

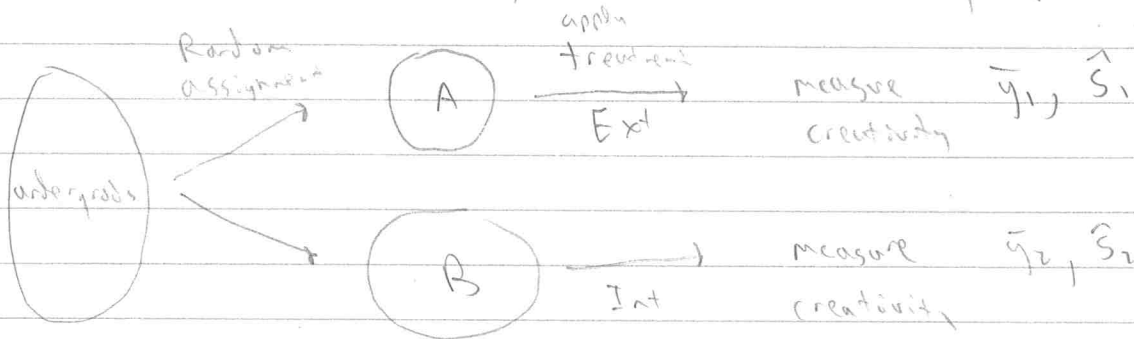
Learning Objectives: Causality, generalizability, sampling distributions, randomization/permutation tests.

Stat 514, Sections 1.1 to 1.4 in Statistical Search

• Case Study 1.1.1:

• Design:

- Conscript 47 undergrads to participate
- Randomly assign 24 to one group and 23 to another
- Apply different treatments to each group
- Measure creativity score within each group
- Compare creativity scores between groups



↑ Can differences in average creativity be attributed to random assignment?

• Case Study 1.1.2:

• Design:

- Get all clerks at bank and measure their sex and income
- Compare new income between the sexes
- Can differences be attributed to random chance "mental model" where randomly assign salaries?

- Two types of Treatment allocations

sign 1

- Observational Study

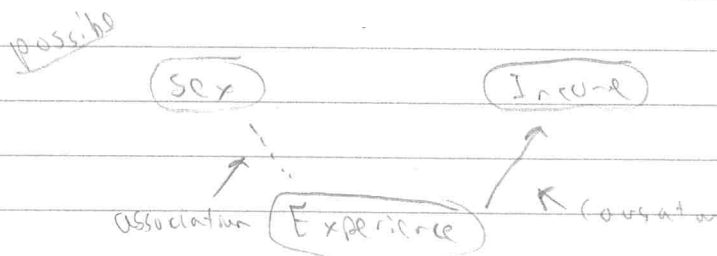
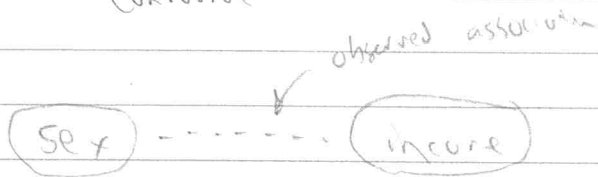
- No treatment imposed on the observational units
- Goal - describe (infer) properties of a population
- Case 1.1.2

I sex not determined by investigators

I want to make statements about clerks at the bank

- cannot claim causation

I a third variable might cause an association
 Confounder



Side Note: - other ways that Education can be incorrect

Hypothesized

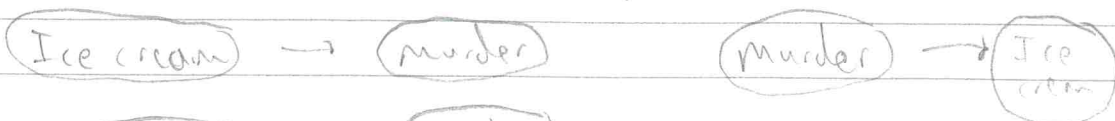


possible

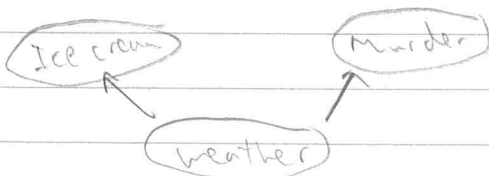


observed Ice cream ↑ and murder rate ↑

Hypothesized



possible



- Causation can never be fully proven by an observational study

- Design 2: Randomized experiment

- ↑ treatment randomly assigned to observational units

- ↑ case 1.1.1

- ↑ can infer causation

- Why can we claim causation?

- ↑ all other variables are at equal levels on average

- ↑ eg, possible to place really good writers in INT group and really bad in EXT, but not really because we made groups by random.

- Do observational studies b/c

- 1.) Goal is prediction, so causation does not matter

- ↑ predict disease, correlating variables etc

- 2.) Combine w/ science to establish causation

- ↑ smoking, animal trials

- How are observational units selected?

- Random sampling

- ↑ results are generalizable to population

- Non-random sampling

- ↑ scope of inference is limited.

- ↑ Random sampling makes sample look like population - some percent women, race, experience level etc.

- Eg. Case 1.1.1: only good writers selected, so results may only hold for good writers.

	Random Experiment	Observational study
Random Selection	causation generalizable	observation possibly non-generalizable
Non-random Selection	causation non-generalizable	observation possibly non-generalizable

• Inference:

• General Setup:

1.) Model with population parameters

2 parameters describe aspects of population we are interested in

Add-treat
+ treatment
model effect

Y_i = creativity score of subject i in Ext-group

Y_i^* = " " " " " " Int "

$Y_i^* = Y_i + \delta$ ← Add-treat treatment Effect Model

↑
parameter, same for all individuals.

$\delta > 0 \Rightarrow$ Int improves creativity

$\delta < 0 \Rightarrow$ Ext improves creativity

$\delta = 0 \Rightarrow$ No effect

2.) Set a hypothesis based on parameters

$H_0: \delta = 0$

$H_A: \delta \neq 0$

Null hypothesis = "simpler state of affairs"

3.) Test statistic

test statistic = function of sample used to measure plausibility of a hypothesis

Expect $\bar{Y}_1 \approx \bar{Y}_2$
 ↑ ↑
 Ext Int
 creativity creativity
 average average

So test statistic = $\bar{Y}_2 - \bar{Y}_1$

4.) Evaluate how rare our test statistic is under the Null hypothesis

RANDOMIZATION DISTRIBUTION SLIDES

• P-value = Probability we would have observed a test statistic as extreme or more extreme than what we observed if H_0 were true.

GO over Twin Study, Case Study 2.1.2

- Paired t-test § 2.2
- Mean Difference in twin study = 0.1987. This difference is random.
- Recall Sampling distribution-

Take a sample of 15 twins

↳ average difference

Take a sample

↳ average difference

Take a sample

↳ average difference

Sampling distribution



Slides on Sampling Dist

- Only one sample actually taken
↳ how can we know anything about it?
↳ statistical theory!
- Thm: If population mean = μ
Population SD = σ
then mean of sampling distribution of average = μ
SD of sampling distribution of average = σ/\sqrt{n}
↳ n = sample size
- Central Limit theorem: Shape of sampling distribution of sample average is more nearly normal for larger and larger n .

- Recap: Even though we don't know the distribution of X , we know the distribution of \bar{X} (for large n).

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

Sample average \nearrow \bar{X}
 μ \leftarrow $E[\bar{X}]$
 σ^2/n \leftarrow $\text{Var}(\bar{X})$ \leftarrow Sample size

- This is useful because we want to know how extreme \bar{X} is
- General setup of statistics (again)

1.) Posit. Model

$$E[X] = \mu \quad (\text{expected difference is } \mu)$$

2.) Set interpretable hypotheses in terms of parameters

$$H_0: \mu = 0 \quad (\text{No difference})$$

$$H_A: \mu \neq 0 \quad (\text{some difference})$$

3.) Come up with test statistic

$$\bar{X} \quad \frac{\bar{X}}{\sigma/\sqrt{n}}$$

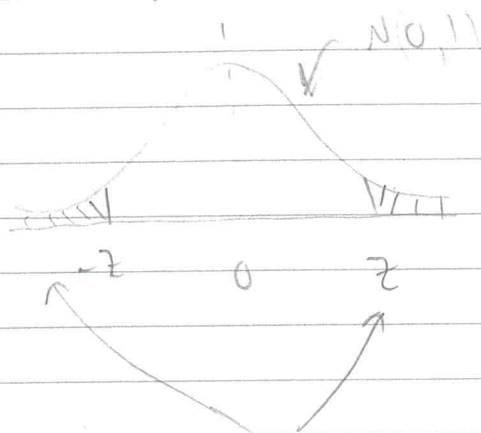
\uparrow larger values of each provide evidence against H_0

4.) Determine how rare test stat is under H_0

under H_0 , $\bar{X} \sim N(0, \sigma^2/n)$

$$Z = \frac{\bar{X}}{\sigma/\sqrt{n}} \sim N(0, 1)$$

So compare Z to a $N(0, 1)$



shaded regions provide evidence against H_0

• Problem: Don't know σ^2

Solution: we can estimate σ^2 w/ s^2

	↑	↑
	pop. variance	sample variance

But then it is better to use

$$\frac{\bar{X}}{s/\sqrt{n}} \sim t_n \quad \text{where } n = n-1 = \text{"degrees of freedom"}$$

• Why t ?

- Estimating σ^2 , so need to account for added uncertainty
- t is more variable than $N(0,1)$ (longer tails)
- You will be more confident in weakness of \bar{X} if you use $N(0,1)$

• R code: `t.test()`

• Two studies: $\bar{Y} = 0.1900$
 $S/\sqrt{n} = 0.0615$
 $N = n - 1 = 14$

$$T = \frac{\bar{X}}{S/\sqrt{n}} = 3.236$$

Compare to t_{14}



$$p\text{-value} = 0.006$$

↑ very rare if null is true!

- So reject H_0 and conclude difference in means

• Confidence Interval for Effect

- Question: What are some plausible values for μ ?
- Use general result

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

- Then we know that in 95% of samples, the following is true

$$-t_{n-1}(0.025) \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{n-1}(0.975) = t_{n-1}(0.975)$$

To solve for μ

$$\Rightarrow -t_{n-1}(0.025) S/\sqrt{n} \leq \bar{X} - \mu \leq t_{n-1}(0.975) S/\sqrt{n}$$

$$\Rightarrow \bar{X} - t_{n-1}(0.025) S/\sqrt{n} \geq \mu \geq \bar{X} - t_{n-1}(0.975) S/\sqrt{n}$$

• Note: $-t_{n-1}(0.025) = t_{n-1}(0.975)$

• 95% Confidence Interval: $\bar{X} \pm t_{n-1}(0.975) S/\sqrt{n}$

• Interpretation: See slides

Two Sample Inference:

Go over Beach Study, Case Study 2.1.1

- Goal: still to explore difference in population means, but now have two independent samples.
- General Setup of Statistics (again)

1.) Posit a Model

Let $X_{1j}, X_{n1} =$ Beach Depths in 1976

$Y_{1j}, Y_{n2} =$ Beach Depths in 1978

$$X_i = \mu_1 + \epsilon_i$$

↑
mean Beach

depth in 1976

↖
some noise, centered

at 0, variance = σ_1^2

$$Y_j = \mu_2 + \delta_j$$

↑
mean Beach

depth in 1978

↖
some more noise,

centered at 0,

variance = σ_2^2

2.) Set interpretable hypotheses

$$H_0: \mu_1 = \mu_2 \quad (\text{same mean beach-depth})$$

$$H_A: \mu_1 \neq \mu_2 \quad (\text{different mean beach-depth})$$

3.) Come up with a test statistic

$$\bar{X} - \bar{Y}$$

↑ values further from 0 support bigger difference

4.) Determine how rare test stat is under H_0

Under H_0 : $\bar{X} - \bar{Y} \sim N\left(0, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$

$$\Rightarrow \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

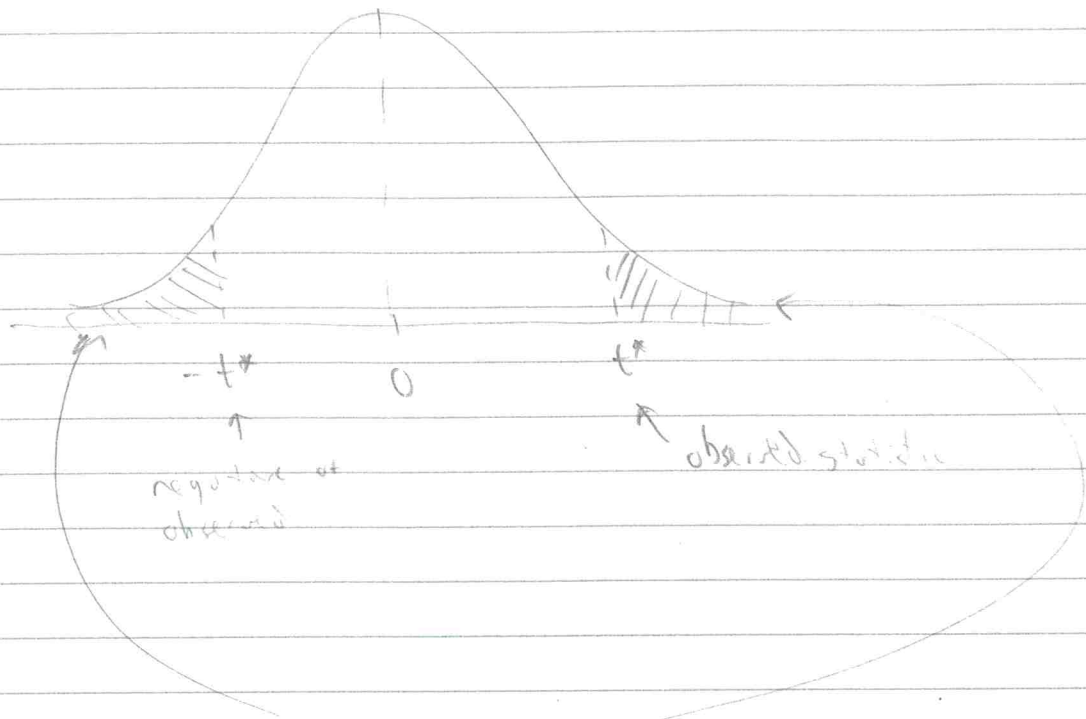
• Problem: don't know σ_1^2, σ_2^2

• Solution: replace with s_1^2 and s_2^2

$$\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t_n$$

n = something weird "Satterthwaite's Approximation"

• We compare our observed test-statistic to the theoretical "null distribution"



proportion of t -stats that would provide as much or more evidence against the null than our observed t -stat.

$p\text{-value} = P_r(|T| \geq |t^*|)$ where $T \sim t_{n-2}$

- Small p -values provide evidence against the null.
- Can also calculate confidence intervals for differences in means:

$$\bar{X} - \bar{Y} \pm t_{\alpha/2}(0.975) \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Some Things to Think About

1.) One-sided tests

↑ set alternatives to be one-sided

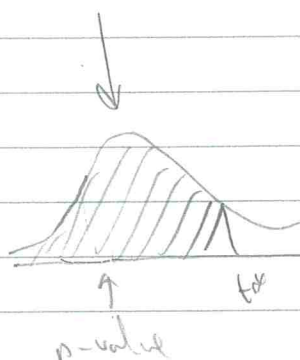
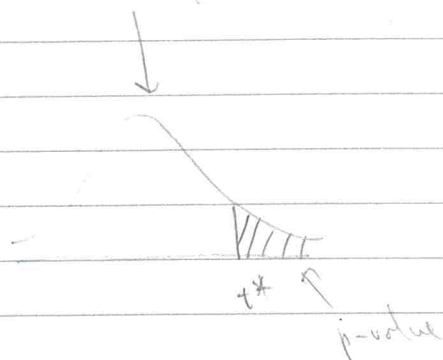
$$H_0: \mu = 0$$

or

$$H_0: \mu = 0$$

$$H_A: \mu > 0$$

$$H_A: \mu < 0$$



Always report if p-value is one or 2 sided

2.) Non-zero Nulls

$$H_0: \mu = a, \quad H_A: \mu \neq a$$

$$t^* = \frac{\bar{y} - a}{s/\sqrt{n}} \sim t_{n-1}$$

3.) "Significance"

XKCD

p-value ≤ 0.05 is not magical, it's stupid.

4.) When reporting p-values, always either report the sample size or the confidence intervals, or both.

Chapter 3: Assumptions of 2-sample t-tools

• Assumptions in order of importance:

- 1.) Independence between observational units
- 2.) Equal Variance Assumption (if assuming equal variances)
- 3.) Normality (no skew / outliers)

- Each of these assumptions can be checked
- Solutions exist for each problem.

1.) Independence:

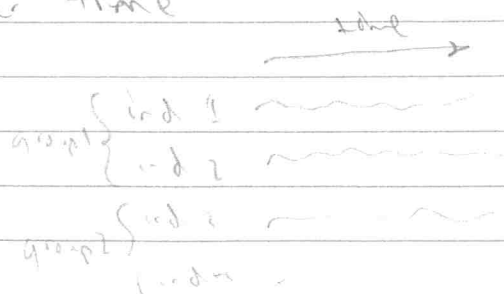
↑ Knowing the value of one observation does not tell you any information on value of second

• Example of violation:

- measure expression levels of a gene in two groups (disease/control)
- individual samples were collected by two technicians
- technician 1 tends to produce samples with higher levels
- cluster effects

• Example of violation:

- measure expression levels on the same individuals over time



• to detect:

- 1.) think carefully about how data were collected
 - ↑ were different responses measured on same subject?
 - ↑ were data collected in groups?
 - ↑ were groups treated differently, unrelated to treatment?
- 2.) Residual plots (chapt. 8)

• Issues:

Recall $\bar{X} - \bar{Y} \sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$

- If samples are not independent, the variance formula is too small
- You have less info than you think
- So p-values are too small
- Confidence intervals are too narrow

• Extreme Example:

You have 4 people, measure each one 50 times

Data:	group 1	Ind 1	0	0	0	...	0
		Ind 2	1	1	1	...	1
	group 2	Ind 3	2	2	2	...	2
		Ind 4	3	3	3	...	3

↑ sample is just size 4, but you think it is 200!

Correct variance estimates: $\frac{\sigma_1^2}{2} = \frac{\sigma_2^2}{2} = \frac{0.25}{2} = 0.125$

Assumed variance estimates $\frac{\sigma_1^2}{100} = \frac{\sigma_2^2}{100} = \frac{0.25}{100} = 0.0025$

• Solutions:

- use more sophisticated stats
- ANOVA (ch 12, 13) for cluster effects
- Longitudinal analysis (ch 15) for serial effects

2.) Unequal Variances:

- Sometimes people assume $\sigma_1^2 = \sigma_2^2$
↑ easier to generalize to more complicated methods.

$$\bar{X} - \bar{Y} \sim N(\mu_1 - \mu_2, \sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right))$$

Estimate σ^2 with $\frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1 + n_2 - 2}$

• To detect:

- make side-by-side boxplots
- only worry if extreme

• Issues:

- Not really estimating σ , so variance estimate of $\bar{X} - \bar{Y}$ is also wrong.
- only really a problem when n_1 very different from n_2 and variances are very different

• Solution:

- assume unequal variances
- could transform data.

3.) Normality (skew):

- t-test assume X_{1j}, X_{2j} and Y_{1j}, Y_{2j} are normal,
or that the sample size is "large enough"

- More skew \Rightarrow larger sample size needed

• Issues:

- Confidence intervals do not have 95% coverage
if hard to predict if conservative or liberal
- p-values do not have correct interpretation

• To detect:

- Make quantile-quantile plots

(slides)

• Solutions:

- transformation
- Do permutation test

4.) Normality (outliers)

• To detect

- histograms, box plots, qq-plots

• Issues:

- results are vulnerable (remove point, change results)

• Solution:

- Run same analysis with and without outliers
- If answers are same, do not remove
- If answers differ, either
 - 1.) Run a robust analysis (ch4)
 - 2.) Report both analyses.

• Transformations:

- many transformations of data are possible
 \sqrt{y} , $\arcsin(y)$, $\log(y)$
- \log is usually the only one used

• When to use \log :

- 1.) values are positive and
- 2.) larger mean \Rightarrow larger variance and
- 3.) data are skewed

\uparrow
 \downarrow makes such data symmetric w/ equal variance

Rainfall Slides

\downarrow Interpreting log-transformations in causal models

• Interpreting log-transformations in observational study model.

$$w_i = \log(Y_i) \Rightarrow Y_i = \exp\{w_i\}$$

$$w_i = \mu + \varepsilon_i$$

$$z_i = \log(X_i) \Rightarrow X_i = \exp\{z_i\}$$

$$z_i = \mu + \delta + \xi_i$$

$$\delta = \text{Average}(z) - \text{Average}(w) = \text{Average}(\log(X)) - \text{Average}(\log(Y))$$

↓ want to interpret on X, Y scale

- Note $\log(\text{Average}(X)) \neq \text{Average}(\log(X))$

E.g.) $X = 2, 4$

$$\log_2(X) = 1, 2, \quad \text{Average}(\log_2(X)) = 1.5$$

$$\text{Average}(X) = 3, \quad \log_2(\text{Average}(X)) = \log_2(3) \approx 1.6$$

- So not true that $E[X] = e^\delta E[Y]$

- But is true that $\text{Median}(X) = e^\delta \text{Median}(Y)$
if $\log(X)$ or $\log(Y)$ or both symmetric.

e^δ = multiplicative difference in medians

ideal: $e^\delta = \frac{e^{\text{Ave}(\log(X))}}{e^{\text{Ave}(\log(Y))}} \Rightarrow e^\delta = \text{Ave}(X) / \text{Ave}(Y)$

$e^\delta = \frac{e^{\text{Med}(\log(X))}}{e^{\text{Med}(\log(Y))}} \Rightarrow e^\delta = \text{Med}(X) / \text{Med}(Y)$

$$\delta = \text{Average}(\log(x)) - \text{Average}(\log(y))$$

$$= \text{Median}(\log(x)) - \text{Median}(\log(y))$$

if distributions of
 $\log(x)$ and $\log(y)$ are
Symmetric.

$$= \log(\text{Median}(x)) - \log(\text{Median}(y))$$

$$= \log \left(\frac{\text{Median}(x)}{\text{Median}(y)} \right)$$

$$\Rightarrow e^\delta = \frac{\text{Median}(x)}{\text{Median}(y)}$$

$$\Rightarrow \text{Median}(x) = e^\delta \text{Median}(y)$$