frequencies e_{ij} are

$$e_{ij} = \frac{r_i c_j}{n}, \quad i=1, \dots, a, \quad j=1, \dots, b.$$

where

$$\begin{aligned}r_i &= \# \text{ of outcomes under } A_i && \text{total of row } i \\ c_j &= \# \text{ of outcomes under } B_j && \text{total of column } j.\end{aligned}$$

LIKELIHOOD RATIO TEST STATISTIC (C7S1259)

\mathbb{U}_1 The likelihood ratio test statistic is then

$$\Lambda = 2 \sum_{i=1}^a \sum_{j=1}^b Y_{ij} \log\left(\frac{Y_{ij}}{E_{ij}}\right)$$

with observed value

$$\lambda = 2 \sum_{i=1}^a \sum_{j=1}^b y_{ij} \log\left(\frac{y_{ij}}{e_{ij}}\right).$$

\mathbb{U}_2 In particular, if H_0 is true and $e_{ij} \geq 5 \quad \forall i, j$, then

$$\Lambda \approx \chi^2_{(a-1)(b-1)}.$$

Proof: degrees of freedom,

$$\begin{aligned}v &= k-1-p \\ &= ab-1-(a-1)-(b-1) \\ &= (a-1)(b-1).\end{aligned}$$

* Since $\alpha_a = 1 - \alpha_1 - \dots - \alpha_{a-1}$, we don't estimate α_a !
Similar with β_b .

\mathbb{U}_3 Our p-value is thus

$$P = P(W \geq \lambda),$$

where $W \sim \chi^2_{(a-1)(b-1)}$.

Chapter 8: Causality or Relationship?

CAUSAL EFFECT (C8S1278)

- We say x has a "causal effect" on y if, when all other factors that affect y are held constant, a change in x induces a change in a property of the distribution of y .
- * this is impractical since we cannot hold all factors that affect y to be constant!
- B₂ we should design studies so that alternative explanations of what causes changes in the distribution of y can be ruled out, leaving x as the causal agent.

REASONS 2 VARIATES CAN BE RELATED

EXPLANATORY VARIATE IS THE DIRECT CAUSE OF THE RESPONSE VARIATE (C8S1281)

- Reason 1: A change in the explanatory variate directly causes a change in the response variate.
 - eg drinking tea & thirst
- B₂ Note that even in this case, we may not see a strong association.
 - eg playing the lottery & winning the lottery

RESPONSE VARIATE IS THE DIRECT CAUSE OF THE EXPLANATORY VARIATE (C8S1283)

- Reason 2: Similar to Reason 1, but now the causal relationship is "flipped"; the response variate directly causes the explanatory variate.

THE EXPLANATORY VARIATE IS A CONTRIBUTING, BUT NOT ONLY, CAUSE OF THE RESPONSE VARIATE (C8S1285)

- Reason 3: The explanatory variate is a contributing cause, but not the sole cause, of the response variate.
 - eg diet & type of cancer

- In particular, we may think we have found a sole cause, when in actuality we have found a necessary contributor to the outcome.

- eg HIV & AIDS
 - we need HIV to get AIDS (so it is a necessary contributor)
 - but HIV is not necessarily the sole cause of AIDS (there might be other factors).

BOTH VARIATES ARE CHANGING OVER TIME (C8S1287)

- Reason 4: Non-sensical associations can result from correlating two variates that are both changing over time.

- eg global avg temp. & # of pirates
 - they both decrease as time increases
 - but are not related in any way

THE ASSOCIATION MAY BE NOTHING MORE THAN COINCIDENCE (C8S1289)

- Reason 5: The association may be nothing more than coincidence.

BOTH VARIATES MAY RESULT FROM A COMMON CAUSE —

CONFOUNDING/LURKING VARIATES (C8S1292)

- Reason 6: An association between 2 variates may be observed because both variates are responding to changes in some unobserved variate or variates.
 - z
 - x
 - y
- B₂ These variates are called "confounding" variates.

SIMPSON'S PARADOX (C8S1294)

- "Simpson's paradox" describes the phenomenon where the association between 2 categorical variables is different than the association after controlling for one or more variables.

Age	Coke	Pepsi
< 30	93% (81/87)	87% (231/270)
≥ 30	73% (192/263)	69% (55/80)
Total	78% (273/350)	83% (289/350)

- for each individual row, coke is bigger
- but for the total, pepsi is bigger!

→ age is a confounding variate.

THE IMPORTANCE OF RANDOMIZATION

(C8S1301)

Randomization is important since it ensures the other variates will be approximately equally distributed across the categories.

ESTABLISHING CAUSATION IN OBSERVATIONAL STUDIES (C8S1304)

To establish causation in observational studies, we need at least the following 4 features:

- ① The association between the 2 variates must be observed in many studies of different types among different groups.
 - this reduces the chance an observed association is due to a defect in one type of study
 - or from a peculiarity in one group of subjects
- ② The association must continue to hold when the effects of plausible confounding variates are taken into account.
- ③ There must be a plausible scientific explanation for the direct influence of one variate on the other variate.
 - so a causal link does not depend on the observed association alone.
- ④ There must be a consistent response:
ie one variate increases (or decreases) as the other variate increases.

CAUSAL INFERENCE (C8S1308)

- Q₁ "Causal inference" is concerned with identifying causal, merely than just associative, relationships.
- Q₂ We are interested in causal effects, and quantifying the effect of a change in some variate on some outcome.

COUNTERFACTUALS (C8S1311)

- Q₁ We write $Y(1)$ for the result of a response variate given one value of the explanatory variate, and similarly $Y(0)$ for another value of the explanatory variate.

Q₂ We can then write the causal effect as

$$Y(1) - Y(0).$$

- Q₃ If we let $Y_i(1)$ & $Y_i(0)$ be the corresponding random variables for unit i , then it follows that the "average causal effect" is roughly

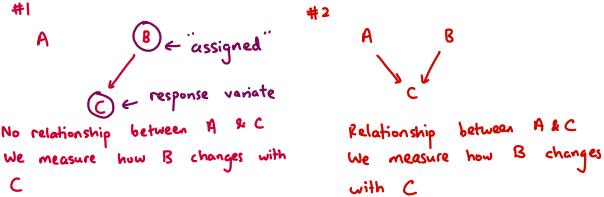
$$T = E(Y(1) - Y(0)).$$

- Q₄ Note for each unit we can only observe one of $Y(0)$ or $Y(1)$; the other is called a "counterfactual" or "potential outcome".

COVARIATE DEPENDENT TREATMENT

(C8S1315)

- Q₁ Consider the following 3 scenarios:



- Q₂ For Scenarios #1 & #2, we can perform a standard analysis:
- ① Let Y_0, Y_1 be the r.v. denoting the value of C given B or $\neg B$ respectively;
 - ② Suppose $Y_0 \sim G(\mu_0, \sigma_0)$ & $Y_1 \sim G(\mu_1, \sigma_1)$.
 - ③ Carry out a test of $H_0: \mu_0 = \mu_1$ using methods from earlier in this course.

- Q₃ However, this may not necessarily work in scenario #3.
→ since A influences B, B's effects are exaggerated.

INVERSE PROBABILITY WEIGHTING (C8S1326)

PROPENSITY SCORE: $\pi(x)$ (C8S1327)

Q₁ Let $\mu_0 = E(Y(0))$ & $\mu_1 = E(Y(1))$.

The "propensity score" is

$$\pi(x) = P(A=1 | X=x),$$

which is the probability A is true given covariate value x.

Q₂ We write $\hat{\pi}(x)$ for our estimate of $\pi(x)$ for a given x.

ESTIMATING μ_0 & μ_1 (C8S1333)

- Q₁ We can estimate $\hat{\mu}_0$ & $\hat{\mu}_1$ by

$$\begin{aligned}\hat{\mu}_1 &= \frac{1}{n} \sum_{i=1}^n \frac{y_i \cdot I_{\{A_i=1\}}}{\hat{\pi}(x_i)} = \frac{1}{n} \sum_{i=1}^n \frac{y_i q_i}{\hat{\pi}(x_i)} \\ \hat{\mu}_0 &= \frac{1}{n} \sum_{i=1}^n \frac{y_i \cdot I_{\{A_i=0\}}}{1-\hat{\pi}(x_i)} = \frac{1}{n} \sum_{i=1}^n \frac{y_i (1-q_i)}{1-\hat{\pi}(x_i)}\end{aligned}$$

where $I_{\{P\}}$ is the indicator function for the statement P.

Q₂ Intuition: we give more weight to more unusual data in $\hat{\mu}_1$, and the other way around in $\hat{\mu}_0$.

Q₃ Note the above $\hat{\mu}_0$ & $\hat{\mu}_1$ are unbiased.

ESTIMATING THE CAUSAL EFFECT (C8S1336)

- Q₁ We can estimate the causal effect

$$\hat{T} = \hat{\mu}_1 - \hat{\mu}_0.$$

Q₂ This is also unbiased.

→ since $\hat{\mu}_1, \hat{\mu}_0$ are unbiased.

ASSUMPTIONS FOR IPW (C8S1342)

- Q₁ Assumptions needed to perform IPW:

- Q₂ Consistency — a counterfactual outcome is equal to that which would have been observed had "B" been different.
- we set

$$Y_i = Y_i(0)(1-A_i) + Y_i(1)A_i.$$

- Q₃ Stable Unit Treatment Value Assumption (SUTVA) — there is no interference in the value of the B's.

- often reasonable, but not always!
- eg vaccines

- Q₄ No Unmeasured Confounders (NUC) —

- the assignment of the B's is independent of the potential outcomes, given covariates.

- we write

$$Y(0), Y(1) \perp\!\!\!\perp A | X$$

- Q₅ Positivity — there is a non-zero probability of assignment to a category of B for all subjects: ie

$$0 < P(A_i=1 | X=x) < 1 \quad \forall x$$

DIRECTED ACYCLIC GRAPHS (C8S1350)

We can use directed acyclic graphs to model hypothesized causal relationships.

eg



This shows we hypothesize screen time causes obesity.

MEDIATOR (C8S1351)

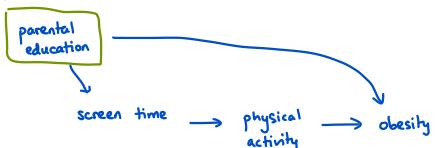
Consider



"physical activity" is a mediator.

CONFOUNDER (C8S1352)

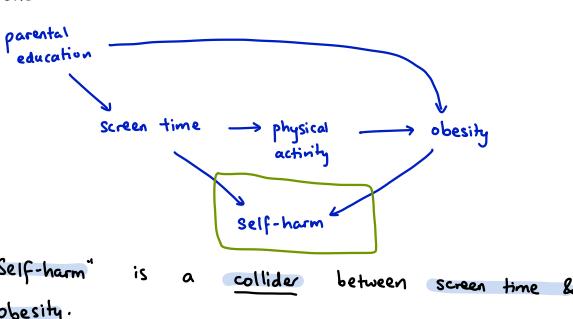
Consider



"Parental education" is a confounder.

COLLIDER (C8S1353)

Consider



"Self-harm" is a collider between screen time & obesity.

DIRECTED PATHS (C8S1356)

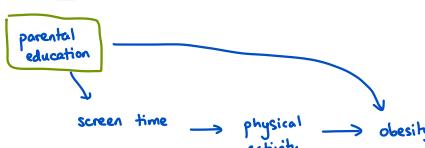
In a "directed path", all arrows point in the same direction.

eg screen time → physical activity → obesity

Any associations represent a causal relationship.

BACKDOOR PATHS (C8S1357)

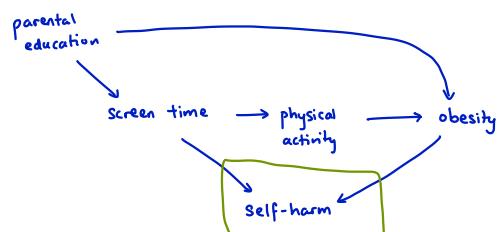
In a "backdoor path", two variables share the same cause.



The association by this path is non-causal.

CLOSED / BLOCKED PATHS (C8S1358)

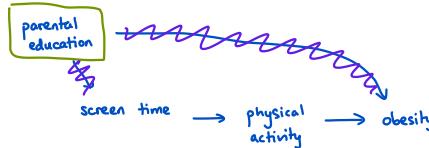
In a "closed path", two variables affect another variable.



CHANGING PATHS (C8S1359)

We can change the status of a path by controlling for or conditioning on a variable for our analysis.

eg



screen time → physical activity → obesity

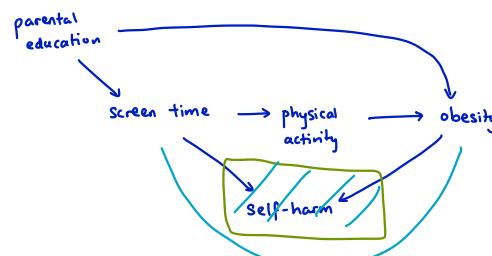
- parental education is a confounder for screen time & obesity
- "including" parental education in our analysis "closes" the backdoor path.

U₂ Note:

- ① Controlling for confounders is good:
 - removes the confounder from analysis
 - so estimate is more accurate.
- ② Controlling for mediators is bad.
 - causes incorrect estimates of the overall association between the variables

screen time → physical activity → obesity

U₃ Conditioning on a collider changes the path from closed to open.

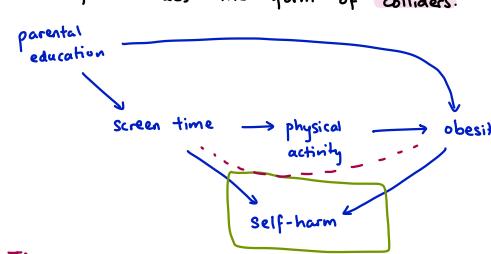


This transmits a non-causal association.

SELECTION BIAS (C8S1367)

Selection bias can be visualized in DAGs.

U₂ This often takes the form of colliders.



- If selected adolescents for the study because of a prior history of self-harm, we can view this as conditioning/controlling that collider.

LIMITATIONS (C8S1369)

U₁ Limitations of DAGs:

- ① They are qualitative; they do not indicate the form, size or direction of a relationship.
- ② They are limited by the information available to form them; it must be complete to give a proper causal interpretation.
- ③ They cannot depict random error; confounding can occur even when treatments are randomized, which is not shown in DAGs.