ECE1724H S2: Empirical Software Engineering Introduction of Quantitative Study



Course Project

- Interim Report: Annotated outline (due 3/11)
- Final Report: (due 4/11)



The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO

•

When to use Survey Research?

- To evaluate the frequency of some characteristic across a population
 - E.g. how many companies use agile methods?
- To evaluate the severity of some condition that occurs in a population
 - E.g. what's the average cost overrun of software projects?
- To identify factors that influence a characteristic or condition
 - E.g. What factors cause companies to adopt new ASE tools?

What type of question are you asking?					
	<u>Exploratory</u> Correlation				
	→Existence:	→Relationship			
	→Description & Classification	Do occurrences of X correlate with occurrences of Y?			
	What is x like? What are its properties?	→Causality Causal			
	 ✤ How can it be categorized? ✤ How can we measure it? ✤ What are its components? →Descriptive-Comparative ✤ How does X differ from Y? 	 ♥ Does X cause Y? ♥ Does X prevent Y? ♥ What causes X? ♥ What effect does X have on Y? →Causality-Comparative 	hip		
Base-rate	 →Frequency and Distribution ♥ How often does X occur? ♥ What is an average amount of X? 	 Does X cause more Y than does Z? Is X better at preventing Y than is Z? Does X cause more Y than does Z under one condition but not others? 			
Just a reminder	 →Descriptive-Process ♦ How does X normally work? ♥ By what process does X happen? ♥ What are the steps as X evolves? 	→Design Des Solution → Des Solution → Des Des Des Des Des Des Des Des	ign		

Agenda for today

• Paper reading presentation for Surveys



- Descriptive statistics
- Hypothesis testing







- Henne, B., Harbach, M., & Smith, M. (2013). <u>Location privacy</u> <u>revisited: factors of privacy decisionsPreview the document</u>. Extended Abstracts on Human Factors in Computing Systems.
- Shklovski, I., Mainwaring, S. D., Skúladóttir, H. H., & Borgthorsson, H. (2014). Leakiness and creepiness in app space: perceptions of privacy and mobile app usePreview the document, Proceedings of the ACM conference on Human factors in computing systems (pp. 2347-2356): ACM.

Research Designs

• Step 1 to **design** empirical research:

adopt a general (and guiding) approach

- Three main approaches:
 - Qualitative
 - Quantitative
 - Mixed Methods

Chapter 1 - Creswell, J. W., & Creswell, J. D. (2017 Research design: Qualitative, quantitative, and mixed methods approaches. Sage publications.

Just a reminder...

- Qual. vs Quant.:
 - Not dichotomies
 - Rather, different ends on a continuum

Qualitative vs Quantitative

- Often:
 - Words (qual) vs Numbers (quant)
 - Open-ended questions (qual interview questions) vs
 - Closed-ended questions (quant hypotheses)
 - Inductive (qual) vs Deductive (quant)
- More completely, differences in:
 - philosophical assumptions
 - strategies of inquiry
 - research methods



Planning Checklist

- 💰 Pick a topic
- Identify the research question(s)
- Scheck the literature
- Identify your philosophical stance
- Identify appropriate theories
- Schoose the method(s)
- Design the study
 - Unit of analysis?
 - Target population?
 - Sampling technique?
 - Data collection techniques?
 - Metrics for key variables?
 - Handle confounding factors

- Critically appraise the design for threats to validity
- Get IRB approval
 - Informed consent?
 - Benefits outweigh risks?
- Recruit subjects / field sites
- Conduct the study
- Analyze the data
- Write up the results and <u>publish</u> them
- o Iterate

Just a reminder...

Agenda for today

- Paper reading presentation for Surveys
- Descriptive statistics
- Hypothesis testing

Measurement





"Measuring programming progress by lines of code is like measuring aircraft building progress by weight."



☐ isaacs / github Output Unwatch Contributions <> Code (!) Issues 1.3k 11 Pull requests 2 Projects 🛄 Wiki Actions U Security Insights Feb Mar Arx May Contribution graph can be harmful to contributors #627 () Open mxsasha opened this issue on Apr 1, 2016 · 197 comments Summary of pull requests, issues opened, and commits. Lear 0 mxsasha commented on Apr 1, 2016 ···· ··· Contributions in the last year A common well-being issue in open-source communities is the tendency of people to over-commit. Many contributors care deeply, at the risk of saying yes too often harming their well-being. Open-source communities are especially at risk, because 235 total many contributors work next to a full-time job. Feb 8, 2015 - Feb 8, 2016 The contribution graph and the statistics on it, prominent on everyone's profile, basically rewards people for doing work on as many different days as possible, generally making more contributions, and making contributions on multiple days in a row

without a break.

Contributing graphs considered ł https://www.hanselman.com/

Stepping away from our work regularly is not only important to uphold high quality work, but also to maintain our well-being. For example, I personally do not generally work in the weekends. That's completely healthy. I take a step back from work and spend time on other things. But in the contribution graph it means I can never make a long streak, even though I do work virtually every day except weekends. So the graph motivates me to work in my weekends as well, and not take breaks. And

Context of Measurement

- The meaning of measurements will vary depending on whether they derive from observation or experiment.
 - Observation: sampling difficulties \rightarrow bias
 - Experiment: a lack of generalizability
- The available measurements are not immediately connected with the phenomena of interest
 - Leading/Lagging -- difficult to use in managing day-to-day operations
- Side-effects -- the numbers are "massaged" or the work process altered in order to produce the "right" results
- Good measurements are **actionable**

(Ch 6) Statistical Methods and Measurement . from F. Shull et al. (eds.), Guide to Advanced Empirical Software Engineering. Springer 2008

Creating Effective Metrics

- Precise, Reliable, Valid
- Well-motivated -- a metric must provide at least a partial answer to a specific question, a question which itself is aimed at some particular research or management goal.

- "What is the expected amount of time for a specific class of defects to go from the initial reported state to the Repaired state?"

- "What percent of all customer reported defects are in the Repaired state within two days of being first reported?"

• a single metric is usually not sufficient

(Ch 6) Statistical Methods and Measurement . from F. Shull et al. (eds.), Guide to Advanced Empirical Software Engineering. Springer 2008

Defining a Metric

- Simple
 - counts (e.g., number of units shipped this year)
 - dimensional measures (e.g., this year's support costs, in dollars)
 - categories (e.g., problem types)
 - rankings (e.g., problem severity)
- Compound
 - ratios (e.g., defects per thousand units)
 - rates (time-based ratios such as number of problem reports per month)
 - proportions or percentages (e.g., proportion of customers responding "very satisfied" to a survey question)
 - linear algebraic combinations (e.g., mean repair cost the sum of all repair costs divided by the total number of repairs)
 - **indices** (dimensionless measures typically based on a sum and then standardized to some baseline value). Whereas

How will you measure things?

Туре	Meaning	Example	Admissible Operations
Nominal Scale	Unordered categories	Gender, Political preferences, Place of residence	=
Ordinal Scale	Ranking of objects into ordered categories (intervals between the values are not necessarily of the same size)	Satisfaction, Happiness, grades	=, <, >
Interval Scale	Differences between points on the scale are meaningful (equal intervals)	Celsius, Fahrenheit Temperature, IQ (intelligence scale), SAT scores	=, <, >, difference, mean
Ratio Scale	Ratios between points on the scale are meaningful (ordered, equal intervals with a zero point)	weight, height, sales figures, ruler measurements, number of children	=, <, >, difference, mean, ratio

Levels of Measurement

Offers:	Nominal	Ordinal	Interval	Ratio
The sequence of variables is established	_	Yes	Yes	Yes
Mode	Yes	Yes	Yes	Yes
Median	_	Yes	Yes	Yes
Mean	_	-	Yes	Yes
Difference between variables can be evaluated	_	-	Yes	Yes
Addition and Subtraction of variables	_	-	Yes	Yes
Multiplication and Division of variables	_	-	-	Yes
Absolute zero	_	-	-	Yes

https://www.questionpro.com/blog/nominal-ordinal-interval-ratio/

Quantitative interpretation



Fig. 10.1 Three steps in quantitative interpretation

C. Wohlin et al., Experimentation in Software Engineering, Springer-Verlag Berlin Heidelberg 2012



Descriptive Statistics

Scatter plot





Descriptive Statistics

- Measure of central tendency (mean, median, mode)
- Measure of dispersion (variance, standard deviation, coefficient of variant)
- Measure of dependency (regression, covariance, correlation coefficient,...)

Measures of Central Tendency



most *representative or typical* of all values in a group "average"

MODE

- most frequent data point
- mode exists as a data point
- unaffected by extreme values
- useful for qualitative data
- may have more than 1 value

MEDIAN

- value that divides ranked data points into halves: 50% larger than it, 50% smaller
- may not exist as a data point in the set
- influenced by position of items, but not their values

 Σ_x

MEAN

- most stable measure
- affected by extreme values
- may not exist as a data point in the set

https://www.slideshare.net/tol edo98/measures-of-centraltendency-and-variability



Measures of Dispersion



Measures of Dependency/association

COVARIANCE



Near Zero Large Negative Large Positive Covariance Covariance Covariance r = 0.4 r = 0 r = -0.4 **Positive Correlation** No correlation Negative

Correlation

- Measure of the relation between two variables:
 - -1 variables perfect inverses (negative correlation)
 - 0 no correlation at all
 - +1 variables are perfectly correlated (they appear on a straight line with positive slope)
- Pearson's r
 - Computed as: $r_{xy} = \frac{\sum (x_i \bar{x})(y_i \bar{y})}{(n-1)s_x s_y}$
 - •x and y are the sample means
 - $\bullet s_x$ and s_y are the sample standard deviations $\bullet n$ is the sample size
 - Assumes variables are interval or ratio scale
 - Is independent of the measurement unit



The Datasaurus Dozon



https://www.r-bloggers.com/2020/10/predicting-classmembership-for-the-tidytuesday-datasaurus-dozen/

Correlation Cousality



http://www.tylervigen.com/spurious-correlations

Divorce rate in Maine correlates with Per capita consumption of margarine

 \equiv

Correlation: 99.26% (r=0.992558)



Data sources: National Vital Statistics Reports and U.S. Department of Agriculture



http://xkcd.com/552/

- For causation
 - Provide a theory (from domain knowledge, independent of data)
 - Show correlation
 - Demonstrate ability to predict new cases (replicate/validate)

The association between early career informal mentorship in academic collaborations and junior author performance

We also find that increasing the momentian of formal mentancies and not



Quantitative interpretation



Fig. 10.1 Three steps in quantitative interpretation

C. Wohlin et al., Experimentation in Software Engineering, Springer-Verlag Berlin Heidelberg 2012

Removal of outliers





Quantitative interpretation



Fig. 10.1 Three steps in quantitative interpretation

C. Wohlin et al., Experimentation in Software Engineering, Springer-Verlag Berlin Heidelberg 2012

Two flavors of Hypothesis Testing

- Parametric tests -- operate on data from a probability distribution, such as the normal distribution or the t -distribution
- Non-parametric tests: distribution free

Normal Distribution



Source: wikipedia - http://en.wikipedia.org/wiki/Image:Standard_deviation_diagram.png

Checking your data is normal

- →Draw a Histogram
- →Compute the mean and standard deviation
- →Superimpose the expected normal curve over the histogram





Fig. 2 Two very different samples with the same mean and standard deviation

(Ch 6) Statistical Methods and Measurement . from F. Shull et al. (eds.), Guide to Advanced Empirical Software Engineering. Springer 2008

How will you measure things?



Туре	Meaning	Example	Appropriate Stat. Test	
Nominal Scale	Unordered categories	Gender, Political preferences, Place of residence		
Ordinal Scale	Ranking of objects into ordered categories (intervals between the values are not necessarily of the same size)	Satisfaction, Happiness, grades	Non-parametric tests	
Interval Scale	Differences between points on the scale are meaningful (equal intervals)	Celsius, Fahrenheit Temperature, IQ (intelligence scale), SAT scores	Parametric tests Non-parametric tests	
Ratio Scale	Ratios between points on the scale are meaningful (ordered, equal intervals with a zero point)	weight, height, sales figures, ruler measurements, number of children		

Hypothesis Testing

- Set up some hypotheses
 - Null hypothesis (H₀) asserts that a relationship does not hold
 - In many cases, this is the same as saying there is no difference in the the means of two different treatment groups
 - Alternative hypotheses (H₁, ...) each asserts a specific relationship
 - Type I error: A false positive (rejecting H₀ when it's true)
 - Type II error: A false negative (accepting H₀ when it's false)
- For the statistical tests
 - P value (we calculate this) probability that a relationship observed in the sample happened by chance
 - Alpha level (selected a priori) a threshold for p at which we will accept that a relationship did not happen by chance (typically 0.01 or 0.05)
 - This allows us to fix the probability of a type I error in advance
 - If $p < \alpha$, we say the result was significant



https://www.simplypsychology.org /type_I_and_type_II_errors.html



P-value



Statistical Power





Effect size

- increase sample size
- 80%
- Estimate how many subjects they need based on estimates of effect size and power

A/B Testing

A/B testing

- Controlled randomized experiment with two variants, A and B, which are the control and treatment.
- One group of users given A (current system); another random group presented with B; outcomes compared.
- Often used in web or GUI-based applications, especially to test advertising or GUI element placement or design decisions.

 3,662 visitors and get 378 conversions: conversion rate = 378/3,663 = 10.3% ± 1.0 %

10% Conversion Rate Standard Error by Sample Size

Hypothesis testing

- H0: Variation B is NOT a meaningful improvement over Variation A.
- H1: Variation B will convert at a higher rate for our overall population than Variation A will.

1.Test says Variation B is better & Variation B is actually better
2.Test says Variation B is better & Variation B is not actually better (type I error)
3.Test says Variation B is not better & Variation B is actually better (type II error)
4.Test says Variation B is not better & Variation B is not actually better

If Variation B performs better in our sample. How do we know whether or not that improvement will translate to the overall population? How do we avoid making a type I error?

P-value

Hypothesis testing

- H0: Variation B is NOT a meaningful improvement over Variation A.
- H1: Variation B will convert at a higher rate for our overall population than Variation A will.

1.Test says Variation B is better & Variation B is actually better
 2.Test says Variation B is better & Variation B is not actually better (type I error)
 3.Test says Variation B is not better & Variation B is actually better (type II error)
 4.Test says Variation B is not better & Variation B is not actually better

Statistical significance cannot be used as a stopping point

Usability: A/B testing

- However, it cannot..
 - Tell you why
 - Let you test drastic redesigns of your website or app.
 - Tell you if you're solving the right/wrong problem.

Which Statistical Test?

https://seeing-theory.brown.edu/

Visualizations

A Share 😈

Since 2014 I've tried to illustrate various statistical concepts using interactive visualizations. This page contains an updated overview of my visualizations.

Maximum Likelihood

An interactive post covering various aspects of maximum likelihood estimation.

Cohen's d An interactive app to visualize and understand standardized effect sizes.

Statistical Power and Significance Testing

Equivalence and Non-Inferiority Testing

Explore how superiority,

non-inferiority, and

equivalence testing

relates to a confidence

An interactive version of the traditional Type I and II error illustration.

Confidence Intervals

An interactive simulation of confidence intervals

Bayesian Inference

An interactive illustration of prior, likelihood, and posterior.

Correlations

Interactive scatterplot that let's you visualize correlations of various magnitudes.

P-value distribution

Explore the expected distribution of p-values under varying alternative hypothesises.

https://rpsychologist.com/viz

References

- (Ch 10) Analysis and Interpretation . from C. Wohlin et al., Experimentation in Software Engineering, Springer-Verlag Berlin Heidelberg 2012

- (Ch 6) Statistical Methods and Measurement . from F. Shull et al. (eds.), Guide to Advanced Empirical Software Engineering. Springer 2008
- (Ch 6) Hypothesis Testing. from MacKenzie. Human-Computer Interaction. Elsevier 2013