# Introduction to Modern Data Analysis

15

PART 3



# A Very Quick Tool Discussion

## Things to think about when you select analysis tools

- A. Capability: What is their functionality + performance – do they have all the techniques, do they have the processing power
- B. Integration: How do they connect to other parts of your pipeline
- C. User-Experience: What is the user experience like what background/level of expertise do you need to operate this tool, how easy is it to use this tool?
- D. Cost short and long term

# Tools for statistical analysis (I)

Python modules Enterprise (aka \$\$\$\$) Specialized Commercial Software (SAS, SPSS)

Other GUI – Excel, PowerBI\*

Other niche – e.g. Julia

R packages

# Tools for statistical analysis (II)

RULE OF THUMB 1: A t-test is a t-test is a t-test.

RULE OF THUMB 2: Do NOT implement any statistical analysis technique by hand (unless for fun/better understanding).

R will 99.99% guaranteed have any statistical technique you want for free

Python probably will have most as well

So tool choice basically comes down to how you want to prioritize/optimize A, B, C and D

## PowerBI and Statistical Analysis

- PowerBI can be used for basic descriptive statistics (e.g. mean, max/min)
- PowerBI has some DAX functions that can be used to carry out some of the calculations required for some types of inferential statistics (e.g. confidence interval for a mean)
- However, this DAX functionality is quite limited! Even Excel is arguably better in terms of its statistical analysis functionality.
- Better option use another tool to generate statistical results, then import and visualize in PowerBI
- Even better option embed R code right into PowerBI

# Statistics in the Modern Age

...

...







#### Present-Day Statistics

- At the moment, statistics is having a bit of a tough time!
- Could this be a growth opportunity for the discipline?
- Perhaps a great time for the democratization of statistics. Desktop data science!

#### Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.



From: Science Isn't Broken - It's just a hell of a lot harder than we give it credit for. (Christie Aschwanden, 2015)

# Does this mean we throw out statistics?

#### No!

- Statistics provides three very specific types of insight or knowledge that we often want, (that other techniques don't do as well):
  - Value (point or interval) estimates
  - Comparison of values, checking conditions (via hypothesis testing)
  - Understanding of associations (A and B are connected in some way)

#### **Physical Sciences**

Relatively straight-forward measurements (weight)

Relatively objective measurements

Measure the same object multiple times: Any difference is error

Independence of objects and measurements

Measurement and environment conditions strictly controlled

(sometimes) simple variable relationships

#### **Social Sciences**

Measurements can be subjective (e.g. happiness?)

High variability across objects of the same type (e.g. humans are different)

Variability is not (just) due to measurement error

Easier for dependencies to arise

Less control over environmental/experiment al conditions

Complex relationships

#### **Applied Data Analysis**

All the problems associated with social sciences, plus...

Data is not originally collected with analysis in mind

Typical no control over data collection conditions

Data is typically observational

Independent/dependent variables not easily controlled

Data collection does not necessarily occur with high quality control

## Alternate Analysis Techniques

The emergence of alternate analysis techniques means that conventional statistics is 'under attack' from many different directions:

**Bayesian Statistics** 

**Causal Analysis (Bayesian Networks)** 

**Machine Learning** 

As a result – much of what you learned about statistics in school may need a bit of an upgrade

# Strategies

#### How to deal with all of this?

- Tools are getting more intelligent
  - Example with categorical data, if you have small counts you often need to switch techniques e.g. if you have less than a certain number of values in a chi squared test, you need to do things slightly differently
  - Good tools will automatically make this adjustment for you
- Have an expert statistician on staff or on call BUT – one versed in modern methods and who is not afraid to get their hands dirty!

# Some Useful Statistical Concepts





## Is It a Sample or Population?

Are we dealing with a sample or a population?

Given a particular dataset, the answer to this question depends entirely on your inference goal

**Examples:** 

- I want to understand this year's finances, and I currently have full data on this year's finances no inference required you have a population
- I want to compare this year's finances with last year's finances no inference required – you have a population
- I want to use this year's data to derive fundamental laws about financial transactions— definitely inference, definitely a sample

# We know that different samples will give us different values



There are 22% blue stars in the sample



There are 33% blue stars in the sample

## A population of samples? Very meta!



Each sample in the population of samples has a particular percentage of blue stars





### Confidence Intervals

- A **confidence interval** calculation tells you that the true value (e.g. true population mean) falls within the range calculated.
- It is based off the value calculated from the sample, some additional measures taken from the sample, the size of the sample and some additional assumptions.
- The biggest one of these assumptions is that we have a good sample!
- Which is to say, a representative sample.

(https://towardsdatascience.com/statistics-are-you-bayesian-or-frequentist-4943f953f21b)



#### Confidence Levels

- We can further ask: how surprised would we be if it turned out this is a bad sample?
- OR We can flip this around and say how confident are we that this is a good sample?
- This depends on a number of things- like what?
- We could say we are 95% confident (or 80%, or 90%) depending on what we picked for our confidence level when calculating the interval.
- The confidence interval is connected to the confidence level

(https://towardsdatascience.com/statistics-are-you-bayesian-or-frequentist-4943f953f21b)





#### Bar plot with 95% confidence intervals



### Statistical Tests

- What is a statistical (significance) test? A method for drawing rigorous inferences from data.
- Usually involves a comparison or check (e.g. is mean A different from mean B, is this distribution normal, is proportion C less than 0.2)

## Hypotheses and Significance

- We come up with a null hypothesis and an alternate hypothesis:
  - Null hypothesis these two traits, A and B are evenly split in the population
  - Alternate hypothesis there is more of trait A than B in the population
- BUT we know that any given sample will almost certainly not be a perfect reflection of the actual population!
- How do we deal with this?
- **Significance** is another strategy to deal with this issue.

#### A Pie Analogy

To help you get a better feel for the reasoning behind significance, consider this story about my brother, pie and hypothesis testing...



Suppose we end up with Sample 4 – it does show a difference from our null hypothesis. How surprised would we be if it turned out this difference was just due to bad luck (we could have ended up with sample 1, after all)

## Significance: Definition and Interpretation

If the sample is big (and the difference is big) and we have reason to believe the sample is representative... we would be pretty surprised if the difference was due to a bad luck sample.

Significance quantifies this intuition.

Technical Definition of Significance Level: the probability of getting results *at least as* extreme as the ones you observed, given that the null hypothesis is correct.

Example: There's an 20% chance that we would get a difference in proportions this big or larger if the difference between the populations was actually 50%.

Interpretation: How surprised will you be if the null hypothesis turns out to be true, under these circumstances.

Significant is not the same as substantial.

## Significance: Sample vs Population

As the sample gets larger and larger it gets closer and closer to being equivalent to the full population\*

With large sample sizes, even very small differences become significant – we would be very surprised if the difference was just by chance.

For this reason: Significant is not the same as Substantial

Once we get to the full population, any size of difference is significant (aka real!). There's no way the difference is just by chance.

IMPORTANT: no need for inferential statistics in the case of populations – also no need for significance, or tests!

#### Tons of Statistical Tests

For a comprehensive table of statistical tests:

Choosing a Statistical Test (Summary and Analysis of Extension Program Evaluation in R, Salvatore S. Mangiafico, 2016 https://rcompanion.org/ handbook/D\_03.html



# Statistical Tests in Applied Situations

Problem : Traditional Vanilla Statistical Tests (e.g. T-Test, Z-Test) *really only work well in scientific experiment contexts.* Specifically, where you have:

Developed hypotheses in advance of collecting data

Selected sample size based on power you think you need, again *in advance* 

Used a sampling method that will likely provide a representative sample of data

Controlled the conditions of the sample selection as well as other experimental elements (environment, hypotheses being tested) to make it easy to draw strong conclusions

Good reason to believe that the system represented by the data conforms to the other required assumptions of the test (e.g. normal distribution, independence)

### Chi Square Test for Independence

- Null Hypothesis: The distribution of the outcome is independent of the groups (no relation between variables)
- Alternative Hypothesis: there is a difference in the distribution of responses to the outcome variable among the comparison groups (relationship between variables)
- Let's pick an *alpha value* of 0.05.
- This tells us that we would be '95% of maximumsurprised" if the null hypothesis was actually true in this situation when our test statistic said it was false.
- we'll reject the null hypothesis if the *p-value*\* is LESS THAN this value. We want our 'minimum surprise level' to be 95%.
- Remember small differences can be significant with large sample sizes

\*the probability of obtaining test results at least as extreme as the results actually observed, under the assumption that the null hypothesis is correct

## Tests Example 1: Chi Square

Size of Institutions Relative to Year

Year	Large	Medium	Small
2017	2285	537	418
2018	2274	491	448
2019	2379	570	455
2020	2385	552	435
2021	1416	331	261

## But suppose we considered this a sample...

(with the population being...?)

Year	Large	Medium	Small
2017	2285	537	418
2018	2274	491	448
2019	2379	570	455
2020	2385	552	435
2021	1416	331	261

#### Code to test for independence

- isize <- read.csv("year\_instituttion\_size.csv") #LOAD DATA</li>
- isize\_sum <- table(isize\$Competition.Year,isize\$Institution.Size..ENG) #FORMAT DATA
- chisq <- chisq.test(isize\_sum) #RUN TEST</li>
- Pearson's Chi-squared test
- data: isize\_sum
- X-squared = 4.8042, df = 8, p-value = 0.7783

The p-value is much larger than 0.05, so there is no dependence between variables.

In other words, year is not associated with size of institution. We cannot use information about year to help us guess proportion of large, small or medium institutions in a given year.

#### Some additional points about statistical tests

Non-Normal Data: If you know you have nonnormal data, turn to non-parametric tests. If you don't, you can do a test to see if your data is normally distributed, and then proceed

Non-parametric tests: These make no assumptions about the distribution of the data so you can use them whenever you want. Unfortunately they can be less definitive and informative.

Transforming your data: Another option is to transform you data from not normal to normal, and then use parametric tests. Alternative To Using Statistical Tests for Hypothesis Testing: Statistical Modelling Statistical modelling tends to get complicated very fast.

BUT the reality is that real world data is usually complicated.

This is particularly the case with OBSERVATIONAL DATA.

What this means is that with observational data you don't really have an easy middle ground – you either:

- do descriptive statistics OR
- you jump to consulting with a professional statistician OR
- get used to doing creating complicated statistical models yourselves (preferably with the help of a tool like R)

#### Two Variable Linear Model

- There are straightforward techniques to let you fit a linear (straight line) model to your data.
- Let your relationship coefficient help you decide if this is a valid (or useful) model



#### Confidence + Prediction Intervals



blue line = linear model grey band = 95% confidence interval red lines = 95% prediction interval

- The confidence and prediction intervals give us a sense of how certain we are in our model.
- Confidence interval tells us what 'y' we can expect on average for a given x.
- Prediction interval tells us the range we can expect 'y' to fall in, for a specific instance of 'x'.

Chi Squared Test - Log Linear Models Chi Squared: Tests to see if there is an interaction between categorical variables

Instead you can use a log linear model to detect and further investigate the associations and interactions between variables



# Case Study



Gender differences in grant and personnel award funding rates at the Canadian Institutes of Health Research based on research content area: A retrospective analysis

Karen E. A. Burns<sup>1,2,3</sup>\*, Sharon E. Straus<sup>1,4,5</sup>, Kuan Liu<sup>5</sup>, Leena Rizvi<sup>2</sup>, Gordon Guyatt<sup>3</sup>

 The Li Ka Shing Knowledge Institute, St. Michael's Hospital, Toronto, Ontario, Canada, 2 The Interdepartmental Division of Critical Care, Department of Medicine, University of Toronto, Toronto, Ontario, Canada, 3 Health Research Methods, Evidence and Impact, McMaster University, Hamilton, Ontario, Canada, 4 Division of General Internal Medicine, Department of Medicine, University of Toronto, Toronto, Ontario, Canada, 5 Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada

\* burnsk@smh.ca

#### Background

Although women at all career stages are more likely to leave academia than men, earlycareer women are a particularly high-risk group. Research supports that women are less likely than men to receive research funding; however, whether funding success rates vary based on research content is unknown. We addressed gender differences in funding success rates for applications directed to one or more of 13 institutes, representing research communities, over a 15-year period.

#### Methods and findings

We retrospectively reviewed 55 700 media bit 207 uncereal award applications submitted to the canadian Institutes of Health Research. We analyzed application success whe according to gender and the primary institute selected by applicants, pooled gender differences in success rates using random effects models, and fitted Poisson regression models to assess the effects of gender, time, and institute. We noted variable success rates among trant applications directed to selected institutes and declining success rates over time. Women securited 21 1% and 44.7% of grant and personnel award applications of peo-

tively. In the pooled estimate, women had significantly lower grant success (risk ratio [RR] 0.89, 95% confidence interval [CI] 0.84–0.94; p < 0.001; absolute difference 3.2%) compared with men, with substantial heterogeneity (I<sup>2</sup> = 58%). Compared with men, women who directed grants to the Institutes of Cancer Research (RR 0.86, 95% CI 0.78–0.96), Circulatory and Respiratory Health (RR 0.74, 95% CI 0.66–0.84), Health Services and Policy Research (RR 0.78, 95% CI 0.68–0.90), and Musculoskeletal Health and Arthritis (RR 0.80,

Other potentially relevant statistical modelling case studies Multilevel Modelling in Repeated Measures of the Quality of Finnish School Life

Performing Learning Analytics via Generalised Mixed-Effects Trees

Racial Bias and LSI-R Assessments in Probation Sentencing and Outcomes

Intersectionality in quantitative research: A systematic review of its emergence and applications of theory and methods

For more stats and ML, check out our Data Science Report Series

 <u>https://www.data-action-</u> lab.com/data-science-report-series/

## Extra Material

#### Proportion Tests – Example 2

 Proportion of French Language Applicants – Does this differ across years?

Year Fr	ench
2017	99
2018	104
2019	107
2020	103
2021	80

#### Proportion Tests – Example 2

- There is clearly a difference in the numbers.
- This is a real difference.
- Is it a substantial difference? We need to ask the SMEs!

Year Fro	ench
2017	99
2018	104
2019	107
2020	103
2021	80

#### Proportion Tests – Example 2

- What if we treat this data as a sample.
- (... of what?)

Year Fro	ench
2017	99
2018	104
2019	107
2020	103
2021	80

• What if we treat this data as a sample (... of what?)

Year	French	Year	French
2017	99	2017	0.20
2018	104	2018	0.21
2019	107	2019	0.22
2020	103	2020	0.21
2021	80	2021	0.16

Null hypothesis – all categories have a proportion (probability) of 0.20

Result of Proportions Test

chisq.test(lang\_tab\_fr)

Chi-squared test for given probabilities

data: lang\_tab\_fr

X-squared = 4.7181, df = 4, p-value = 0.3175

**Conclusion: We can't reject the null hypothesis**.

### What would it look like if there was more of a difference?

#### Consider this fake data

Year	Esperanto
2017	10
2018	15
2019	600
2020	31
2021	200

Pearson's Chi-squared test data: fake\_data X-squared = 1237.7, df = 4, p-value < 2.2e-16 Conclusion: Reject the null hypothesis!

#### Proportion Tests (Example 3)

- Percent French language applications across all years: 0.03%
- We want to know is this number of applications representative, relative to the prevalence of French in other contexts (e.g. in Canada)?
- What additional data can we use to try to answer this question?
- In Canada, people who say their mother tongue is French: 21% of population.
- A possible approach? Could we ask: is there a significant difference between percentage of French Language Applicants and people who say French is their mother tongue?
- (Does it even make sense to ask this question?)

#### **One-proportion test (Z-test/Chi-Square)**

- Possible Strategy One: We are dealing with two populations! We don't need to do statistics. So YES there is an obvious difference.
- Possibly Strategy Two: We can deal with this as a sample of some larger population
- Suppose we decide to compare it to a 'constant' but that constant value is based on the 21% statistic
- Several difficult questions arise:
  - does it make sense to compare the university applicant population with Canadian citizens?
  - Is our sample from some population like "People who could have applied in French but chose not to apply for some reason?
- Let's go ahead and see anyway we can use a one-proportion test to see if the value is significantly different from 21%

#### R-Code

```
English French

14744 493

test <- prop.test(

+ x = 493, # number of successes

+ n = 15237, # total number of trials (14744 + 493)

+ p = 0.21

)
```

#### Result

1-sample proportions test with continuity correction

- data: 493 out of 15237, null probability 0.21
- X-squared = 2897.3, df = 1, p-value < 2.2e-16
- alternative hypothesis: true p is not equal to 0.21
- 95 percent confidence interval:
- 0.02963027 0.03531913
- sample estimates:
- p
- 0.03235545

Some additional points about statistical tests Non-Normal Data: If you know you have non-normal data, turn to non-parametric tests. If you don't, you can do a test to see if your data is normally distributed, and then proceed

Non-parametric tests: These make no assumptions about the distribution of the data so you can use them whenever you want. Unfortunately they can be less definitive and informative.

Transforming your data: Another option is to transform you data from not normal to normal, and then use parametric tests.

#### Ordinal Data

#### PAIN MEASUREMENT SCALE



# Ordinal Data Best Practices

Try to avoid it! But that's probably not realistic...

Could just treat as categorical, then do proportions

Does it make sense to take the mean of ordinal data values? It's certainly possible to do, but...

Some argument in social science literature has been made (on the basis of some evidence) that it is can be acceptable to treat ordinal data as if it were numeric for the purposes of analysis, under some circumstances.

The more fine-grained, the closer you are to numerical (but... heaping!)

Non-parametric tests. Or use a Lickert scale/Likert item approach.

# Side note: Meta-Analysis

What if you only have the results of statistical analysis, and not the raw data?

E.g. suppose you just have: sample size, mean and confidence interval

This is a common situation in meta-analysis studies.

There are packages (e.g. meta) that are designed to carry out comparisons when you just have this info.

However, if you have raw data that you wish to compare to a result, you might consider just treating the result as a baseline hypothesis for a plain-jane statistical test.