

Stats Group: Last Lecture

Objective:

- High level overview of course, drill principles

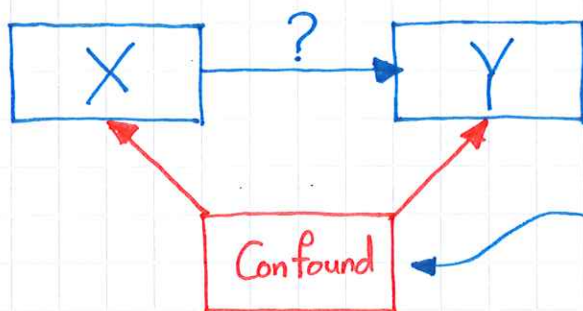
Experiments

Why perform an experiment?

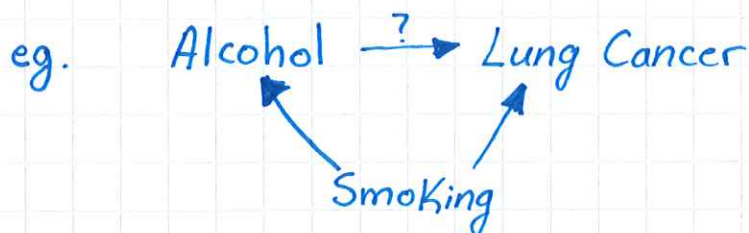
- We have questions and want answers!

What form do our questions take on?

Does treatment (X) cause outcome (Y)?



Confounds can cause us to draw the wrong conclusions!



How does a researcher control for confounds?

- By randomly assigning treatment to experimental units

↑
Spread confounds
between groups

↑
Control assignment of
treatment to group

Given a question, how does an experiment help?

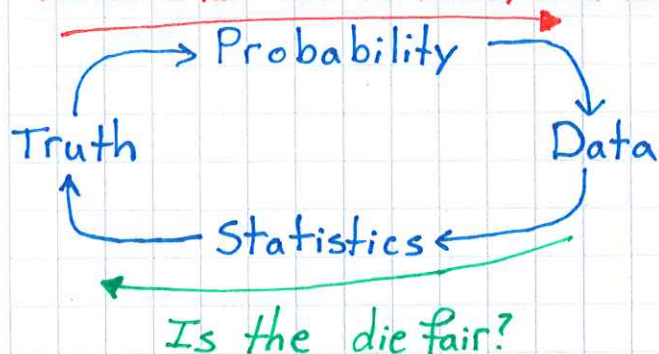
- Collect data from the natural world.
- Ask Mother Nature!

Statistics

How can we answer our question given data?

- Build a probability (mathematical) model!

"You role a fair die 20 times, what is..."

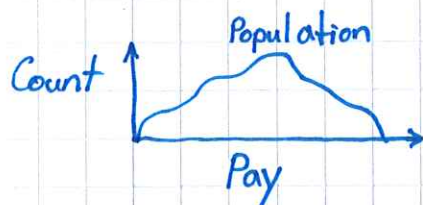


Statistics produces truths from data!

eg. Hypothesis: MSc are paid more than BSc.

The ask: $H_0: Y_{MSc} = Y_{BSc}$

How do you answer a question like that?



sample

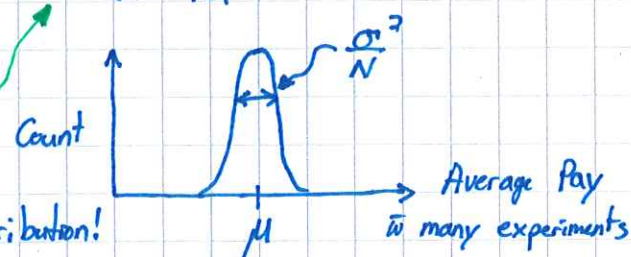
$X_1 = 20\,000$
 $X_2 = 30\,000$
!
 $X_N = 100\,000$

Average

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

Can we say anything about \bar{X} ?

Yes! As $n \rightarrow \infty$, $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$
(# samples) for any population distribution!



We have a probability now! We can ask questions like:
"What is the probability $Y_{MSc} = \$100\,000$?"

$$p\text{-value} = P(Y_{MSc} > \$100\,000)$$

Notice: We are building models of the mean

Quickly: What is the difference between a statistic and a parameter?

- Parameter is a value in our model (μ)
- Statistic is an estimate of the parameter (\bar{X})

T-test

eg. MSc are paid more than BSc. How do we test this?

- Ask 20 MSc their pay, ask 20 BSc their pay
- Compute the averages: \bar{X}_{MSc} , \bar{X}_{BSc}
- Run a t-test!

$$H_0: \bar{X}_{MSc} = \bar{X}_{BSc}$$

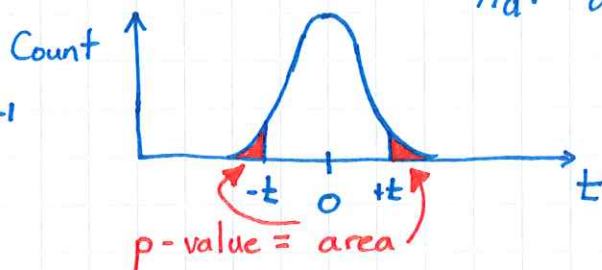
$$H_a: \bar{X}_{MSc} \neq \bar{X}_{BSc}$$

$$d = \bar{X}_{MSc} - \bar{X}_{BSc} \rightarrow H_0: d = 0$$

$$H_a: d \neq 0$$

$$t = \frac{d - 0}{SE(d)} \rightarrow T_{n-1}$$

$$SE(d) = \frac{\sigma}{\sqrt{n}}$$



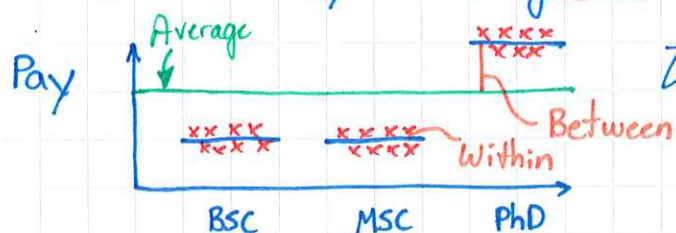
Note: Slightly more complicated because we don't have σ , we have $\hat{\sigma}$

ANOVA

eg. PhDs, MSCs, and BSCs are all paid the same.

Why can't we just run multiple tests?

- We get multiple comparisons.
- The probability of finding a difference when one isn't really there grows with every test.



We now ask a different question:

Is the variance between the groups different from the variance within the groups?



What is $\sigma^2_{\text{between}}$ relative to σ^2_{within} ?

$$F = \frac{\sigma^2_{\text{between}}}{\sigma^2_{\text{within}}} \sim F_{K-1, n-K}$$

$$\sigma^2_{\text{between}} = MS_{\text{between}} = \frac{1}{K-1} \sum (\bar{Y}_K - \bar{Y})^2$$

$$\sigma^2_{\text{within}} = MS_{\text{within}} = \frac{1}{n-K} \sum (Y_{ik} - \bar{Y}_K)^2$$

In General,

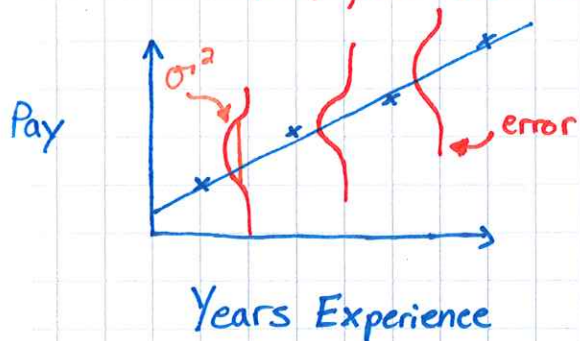
$$\sigma^2 = MS = \frac{SS}{df}$$

General Linear Model (I jumped the gun!)

Regression

ex. Pay increases with work experience. How do I test this?

- Do I run an ANOVA with 0, 1, 5, 10, 20, 30 years experience?
- No, you do regression!



$$\text{Pay} = 10\,000 + 15\,000 \cdot (\text{Years of Experience})$$

$$Y = \beta_0 + \beta_1 X + N(0, \sigma^2)$$

Parameters: β_0, β_1 Statistics: $\hat{\beta}_0, \hat{\beta}_1$

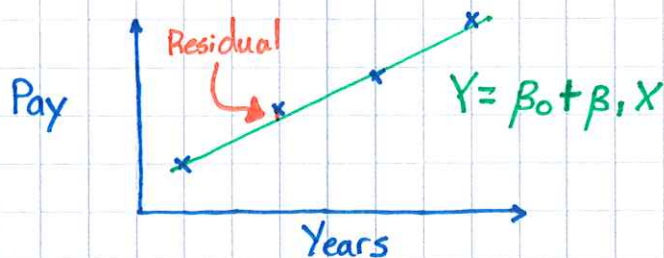
We built a model of the mean! Can ask questions like: $\beta_0 \stackrel{?}{=} 0$

How do we know if this is a good model?

Is $Y = \beta_0 + \beta_1 X$ (Model 1)
better than $Y = \beta_0$ (Model 2)

Notice: They are nested!

$$\text{Model 1} = \text{Model 2} + \beta_1 X$$



We can compare their residuals like ANOVA!

$$F = \frac{\text{Additional Variance Explained}}{\text{Unexplained Variance}}$$

$$F = \frac{\left(\frac{SS_{\text{model 2}} - SS_{\text{model 1}}}{df_1 - df_2} \right)}{MS_{\text{model 1}}}$$

Annotations:

- Additional Vari Sum of Squares (points to the numerator)
- Additional predictors (points to the denominator of the fraction)
- Variance in new model (points to the entire F-statistic)
- MS model 1 (points to the denominator of the F-statistic)

Wait, maybe t-tests and ANOVA ain't so different...

General Linear Model

It turns out that t-tests, ANOVA, and regression are the same thing...

We have two basic principles: means and residuals

Means: All our models are predicting averages

Residuals: What our model cannot explain is variance

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + N(0, \sigma^2)$$

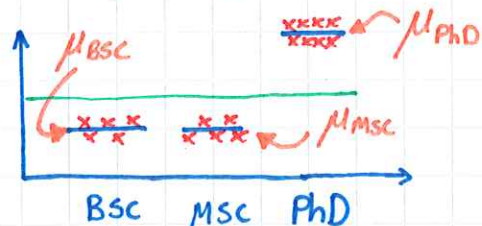
outcome

treatment

Unexplained

- error, noise, sampling.

eg. PhD, MSc, BSc are all payed the same.



Can we make this a regression?

$$Y = \beta_0 + \beta_1 MSc + \beta_2 PhD$$

$$MSC_{Model} = \begin{cases} 1, & \text{if } MSC \\ 0, & \text{otherwise} \end{cases}$$

$$PhD = \begin{cases} 1, & \text{if PhD} \\ 0, & \text{otherwise} \end{cases}$$

<u>Candidate</u>	<u>Indicator</u>	
	<u>MSC</u>	<u>PhD</u>
BSC	0	0
MSC	1	0
PhD	0	1

$$Y = \beta_0$$

$$Y = \beta_0 + \beta_1$$

$$Y = \beta_0 + \beta_2$$

$$\mu_{\text{Bsc}} = \beta_0$$

$$\mu_{MSC} = \beta_0 + \beta_1$$

$$\mu_{PhD} = \beta_0 + \beta_2$$

They're all the same!

More GLM

ex. Can a linear model explain this data?



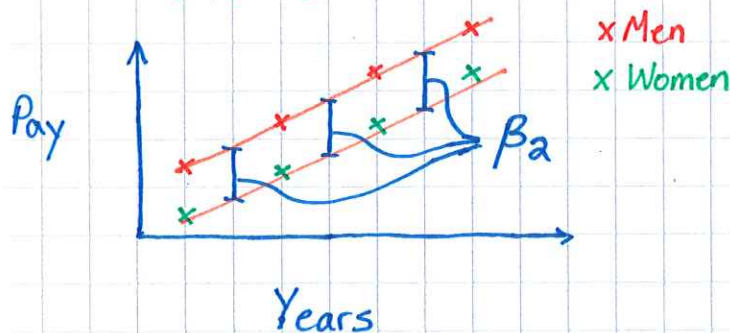
Well of course!

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2$$

The model is linear in the parameters not the treatment!

ex. What if my treatment is categorical and continuous?

eg. Pay increase with work experience is different for men and women.



$$Y = \beta_0 + \beta_1 \text{Years} + \beta_2 \cdot \text{Men}$$

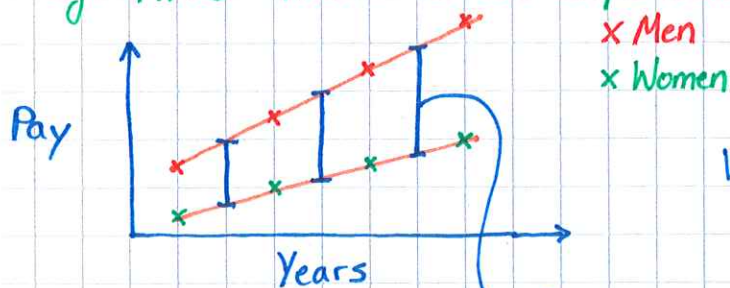
~~$$\text{Men: } Y = \beta_0 + \beta_1 \text{Years}$$~~

~~$$\text{Women: } Y = (\beta_0 + \beta_2) + \beta_1 \text{Years}$$~~

$$\text{Men: } Y = (\beta_0 + \beta_2) + \beta_1 \text{Years}$$

$$\text{Women: } Y = \beta_0 + \beta_1 \text{Years}$$

eg. AND there is an experience effect!



$$Y = \beta_0 + \beta_1 \text{Years} + \beta_2 \text{Men} + \beta_3 \text{Men Years}$$

$$\text{Women: } Y = \beta_0 + \beta_1 \text{Years}$$

$$\text{Men: } Y = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \text{Years}$$

Increases over time!
 $\beta_2 + \beta_3 \text{Years}$

A Model of Everything?

What does it model?

What if Y isn't continuous?

What if we have uninteresting effects?

Means

Generalized Linear Models

Mixed Effects Models

t-test
(means)

ANOVA
(residuals)

(Really, they're on)

Generalized Linear Models

$$Y = X\beta + N(0, \sigma^2)$$

(Y non-continuous)

Generalized Linear model

$$g(Y) = X\beta + N(0, \sigma^2)$$

(Uninteresting Parameters)

Mixed Effects Models

$$g(Y) = X\beta + ZU + N(0, \sigma^2)$$

I hope after 10 weeks you understand up to here

- For assumptions, see online
- For understanding the models, I hope this has helped