L project for equity in data science









What is the average size classroom?





The average classroom is 3 students per class.





The average classroom is 4 students per class.







Both are correct.







Teacher Perspective 1 + 3

Student Perspective 1 + 3 + 3 + 3 2 2 2 2



3 + 3 + 3 + 5 + 5 + 5 + 5 + 5 = 35 35/9 = 4



WeAllCount.com Heather Krause, PStat @datassist





Feminist Data Analysis

Step by step processes to avoid racism, sexism, homophobia and more in data and analysis.



Feminist Data Analysis:





Locus of Power

World View

Unit of Analysis

Assumed 'Normal'





"I refuse to accept the idea that we can simply shoehorn women into a global economy that is exploiting them and then celebrate it as women's economic empowerment."

- Winnie Byanyima, Executive Director, **Oxfam International**



Sources of bias can be identified in each step of the data life cycle.



Funding



Motivation





Project Design

& Distribution



Funding





The source of funding impacts the actual outcomes of RCTS.

Column

The dark side of education research: widespread bias

Johns Hopkins study finds that insider research shows 70 percent more benefits to students than independent research



Proof Points

Column by JILL BARSHAY



March 18, 2019



The Library of Missing Datasets

















Questions for Funding Step:

Who is funding this research?

What research is not being funded?

How would I design, conduct, analyze this data if I was independently wealthy?

with humans that will contribute data to this project?

- What are the demographic profile of the funders compared



Motivation





Why is this data project being done?



Tension between explicit purpose (understand if this works) and implicit purpose (have a good report for the annual general meeting).



NYC & Rats



Question for Motivation Step:

What is the stated purpose of this data project? What are the other hidden purposes behind this data? Who stands to benefit from this data project? Who stands to lose from this data project? Who wants to know what? How could we align the explicit and implicit purpose?



Project Design





Project Design is the phase where the WHY becomes the HOW

Critical step in data equity











Sample design based on definitions whose definitions?







Study Up



RCTs: The Gold Standard ... of What?



Unbiased estimate of average treatment effect of the population.

*What we don't find is that a small number of people had their income increase \$1000 and a big group of people had their average income decrease by \$100.

We found that the population average income increased \$100*

Questions for Project Design Step:

Who is deciding the methods we're using in this data project?
When are results needed - how will this impact different power dynamics?
How good is our method at capturing unintended effects?
Who is benefitting the most from this design - who is this easiest for?
What level(s) are you looking at with your project design?
Is the method a good fit with the community of the research? (Do they both believe in linear time?)



Data Collection & Sourcing



Data Biographies at the bare minimum must accompany each dataset you are using:



Why

How

Where



Who is the head of YOUR household?





Group #1

Unmarried Mother | Her Father | Her Son

Married Woman | Her 3 Children Husband Lives Abroad



Group #4 Two Unmarried Women | Their Kids | One Brother



Married Woman | Her 3 Children Husband Lives Abroad



Group #5

Disabled Father | Married Couple | Two Kids



AVERAGE INCREASE IN HOUSEHOLD SAVINGS











Measuring social constructs (Demographics) Who is constructing and how. Which ones are fluid?

And how are you going to use this?





Additive vs Intersectional Analysis **Odds of Getting Referral (Additive)**



Black Men

Black Woman


Additive vs Intersectional Analysis Odds of Getting Referral (Intersectional)





White Women

Black Men

Black Woman



Who collects the data matters. Data collected by higher status enumerators results in different results.





What counts as work?

The following six things do not count as work:

- 1) The cleaning, decoration and maintenance of the dwelling occupied by the household, including small repairs of a kind usually carried out by tenants as well as owners;
- 2 The cleaning, servicing and repair of household durables or other goods, including vehicles used for household purposes;
- **3** The preparation and servicing of meal;
- 4 The care, training and instruction of children;
- 5 The care of sick, infirm or old people;
- **6** The transcription of members of the household or their goods.



Questions for the Data Collection Step:

What are the core concepts we're measuring?

What are three different ways to measure each concept?

Who is collecting the data?

Who owns the data?

Who is defining success?

If collecting from beneficiaries, can they receive the treatment without consenting? If not, is it informed consent?

What is the burden being places on people contributing data?



Analysis







Methods Matter ALOT.



Chance of having a low birthweight baby is







Chance of having a low birthweight baby



Ethnic Group A



Ethnic Group B



Chance of having a low birthweight baby when taking into account **means community of residence is between**







Chance of having a low birthweight baby when taking into account community of residence and state of residence is between







Amount related to:





Chance

Individual





State

Community



Chance of having a low birthweight baby when taking into consideration community and state







Ethnic 15% Group B





21%

Contextual Variables BIRTH WEIGHT = Person **Person+Ethnicity Person+Ethnicity+Community Person+Ethnicity+Community+State** Person+ Ethnicity+Community*State+C(Community)+C(State)





Your world view determines how you model success. Productivity increases over time



PRODUCTIVITY (LITERS PER COW)



Your world view determines how you model success. Hours of farm work increases over time



		•••••
	 	••••••
	 	•••••
Time 2	Time 3	

HOUR PER DAY ON DAIRY



Your world view determines how you model success. Productivity controlling for increase in work time



PRODUCTIVITY HOLDING HOURS STEADY



Your world view determines how you model success. Productivity controlling for increase in work time





PRODUCTIVITY HOLDING HOURS STEADY



Questions for the Analysis Step:

How are we defining success?

questions?

on your world view?

cultural setting you're in?

- Are you using biased data to set up prediction that will be biased?
- What assumptions are your statistical methods making?
- What other statistical methods could you use to ask the same
- What are you controlling for in your data that might be based only
- What are you not controlling for that might be important to the



Interpretation



Control Group





Project Participants

Avg. Monthly Income (USD)



















Rate of students who experienced one suspension or more, by racial/ethnic group



PERCENT OF STUDENTS WHO EXPERIENCED ONE OR MORE SUSPENSIONS



Relative rate ratios comparing the rates of students who experienced one suspension or more in specific racial/ethnic groups with the rate among White students who experienced one suspension or more.



RELATIVE RATE RATIO



Comparison of two compositions: Proportion of the student group who were suspended and proportion of the group in the student population, by racial/ethnic group.



Composition of student enrollment Composition of suspended Students



The relative difference in composition between the proportion of students who experienced suspensions in each racial/ethnic group and proportion of the group in the total student population.







Ways to think about fair

group agree.

- Equal False Negative Rates: the fraction of positives which are marked negative in each group agree.
- **Equal False Positive Rates:** the fraction of negatives which are marked positive in each group agree.
- Equal Positive Predictive Values: the fraction of those marked positive which are actually positive in each
- Statistical Parity (equal positive decision rates): the fraction marked positive in each group should agree.



Questions for the Interpretation Step:

What are three alternative explanations for your results?

would their questions be?

way?

- What would your interpretation be if your results were reversed?
- How would you describe these results to a five year old? What
- How can you break out your interpretation in an intersectional
- How does uncertainty and probability affect these results



Communication & Distribution




























Data Viz "best practices" are not culturally universaly.



TIME SPENT ON DAIRY ACTIVITIES PER DAY





Lastly, consider the mediums used to distribute information. All distribution systems come with compromises.



Questions for the Communication Step:

Who gets to see the findings? Really?

Are they in English only?

Is the graphic design reflect the cultural and cognitive preferences of which power group?

Where is the locus of power in the visual and written communication?

Does the result and communication require internet access?



Sources of bias can be identified in each step of the data life cycle.



Funding



Motivation





Project Design

& Distribution





DEMYSTIFY. DEMOCRATIZE. DEMONSTRATE.



We All Count Tools

We All Count believes that the world is a little too full of people pointing out problems without offering solutions. WAC is committed to providing practical resources to help anyone who wants to make their data science more equitable.





















Thank you. WeAllCount.com Heather Krause, PStat heather@idatassist.com @datassist



