

# Case-control simulation meeting

12/12/17

Genesis of this idea

Straightforward ways to combine data sources, weights

Basic - sim made up pop

outcome / disease incidence  
covars

show that we can  
recover RR, IDR, OR

- ↑ complexity ↓
  - 1. all cases, SRS controls - show this works
  - 2. " " , probability-based sample of controls
  - 3. " " # stratified/complex sampling
  - 4. diff source pops - assumptions  
around W, modifiers, transportability

NSF grant app - potential - due Feb 28

→ creation of data resources - need to be public

→ advancement of methodology

social, econ, behav. sciences

↳ all prob using big, admin  
data, even if not  
case-control

- Jen to think about whether or not she wants to  
do this

→ will talk to Will, psychologists, sociology

Start w/ <sup>several</sup> categorical, binary, cont W's - age, race,  
- related sex, county  
- exposure - binary zip, household,  
- outcome - binary urban/rural

Catherine, Chris, Ellie - we can meet + work  
on basic structure

Samir - interpreting epi evidence  
case-control studies - characteristics  
of control group

after 1<sup>st</sup> week of Jan

Would be good to check what's been done  
on this - Catherine to do a lit search  
first

Patrick - thinking  
start w/ a known case-control study  
what if we didn't have these controls  
Build simulations off of real data  
→ simulate outcome only - use  
real covariate data

Jen wants to start from cohort + show that all  
the case-control designs

Would be good to have an applied example - use cohort  
data as base  
BioLink - publicly available datasets  
ERIC  
Framingham

Start ~~basic~~ basic - given that we don't have  
a lot of time left on grant

More complexity  
- pair matching  
- freq. "

Start w/ generic simulation

n = 1 mil or 100,000

DF = 5%-ish

How many controls - 4:1 - diminishing returns

Start w/ cors from blown up ACS + then shrink it down  
from 1 year to something manageable

## Computing

### Variability

Treat population as truth

Repeat sampling many times to characterize the distribution of point estimates + SE ests from each run

Patrick will ask Burke for grizzly bear access for Catherine

Timeline - as soon as we can.

Parameters? Assume conditional not marginal?  
OR/RR?

Report design ~~weights~~  
factor / relative efficiency

Possible extensions: measurement error in the number of controls (needed for case-control weights)