

- Careful evaluation of Doomsday is important to know

the background of the data & how to best handle it without bias

with this

LECTURE 2 : INFERENCE STATISTICAL ANALYSIS

WIKI

BAG A or BAG B

- ~~Hence to choose blindly between~~ Based on the context a Mr bag ~~given~~ shown, you are to decide which bag you're with.
- Bag A or Bag B? Anyone you get, you pay for it
- Make your decision ~~based~~ ^{using} a frequency plot
 - Revision rule : Look @ the values that are possible for our show of car data and look for which is more likely & reject / no reject the null
 - You get \$10 = More likely to be bag A \rightarrow Not reject null
\$60 = less likely for bag A \rightarrow Reject Null
\$30 = Hard to tell

$$\text{Type I error} = \frac{2}{20} \times 100\% = 10\%$$

$$\text{Type II error} = \frac{6}{20} \times 100\% = 30\%$$

$\alpha = \text{Type I error} = \text{significance level}$
 $\beta = \text{Type II error}$

Based on our decision rule, our $\alpha = 0.10$

- Often ~~we~~ we pick the α option, then we look @ our data to get a p-value \rightarrow how likely is it to get the result you got or more extreme. This is done if instead of being our decision rule of the value we got $\frac{1}{2}$ considering that for every different test

+ lets say we pick \$40 from the withdraw bag, do we reject or fail to reject? lets look @ the p-value = $\frac{4}{20} = 0.2 \approx 20\%$ chance of getting 40 or more. We'd have to ~~fail~~ to reject the null

. ~~Universal~~ Decision rule \rightarrow Reject Null if $P \leq \alpha$

If the likelihood of making an ~~Type II~~ Type I error is smaller than the Type I error

chance of getting 440 or more. We'll have to ~~fail~~ to reject the null

- Universal Decision rule \rightarrow Reject Null if $P \leq \alpha$

If the likelihood of making an ~~Type II error~~ is smaller than the Type I error

Wk 1

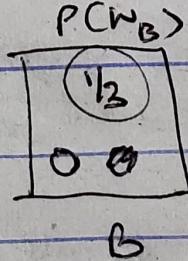
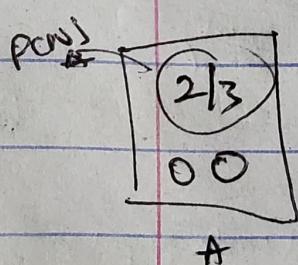
Introduction To BAYESIAN

- There are 2 primary frameworks in stats \geq "Frequentist" \leq "Bayesian" stats
- In frequentist stats, an answer to a question can either be correct (100%) or incorrect (0%). Any notion of belief can only begin in Bayesian

Reason \rightarrow In freq \rightarrow Probabilities made in world

\rightarrow In Bayes \rightarrow Probability made in your head

- You can ^{only} update probabilities in Bayesian statistics



- 2 chances of picking W from bag A
- 1 chance of picking W from bag B

- In Bayesian stats you can update probabilities based on the new information you get

WSK¹

THE PYTHON STATISTICS LANDSCAPE

~~THE TUTORS OR TREAT? LANGUAGE AND NOTATION~~

* 8

- List comprehensions: $\text{My_list} = [\text{expression} \mid i \text{ in } \text{input_list}]$
- Your documents should explain the function, arguments & output
- Lambda functions \rightarrow anonymous functions
lambda arguments : expression like $(\lambda x: x**2)(5)$

WPS

ESTIMATION & POPULATION PROPORTION W/ CONFIDENCE

$$\text{Conf Interval} = \text{Best Estimate} \pm \text{Margin of Error}$$

~~Estimation~~ Whenever we have a statement that we want to infer on we have to come up with ⁽¹⁾ Population & ⁽²⁾ Parameters & infer

$$SE = \frac{\sigma}{\sqrt{n}} \text{ or } \sqrt{\hat{p}(1-\hat{p})} \rightarrow \text{proportion}$$

Wk2 Estimation of Population Proportion w/ Confidence

Conf Interval = Best Estimate \pm Margin of Error

~~Estimator~~ Whenever we have a statement that we want to infer on

We have to come up with ⁽ⁱ⁾ Population & ⁽ⁱⁱ⁾ Parameters of interest

$$SE = \frac{\sigma}{\sqrt{n}} \quad \text{or} \quad \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \rightarrow \text{proportion}$$

- What does CI mean? \rightarrow A range of reasonable values for our param

\hookrightarrow We believe with 95% confidence our population param to be within the CI

Wk2 Understanding Confidence Intervals

- We use our sample proportion to construct the confidence interval

\hookrightarrow Statistics
but the CI isn't referred to a "confidence interval for a sample statistic" but its a range of values to help us estimate the Population Param with a high level of confidence.

- How to best interpret? \rightarrow (i) State the confidence level used
 Your confidence interval gotten.

- (ii) Convey that the data used to estimate the pop param
- (iii) Include that rate for all such populations
- (iv) Indicate the sample size.

Best 10 q

Wk

- H

le

- The population param is unknown hence the need for the study to estimate its value

NOTE: If it's incorrect to ~~interpret~~ the confidence level as a probability or chance

→ It's not a 95% chance that our pop param is in the interval

This is because that ~~interval~~ is fixed & our population param is fixed

The confidence interval refers to our confidence in the ~~method~~ ^{be random}
statistical procedure used to make the confidence interval

- We are confident that with our method used, 95% of the time, it'll contain our pop. parameter

as we only did this once & we used an ^{unbiased} estimate with a representative study & with a 95% conf interval, we are sure that if 95% of the time, ~~the~~ interval gotten will contain the population param.

WK2 Assumptions for a single population proportion conf. interval

- Assumpt. 1: Our sample is a SRS

Study we have a \rightarrow w conf interval, we are sure that if 95% of the time, ~~the~~ the interval gotten will contain the population param.

WK2 Assumptions for a SINGLE POPULATION PROPORTION CONF. INTERVAL

(i) Assumpt. 1: Our sample is a SRS

(ii) Assumpt - 2: We need a large enough sample size so that we can approx our sample dist with a normal curve

How do we check (i) \rightarrow ^{Analyze} Make sure how the sample was collected to make it representative. If complex resampling mentioned in last cont weightage, etc.

WK2 CONSERVATIVE APPROACH & SAMPLE SIZE CONSIDERATION

- Conservative approach is used when we "doubt" our point estimate
- Margin of Error is dependent on \rightarrow confidence level
Sample size
- ~~Usually we pick on~~ Conservatively; you use your desired MoE & conf. level to find out a sample size that will give you your desired MoE
 - What sample size would be need to have a 95% conservative confidence interval with a Margin of Error of only 3% (0.03)?

WK2

* Split screenshot question to each
topic e.g. concrete ...
by

for best estimate

$$\hat{p} = 0.5$$

Practice Quiz

$$\text{Conf. at } = 95\% \quad ; \quad MoE = ?$$

$$SE = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{1}{4n}} = \frac{1}{2\sqrt{n}}$$

$$\text{Conf. Int. } MoE = Z * \frac{1}{2\sqrt{n}} = 1.96 * \frac{1}{2\sqrt{n}}$$

$$MoE_{\text{exact}} = \frac{1.96}{4\sqrt{n}} * \frac{1}{2\sqrt{n}} = \frac{1}{4n} (0.98)$$

$$0.03 = \frac{1}{4n} \rightarrow 0.03^2 = \frac{1}{n} \times 0.98^2$$

$$n \geq 1112 \quad \leftarrow n \geq 1111 + 1$$

$$n \geq 1068 \quad \leftarrow n \geq 1067 + 1$$

WK2

Quiz

Q3

$$MoE = 2 * \frac{1}{2\sqrt{n}} = \frac{1}{\sqrt{232}}$$

Q6

$$MoE = Z * \frac{1}{2\sqrt{n}}$$

$$0.04 = \frac{1.96}{2\sqrt{n}} \rightarrow (\sqrt{n})_c = \left(\frac{1.96}{0.04} \right)^2$$

Q7

$$MoE = \frac{Z}{2\sqrt{n}} = \frac{2.326}{2\sqrt{n}}$$

$$n = 600.25 \rightarrow n \geq 601$$

$$n = \left(\frac{2.326}{2 \times 0.03} \right)^2 \approx 1502.87$$

$$n \geq 1503$$

WK2

ESTIMATING A DIFF. IN POP. PROPS W/ CONFIDENCE *

$$Q.F \quad n_{\text{reqd}} = \frac{z}{2\sigma} = \frac{2.326}{2\sqrt{n}}$$

$$0.04 = \frac{1.96}{2\sqrt{n}} \rightarrow \frac{1.96}{2} = \frac{1.96}{2(0.04)}$$

$$n = 600.25 \rightarrow n \geq 601$$

$$n = \left(\frac{2.326}{2 \times 0.03} \right)^2 \geq 1502.81 \\ n \geq 1503$$

WK2

ESTIMATION A DIFF. IN POP. PROPS w/ CONFIDENCE

Best Estimate \pm Margin of Err

$$\hat{P}_1 - \hat{P}_2 \pm Z_{\text{margin}} \cdot SE_{(\hat{P}_1 - \hat{P}_2)}$$

From slide, $\hat{P}_1 = 0.55$... 55% of parents in our sample had kids by

$\hat{P}_2 = 0.37$... 37% of parents of black kids have had ^{some} swimming lessons

$$\hat{P}_1 - \hat{P}_2 \pm Z_{\text{margin}} \times \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{n_1} + \frac{\hat{P}_2(1-\hat{P}_2)}{n_2}}$$

What 95% conf

$$\rightarrow (0.1123, 0.2477)$$

Interpretation: With 95% confidence, the population proportion of parents of white kids who have taken swimming lessons is 11.23% to 24.77% higher than the population proportion of parents with black children who have taken swimming lessons

Wk 2

INTERPRETATIONS & ASSUMPTIONS for TWO CONFIDENCE POPULATION PROPORTION INTERVALS

Note if D is an estimate - our confidence interval for the difference is estimate, then there is a ~~possibility~~ ^{possibility} that there might be no difference.

- Assumptions:
- (i) We have two independent random samples \rightarrow in the sampling process
 - (ii) We need large enough sample [We need at least 10 year's & 100s 10 n's for each sample]

Wk 2

ESTIMATION A POPULATION MEAN w/ CONFIDENCE

- Now our response is quantitative instead of categorical (proportion)
- means

Graphical Summary

- Histogram: Shows that ~~and the~~ ^{value} far below 70 as to above 110
we see a right skewed in our data

- QQ plot = Also shows us \uparrow with the tail @ the lower end off.

Numerical summary \rightarrow min, max, mean, SD, n

Wk 2

Estimation A Population Mean w/ Confidence

- Now our response is quantitative instead of categorical (proportions)

means

Graphical Summary

- histogram: Shows that ~~we're~~ ^{set} far below 70 up to above 110
we see a right skew in our dist

- QQ plot = Also shows us ↑ with the tail @ the lower end

Numerical summary → min, max, mean, SD , n off.

↳ The measures varied from the mean of X with an average
of SD in .

- Our best estimate is the sample mean but it'll vary if # ^{with greater} ~~greater~~

Sampling

Q: Does the response variable need to be normal for the sampling
mean dist to be normal?

- Standard Error of the Sample Mean = $\frac{\sigma}{\sqrt{n}}$ depends on how much varab & how much data we have

We see a right skew in our dist

- QQ plot = Also shows us ↑ with the tail @ the lower end

Numerical summary → min, max, mean, SD, n \bar{x} .
↳ The measures varied from the mean of X with an average SD in.

- Our best estimate is the sample mean but it'll vary ^{with} ~~the~~ repeated sampling

Q // Does the response variable need to be normal for the sample mean dist to be normal?

↳ Standard Error of the Sample Mean = $\frac{\sigma}{\sqrt{n}}$ depends on how much variah & how much data was sampled

Drawback: We can never calculate it because

we don't know the population SD (σ). ∴ We need to come up with an estimate.

► Estimate Standard Error of the Mean = $\frac{s}{\sqrt{n}} \xrightarrow{s \rightarrow \text{sample SD}}$

95% Confidence Interval

$$\bar{x} \pm t^* \left(\frac{s}{\sqrt{n}} \right)$$

- Using t statistic? - Here you don't need a large sample size to conduct a confidence interval for a population mean unlike the requirement for that of proportions.
- And since we're using (S) ; we'd have to adjust for multiplicity determining the distance from the mean and t comes from a t distribution with $n-1$ degrees of freedom
- The t statistic varies based on the degree of freedom which depends on the sample size.

E.g. For 95%

$$n = 25 \rightarrow t^* = 2.064$$

$$n = 1000 \rightarrow t^* = 1.962$$

S differs depending on the sample size. It gets closer to the 6 as n increases. So t* account for this.

Interpreting the Confidence Level

- It's giving us a range of reasonable vals for our parameter

dependent on the sample size.

E.g. for 95%

$$n = 25 \rightarrow t^* = 2.064$$

$$n = 1000 \rightarrow t^* = 1.962$$

on the samples
It get closer to the
true θ as n increases
So t^* account
for that

Interpreting the Confidence level:

- It's giving us a range of reasonable vals for our param

e.g. With 95% confidence, the population mean Cartwheel Dist for all adults is estimated to be between 76.26 inches & 88.70 in.

We have context, what param we're estimating & Conf level
our confidence is in the procedure we're using

we would expect 95% of those resulting intervals to contain the population mean cartwheel distance

Assumptions for Confidence Interval for Population Mean

- data considered from a random sample

- population of responses is normal (large sample size helps)

WK2

Estimating a mean difference for paired data

- With paired data we treat the simultaneously - because they're measured very

L.e.g.: How someone writes w/ both hands, pre & post treatment of an individual
maternal outcomes data ex. twin data, due to blood relation, by age or gender or weight.

variable: difference of measurements b/w pairs.

Population = all identified firms

Sample of interest = μ_d (diff = older - younger)

Goal: Construct 95% interval for the mean diff of ...

$$\text{CI golden} = (0.0025, 0.1652) \text{ yrs}$$

Interpreting the CI: W/ 95% the pop mean diff of the older tennis

(as the younger tennis self reported education extended to be longer

0.0025 yrs and 0.1652 yrs

wth 95% conf

variable: difference of measurements b/w pairs.

Population = all identified tuis

mean of intercs = μ_d ($\text{diff} = \text{older} - \text{younger}$)

goal: Construct 95% interval for the mean diff of ...

CI gotten = $(0.0025, 0.1652)$ yrs

Interprets the CI: w/ 95% the pop mean diff of the older tuis

(as the younger tuis self reported educatn & estimate to be b/w

0.0025 yrs and 0.1652 yrs.)

Means that both ends are on the +ve end, we expect - the older tuis

has more education. With values being close to 0, that educatn jump might
not be that large.

Assumptions → random sample of identified tuis

→ population of diff is normal (or large enough sample size)

Estimation A DIFFERENCE in Population MEANS w/ CONFIDENCE

Wk 2

INDEPENDENT GROUPS

Interprets the CI : W(95%). The pop mean diff of the older tennis
group the younger tennis self reported educate & estimate to be b/w
0.0025 yrs and 0.1652 yrs.

Means that both ends are on the +ve end, we expect - the older tennis
group has more education. With values being close to 0, that education jump might
not be that large.

- Assumptions
- > Random sample of identified tennis players
 - > Population of diff is normal (or large enough sample size)

Wk 2

ESTIMATION A DIFFERENCE IN POPULATION MEANS W/ CONFIDENCE

<For independent groups>

- Sample based on gender
- o Population: Mexican-American adults (18-29 yrs)
 - o Parameter of interest ($\mu_1 - \mu_2$): BMI (kg/m^2) \rightarrow A quantifiable & homogeneous variable

2 approaches for deriving the MoE

(i) Pooled approach : If variances are close to each other, we can assume them as equal

(ii) Unpooled approach : Assumption of equal var. dropped

• for df for finding t^* use the smaller of $n_1 - 1$ and $n_2 - 1$

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\rightarrow \text{est. standard Err} < \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$\therefore df = n_1 + n_2 - 2$$

L. This is only if pop. variances are equal

QQ plot compare sample quantiles w/ quantiles of a normal distribution

• We compare (A.R. \times box plots) \pm 1 sds of both groups to find

but if we use pooled / unpooled variance

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\rightarrow \text{est. standard Err} = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

. df = $n_1 + n_2 - 2$

. This is only used if pop. variances are equal

QQ plots compare sample quantiles w/ quantiles of normal data

• We ~~use~~ compare (AR or box plots) \pm 2 sds of both groups to find

out if we use pooled (unpooled variance)

Based on CI (-0.385, 1.865) ^{with 95% confidence}, there's a possibility of there being no difference in population means.

• Proportions have a special property, that is the variance is completely determined by the mean.

?? ??
WR? Normal distribution: An introductory course to PDF

& CDF → <copy & past from website>

QQ plot compare sample quantiles w/ quantiles of a normal distribution

- We compare (QRs / box plots) $\pm 1.5 \text{ sds}$ of both groups to find out if we use pooled (unpooled variance)

Based on CI (-0.385, 1.865) ^{with 95% confidence}, there is a possibility of there being no difference in population means.

- Populations have a special property, that is the variance is completely determined by the mean.

?? ?? Wk 2 Normal Distribution: An introductory Guide to PDF

§ CDF → <copy & past from website>

- A normal distribution (Gaussian dist) is a continuous probability distribution for real-valued variables. For a normal dist, other values other than the mean are less probable.

- CLT → Laplace observed that even if a population doesn't follow a normal distribution, as the number of the samples taken increases, the distribution of the sample means tends to be a normal distribution.

The probability density function (pdf) & cumulative distribution function (cdf) helps us determine the probabilities & range of probabilities when data follows a normal distribution.

- ~~PDF~~
- The CDF is the integral, from left to right, of the PDF
 - The pdf is a statistical expression that defines a probability distribution (the likelihood of an outcome) for a discrete random variable as opposed to a continuous random variable. When pdf is graphically represented, the area under the curve = the interval in which the variable will fall

This bias of an estimator is the difference b/w the estimator's expected value and the true value of the parameter being estimated.

Why we use $(n-1)$ instead of n when calculating sample variance?

→ Variance is the average of the sum of squares of the diff. of the obsv. from the mean. So, when we use the sample mean as an approx of the pop. mean for calcul. the sample variance, the num (i.e. \sum squared diff. from the mean) can be small @

(the likelihood of an outcome) for a discrete random variable as opposed to a continuous random variable. When pdf is graphically portrayed, the area under the curve \equiv the interval in which the variable will fall.

- Bias: The bias of an estimator is the difference b/w the estimator's expected value and the true value of the parameter being estimated.

Why we use $(n-1)$ instead of n when calculating sample variance?

→ Variance is the average of the sum of squares of the diff. of the obsv. from the mean. So, when we use the sample mean as an approx of the pop mean for calcul. the sample variance, the num (i.e. \sum squared diff. from the mean) can be small @ times. In those cases, we will get smaller sample variances. Hence when we divide the sample var. by n , we underestimate (i.e. get a biased value) the population var. In order to compensate for this, we make the denominator of the sample var $(n-1)$ to obtain a larger value. This reduces the bias of the sample variance as an est. of the pop. var.

$$-1/2 \cdot \left(\frac{x-\mu}{\sigma}\right)^2$$

the sample variance, the sum (i.e. \sum squared dist. from the mean) can be small @ times. In those cases, we will get smaller sample variances. Hence when we divide the sample var. by n , we underestimate (i.e. get a biased value) the population var. In order to compensate for this, we make the denominator of the sample var. ($n-1$) to obtain a larger value. This reduces the bias of the sample variance as an est. of the pop. var.

Standard Normal pdf: $f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi \cdot \sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{x-\mu}{\sigma}\right)^2}$

A standard normal distribution has a pop. mean of 0 & standard dev = 1.

If we subtract the mean from each value in the sample & divide by the standard deviation

Reasons why we would do this? (i) The data is widely spread out

(ii) To put all our variables on the same scale because most times our data might be measured in diff. scales (units) and might affect our interpretation

• If Z is standard $\text{D}\sigma$, $f(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$

• CDF \rightarrow It's a function derived from the probability density function for a continuous random variable. It gives the prob. of finding the random variable @ a value less than or equal to a given cutoff i.e., $P(X \leq x)$

Note: Chap in "How to build a 95% CI, what is t^* " from week 2 went into notes.

Wk3 ~~1~~ Section 1: A test for a population proportion $\frac{1}{2}$ DOP. Proportion

Why carry out hypothesis tests? \rightarrow we have a question about things

Let the data to support the claim concerning the question

- o You make your hypothesis first before data collection to avoid influencing our beliefs after data collection

$$H_0: p = 0.52 ; H_a: p > 0.52$$

• state what p is \leftarrow note $: \alpha = 0.05 \rightarrow$ cut off pt at which we found stat. to be significant

- o Next a significance level: $\alpha = 10\%$ \rightarrow good

Note: Chap in "How to build a 95% CI, what a ~~fun~~ went
into notes". Testing a one

Wk3 ~~stat~~ SETTING UP A TEST FOR A POPULATION PROPORTION. $\frac{1}{2}$ DOP. Propose

Why carry out hypothesis? - We have a question don't know
[or use data to support the claim concerning the question]

- You make your hypothesis first before data collection to avoid influencing your beliefs after data collection

$$H_0: p = 0.52 \quad ; \quad H_a: p > 0.52$$

• State what p is \leftarrow note ; $\alpha = 0.05 \rightarrow$ cut off pt at which we found smth. to be significant

- Next survey collection; $n = 1018$; $\hat{p} = 56\%$

{ Assumption check = ⁽ⁱ⁾ Random sample ⁽ⁱⁱ⁾ Large sample size to ensure ⁽ⁱⁱⁱ⁾ $n \cdot p \geq 10$
 $n(1-p) \geq 10$

- Run hypothesis test

$$\text{Test statistic} = \frac{\text{Best estimate} - \text{Hypothesized est.}}{\text{SE of estimate}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

SE of estimate

beliefs after data collection

$$H_0: p = 0.52 \quad ; \quad H_A: p > 0.52$$

- State what p is \leftarrow note ; $\alpha = 0.05$] \rightarrow cut off pt at which we found smth. to be significant

- Next survey collection: $n = 1078$; $\hat{p} = 56\%$

[Assumption check = ⁽ⁱ⁾ Random sample

⁽ⁱⁱ⁾ Large sample size to ensure roughly

$$n \cdot p_0 \geq 10$$

$$n(1-p_0) \geq 10$$

- Run hypothesis test

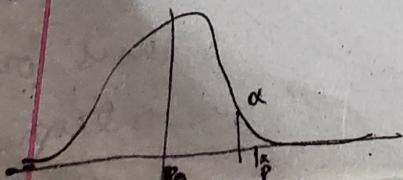
$$\text{test statistic} = \frac{\text{best estimate} - \text{hypothesized est.}}{\text{SE of estimate}} = \frac{\hat{p} - p_0}{\text{SE}}$$

$$\text{S.E. } (\hat{p}) = \sqrt{\frac{p_0 \cdot (1-p_0)}{n}}$$

$$Z = \frac{0.56 - 0.52}{0.0157} \approx 2.555$$

\hookrightarrow null standard error [under the H_0]

* I test statistic because we're under a normal dist. When very proportions, we use a norm dist. as long as n is large enough



Interpret the Z-test statistic: Our ~~observed~~ sample proportion is 2.55% null standard errors above our hypothesized population proportion.

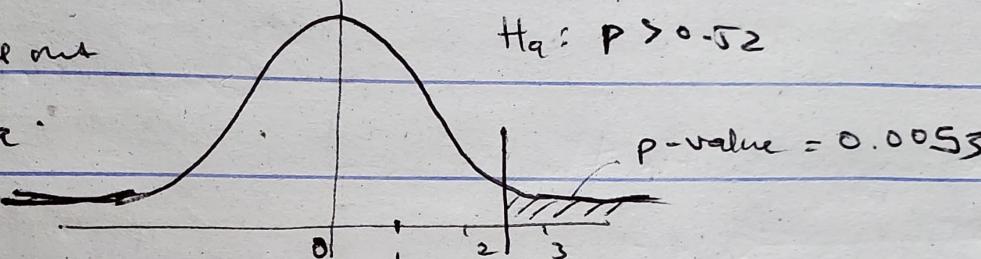
A Z-test statistic

- With the Z-test statistic, we find the p-value

Note that p-value α

the same as P used as H_a .

That our ^{null} proportion



Since $p\text{-value} < \alpha$; we reject H_0 .

- Conclusion statement: There is sufficient evidence to conclude that the population proportion of parents ...

WK3 SETTING UP A TEST OF DIFFERENCE IN Population Prop.

- Population = All parents of black kids ($6-18$ yrs) \neq all parent of hispanic children $<6-18>$

- Parameter of interest $\rightarrow P_1 - P_2$ (1 -black kids, 2 -hispanic kids)

- significance level = α ; $\alpha = 0.10$

since $p\text{-value} < \alpha$; we reject H_0 .

- Conclusion statement: There is sufficient evidence to conclude that the population proportion of parents ...

Wk3 SETTING UP A TEST OF DIFFERENCE IN Population PROP.

- population = All parents of black kids $\langle 6-18 \text{ yrs} \rangle$ & all parents of hispanic children $\langle 6-18 \text{ yrs} \rangle$

- Parameter of interest $\rightarrow P_1 - P_2$ (1 -black kids, 2 -hispanic kids)

- significance level $= 10\%$; $\alpha = 0.10$

$$H_0: P_1 - P_2 = 0 \quad ; \quad H_a: P_1 - P_2 \neq 0$$

- survey results $n_1 = 247$; $\hat{P}_1 = 91/247$

$$n_2 = 308 \quad ; \quad \hat{P}_2 = 120/308$$

- Check our assumptions

- Calc best estimate of the parents; $\hat{P}_1 = 0.37$, $\hat{P}_2 = 0.39$

Children $< 6-18 >$
age

- Parameter of interest $\rightarrow P_1 - P_2$ (1 -black kids, 2 -Hispanic kids)
- Significance level $= 10\%$; $\alpha = 0.10$

$$H_0: P_1 - P_2 = 0 \quad ; \quad H_a: P_1 - P_2 \neq 0$$

Survey results
 $n_1 = 247$; $\hat{P}_1 = 91/247$
 $n_2 = 308$; $\hat{P}_2 = 120/308$

- Check our assumptions
- Calc best estimate of the params ; $\hat{P}_1 = 0.37$, $\hat{P}_2 = 0.39$
- $\hat{P}_1 - \hat{P}_2 = -0.02$

M2

TESTING A DIFF. IN POPULATION PROPORTIONS

$$\text{test statistic} = \frac{(\hat{P}_1 - \hat{P}_2) - 0}{SE_{(\hat{P}_1 - \hat{P}_2)}} \quad SE_{(\hat{P}_1 - \hat{P}_2)} = \sqrt{\hat{P}_1(1-\hat{P}_1)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

$$SE_{(\hat{P}_1 - \hat{P}_2)} = \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{n_1} + \frac{\hat{P}_2(1-\hat{P}_2)}{n_2}} \quad \approx \text{test stat} \approx -0.4831$$

$$= 0.0414$$

• Check many of 2 stats in note

• P-value (z) = 0.63

↳ two tailed

• since $P\text{-value} > \alpha$; we fail to reject null hypothesis

• Alternative approaches to this question

↳ ② Fisher's Exact test = for ~~too~~ smaller ~~small~~ sample size

P-value, P-hacking & more

• P-value - level of surprise. Talks about the surprise of the data. Low = higher exp of surprise,

• P-value doesn't reveal the level of the truth of the data.

• Having a low value could mean statistical significance but not necessarily practical significance

↳ You're almost guaranteed to get ~~small~~ low P-value w/ large data

• P-hacking - whenever we can no longer interpret the p-value by its mathematical definition. You need to know all done on the

higher exp of surprise,

- P-value doesn't reveal the level of the truth of the data.
- Having a low value could mean statistical significance but not necessarily practical significance
 - You're almost guaranteed to get ~~low~~ small p-value w/ large data
- P-hacking - whenever we can no longer interpret the p-value by its mathematical definition. You need to know all done on the ways of hacking : many multiple tests on a data
 - consequences : leads to uninformative studies
- p-value hacking also affects all the other estimates
- For pure p-value, pre-register, ~~if~~ write down all the choices, methodology that you want to follow. This helps you follow a more confirmatory test than an exploratory test (which ^{leads to} ~~leads to~~ more p-hacking)
 - Include all the steps taken in the study, not just what ~~leads to~~ to the p-value

WK3

ONE MEAN: TESTIN ABOUT A POPULATION MEAN w/ CONFIDENCE

- **p-hacking** - whenever we can no longer interpret the p-value by its mathematical definition. You need to know all done on the
 - ways of hacking : many multiple tests on data
 - consequences : leads to uninformative studies
- p-value hacking also affects all the other estimates
- For pme. p-value, pre-register, ~~and~~ write down all the choices, methodology that you want to follow. This helps you follow a more confirmatory test than an exploratory test (which ^{leads to} more p-hacking)
 - Include all the steps taken in the study, not just what ~~leads~~ to the p-value

WK3

OVERVIEW: TESTINg ABOUT A POPULATION MEAN w/ CONFIDENCE

population : All adults

param & vars : population mean CNT, μ

- Define hypotheses $\rightarrow H_0: \mu = 80$, $\alpha = 0.05$
 $H_1: \mu > 80$

Significance level: \rightarrow the results should be so unusual that we would see them no more than 5%

$$n = 25, \text{ min} = 63, \text{ max} = 115, \bar{x} = 82.48, s = 15.06 \text{ m}$$

Example assumptions — (i) Random sample

(ii) Normal distribution for all distances

- histograms show right skew but w/ 25 obs we can use the CLT.

Is the sample mean of 82.48 a sign. greater than hypothesized mean of 80 m. $SE \text{ of the sample mean} = \frac{s}{\sqrt{n}}$ — we don't have pop std of

estimated SE st
the obs $= \frac{s}{\sqrt{n}}$

$$t\text{-statistic} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = 0.82$$

our sample mean is only 0.82 estimated ST above the null value of 80 m

$$n = 25, \text{ min} = 63, \text{ max} = 115, \bar{x} = 82.48, s = 15.06 \text{ m}$$

Example assume

(i) Random sample

(ii) Normal distribution for CW distances

- histograms show right skew but w/ 25 obs we can use the CLT.

Is the sample mean of 82.48 m sig. greater than hypothesized mean

of 80 m. $SE \text{ of the sample mean} = \frac{s}{\sqrt{n}}$ — we don't have pop std

estimated SE st
the obs $= \frac{s}{\sqrt{n}}$

$$\text{t-test statistic} = \frac{\bar{x} - \mu}{SE} = \frac{82.48 - 80}{15.06/\sqrt{25}} = 0.82$$

our sample mean is only
0.82 estimate ST above
the null value of 80 m

Determine p-value | very t-test dist table; p-value = 0.21

since $p\text{-value} > \alpha (0.21 > 0.05)$; we fail to reject the null hypothesis

- Based on our est. sample mean (82.48 m), i.e. . . .

of 80 in.

SE of the sample mean = $\frac{\sigma}{\sqrt{n}}$ — we don't have pop std of

estimated SE st

$$\text{the aa} = \frac{s}{\sqrt{n}}$$

$$\frac{1}{\sqrt{n}}$$

$$\cdot t\text{-test statistic} = \frac{\bar{x} - \mu}{SE} = 0.82$$

our sample mean is only
0.82 estimated SE above
the null value of 80 in

o Determine p-value | very t-test dist table ; p-value = 0.21

since $p\text{-value} > \alpha (0.21 > 0.05)$; we fail to reject the null hypothesis

- Based on our est. sample mean (82.48 in), etc...

- Looking @ 90% CI estimate

WR3 Testing a population mean different

o In what way is our data paired? How wif paired?

- Population = all hours

| param of later = pop mean diff + qual
Gained
↳ Med

$$\alpha = 0.05$$

$$H_0: \mu_d = 0 \quad H_1: \mu_d \neq 0$$

In

$$\cdot t\text{-test statistic} = \frac{\bar{x} - \mu}{SE} = 0.82$$

our sample mean is only 0.82 estimate SE above the null value of 80 in

- Determine p-value | very t-test dist table ; p-value = 0.21
since $p\text{-value} > \alpha (0.21 > 0.05)$; we fail to reject the null hypothesis

- Based on our est. sample mean (82.48 vs) , we ...

- Looking @ 90% CI estimate

WR3 Testing A Population Mean Difference

- In what way is our data paired? How wif paired?

- Population = all rows

$$\alpha = 0.05$$

param of inter = pop mean diff + qual
 ↳ Md

$$H_0: \mu_d = 0 \quad H_1: \mu_d \neq 0$$

- Assumptions - Random sample of diff

normal dist. of (sample size of 25+) | we have a rough bell shape

$$n=20, \bar{x}_d = \$17.30, s_d = \$28.49$$

$$t\text{-test statistic} = \frac{\bar{x}_d - \mu_d}{SE_d} = \frac{17.30 - 0}{\frac{28.49}{\sqrt{20}}} = 2.72$$

p-value = 0.014 ; p-value < 0.05 - reject H₀

95% CI : (\$3.97 ; \$30,63)

Wk3

TESTING FOR DIFFERENCE IN POPULATION MEANS (FOR INDEPENDENT CATEGORIES)

~~-opposite~~: Mex.-Amuræ adults in N.Y.

-^o Param & inter ($M_1 - M_2$): Body Mass Index

→ Define hypotheses ; $H_0: \mu_1 - \mu_2 = 0$; $H_a: \mu_1 - \mu_2 \neq 0$

$$X_1 = 23.57, \bar{X}_2 = 22.83; n_1 = 258, n_2 = 239$$

$$S_1 = 6.24 \quad S_2 = 6.43$$

→ Check Assumptions → (i) Simple random sample

(a) iid (iii) Normality or not large

~~— B. Calc: test stages~~

$$* \text{t-test statis} = (\bar{x}_1 - \bar{x}_2) - 0$$

$$s_1 = 6.24$$

$$s_2 = 6.43$$

→ Check Assumptions (i) Simple random sam

(ii) iid (iii) Normality or n is large

→ Calc. test statistic

$$\text{t-test statistic} = \frac{(\bar{x}_1 - \bar{x}_2)}{\text{SE}_d} = 0$$

$$\text{Best estimate: } \bar{x}_1 - \bar{x}_2 = 0.74$$

$$\text{SE}_d \text{ (approximated approach)} = \sqrt{\frac{s_p^2}{n_1} + \frac{1}{n_2}} \quad t = 1.30,$$

$$\text{SE}_d = \text{using pooled approach} \quad <1 \text{ expects the 10% b/w p & q}> = s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

" In terms of actual calculation, the CLT doesn't differ much from assuming that the data is normal. The key difference is that w.r.t the CLT, you do not need to assume that the data is non-dist-free. Instead, you can rely on the fact that as the sample size increases, the samp. dist. of the sample mean becomes more. norm., regardless of the shape of the original dist.

- o Never p-value !: p-value = 0.19 | we fail to reject Null; $p > 0.05$
- o Chp in the 8th paragraph under "hypotheses Test": Other Cows > Nandos
at their
→ this gives us more insight as to why t-test is used over Z for normal
dist.

~~Practical Notes~~

Q4 The importance of good Research Questions for Sound Inference

- o Data is everywhere and easy to be gotten but still we need a well formulated research quest that informs our study

Key aspects of Ques

Is research question descriptive

→ Target population & info - Is Research question for analysis

→ Has the quest been asked b4? What will the new study add?

→ Are the variables ready available & measurable & they reflect exactly the concepts we want

→ don't check aw & Key aspects?

Key aspects of Ques

- ↳ Is research question descriptive
- ↳ Target population & units → Is research question for analysis
- Has the question been asked before? What will the new study add?
- Are the variables ready available & measurable & they reflect capture the concepts we want

Bad Questions → doesn't check all 4 key aspects

or measure? (e.g., attendance, assign; * mean (sum) - ^{income})

Good questions

Wk4

~~DESCRIPTIVE~~ ~~DESCRIPTIVE~~ WF Examples for some variables using CI

- ~~Remember~~ remember that the confidence interval estimates for an ~~is~~ actual population parameter

Wk4

Descriptive WF Examples for some variables using HYP. TEST

Wk4

Comparing means for two independent samples : An example