# hdPS adjustment for analyzing electronic healthcare data:

Overview, recent advances & open questions

Ehsan Karim

*ehsan.karim@ubc.ca*

Feb 19, 2019; BIRS

# Popularity of hdPS

- High-dimensional propensity score (hdPS)
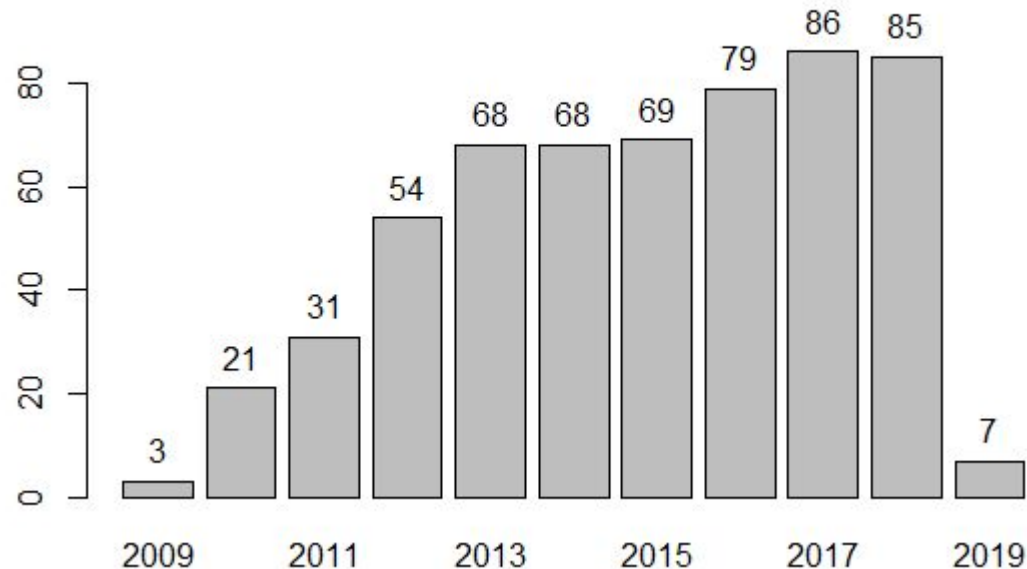- Unmeasured confounding

## Citation of Schneeweiss et al. (2009)

| Year | Citations |
|------|-----------|
| 2009 | 3 |
| 2010 | 21 |
| 2011 | 31 |
| 2012 | 54 |
| 2013 | 68 |
| 2014 | 68 |
| 2015 | 69 |
| 2016 | 79 |
| 2017 | 86 |
| 2018 | 85 |
| 2019 | 7 |

# General Idea of hdPS

Administering health care data:



Longitudinal patient records: diagnostic and procedural information

# Proxy measures of U

| Unobserved confounder | Observable proxy measurement | Coding examples |
|---|---|---|
| Very frail health | Use of oxygen canister | CPT-4 |
| Sick but not critical | Code for hypertension during a hospital stay | ICD-9, ICD-10 |
| Health-seeking behavior | Regular check-up visit; regular screening examinations | ICD-9, CPT-4, #PCP visits |
| Fairly healthy senior | Receiving the first lipid-lowering medication at age 70 years | NDC, ATC, Read |
| Chronically sick | Regular visits with specialist, hospitalization; many prescription drugs | #specialist visits, NDC, ATC |
| Outcome surveillance intensity | General markers for health care utilization intensity | #visits, #different drugs |

# Type of variables

# Type of variables

unmeasured confounder (U)

proxy of the unmeasured confounder (P)

measured confounder (L)

risk factor (R)

intervention (A)

mediating variable (M)

outcome (Y)

instrument (I)

common effect (C)

Effect of outcome (E)

# Type of variables

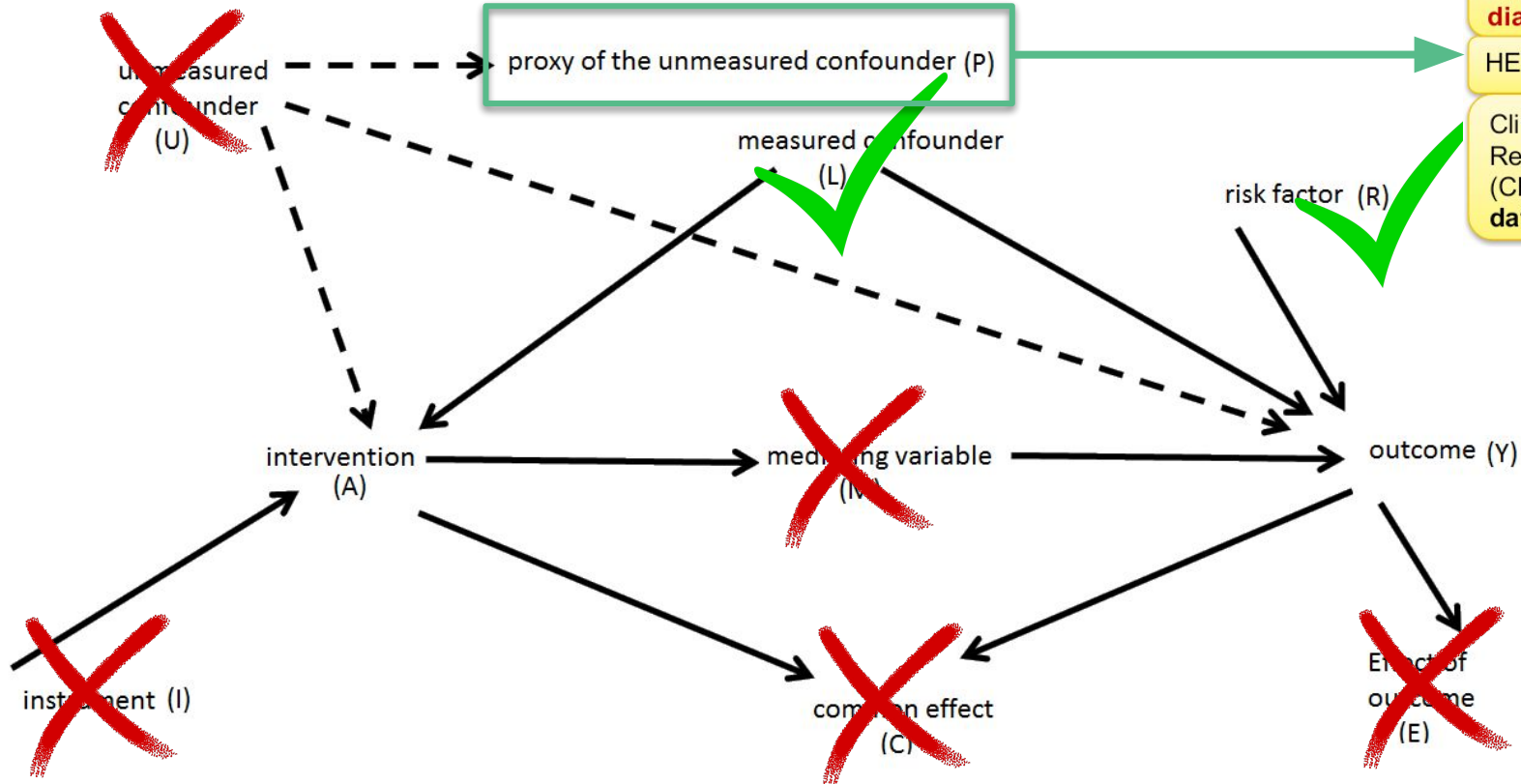Investigator specified covariates: L + R
High-dimensional covariates: P
4 data dim×200×3=2,400 binary variables



general practice data (diagnoses, referrals, immunizations, laboratory tests)

Hospital Episode Statistics (HES) diagnosis data

HES procedure data

Clinical Practice Research Datalink (CPRD) medication data

# Amount of confounding due to an unmeasured confounder

1. Assumptions:
   - $p_{u1}$ = prevalence among treated
   - $p_{u0}$ = prevalence among untreated
   - $p_{uY1}$ = prevalence among dead
   - $p_{uY0}$ = prevalence among alive
   - $RR_{uY}$ = $p_{uY1}$ / $p_{uY0}$

2. Adjusted RR:
   - $RR_{adj} = RR_{obs} \dfrac{p_{u1}(RR_{uY}-1)+1}{p_{u0}(RR_{uY}-1)+1}$

3. Amount of Bias / confounding due to u:
   - $Bias_M = \dfrac{p_{u1}(RR_{uY}-1)+1}{p_{u0}(RR_{uY}-1)+1}$

Let
Y = outcome
    (dead/alive)
A = treatment
$RR_{obs}$ = Crude RR

u = unmeasured
    confounder
    (say, healthy
    eating)
$RR_{adj}$ = ??

**Bross formula** 1966

Spurious effects from an extraneous variable
IDJ **Bross** - Journal of chronic diseases, **1966** - Elsevier
Abstract Spurious effects from an extraneous variable are a troublesome problem in many
areas of research in the biological and behavioral sciences. While investigators have
recognized intuitively that there is a relationship between the size of an effect and its chance
of being spurious, current textbooks do not contain any explicit statement of this relationship.
In this paper one such statement, the Size Rule, is developed. The application of this rule ...

# Prioritization in hdPS

1. ~~Assumptions:~~
   » $p_{u1}$ = prevalence among treated
   » $p_{u0}$ = prevalence among untreated
   » $p_{uY1}$ = prevalence among dead
   » $p_{uY0}$ = prevalence among alive
   » $RR_{uY}$ = $p_{uY1} / p_{uY0}$

2. Adjusted RR:
   » $RR_{adj} = RR_{obs} \dfrac{p_{u1}(RR_{uY}-1)+1}{p_{u0}(RR_{uY}-1)+1}$

3. Amount of Bias / confounding due to u
   » $Bias_M = \dfrac{p_{u1}(RR_{uY}-1)+1}{p_{u0}(RR_{uY}-1)+1}$

Bross formula 1966

**Spurious effects** from an extraneous variable
IDJ **Bross** - Journal of chronic diseases, **1966** - Elsevier
Abstract Spurious effects from an extraneous variable are a troublesome problem in many areas of research in the biological and behavioral sciences. While investigators have recognized intuitively that there is a relationship between the size of an effect and its chance of being spurious, current textbooks do not contain any explicit statement of this relationship. In this paper one such statement, the Size Rule, is developed. The application of this rule …

Replace u by p (binary) and calculate

Amount of Bias / confounding due to p:

» $Bias_M = \dfrac{p_{p1}(RR_{pY}-1)+1}{p_{p0}(RR_{pY}-1)+1}$

Sort based on magnitude of rank-score (descending):

» $|Log(Bias_M)|$ [biased-based, Bross 1966]

PS model based on
- L
- R
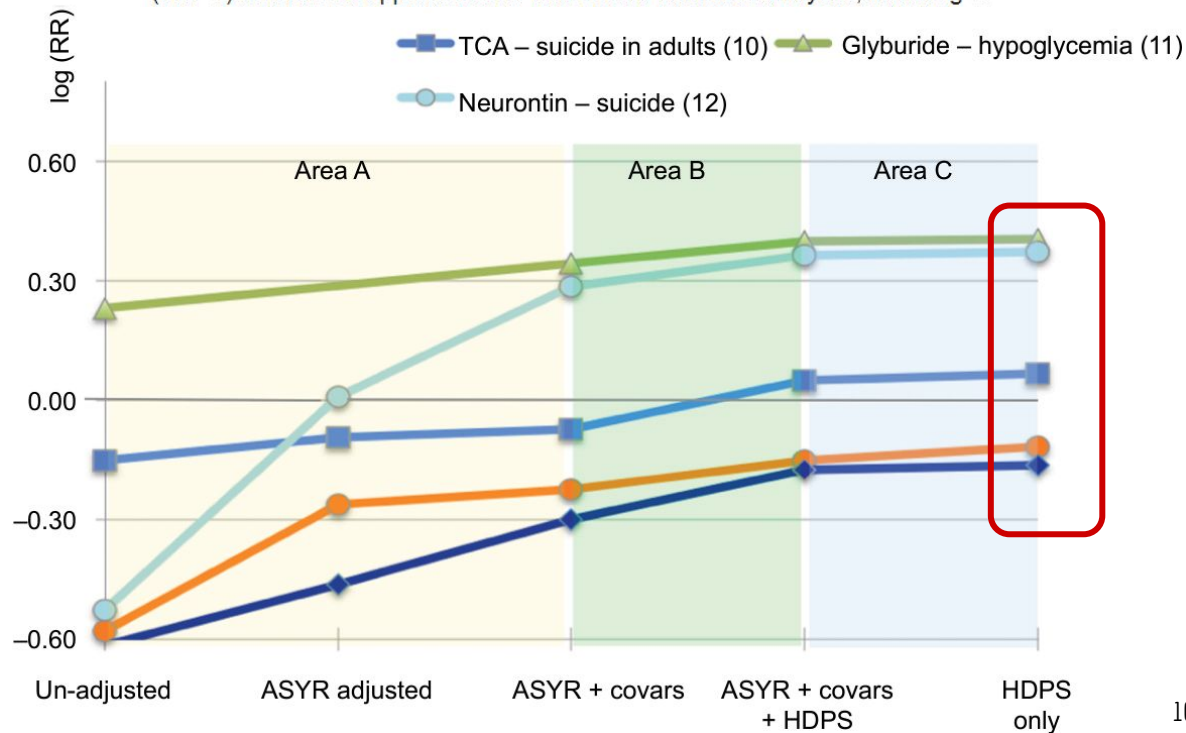- Select p variables (top 500)

9

# Performance of hdPS

**Sequential addition of covariates** vs. **change in effect estimate**

"This strongly suggests that ==even without the investigator-specifying covariates== for adjustment, the ==algorithm alone optimizes confounding adjustment=="."

10

# Limitations / Extensions

- Bivariate adjustment
  - (multivariate adjustment instead of Bross? )
  - (Collinearity?: Ridge/LASSO)
- Mis-specification
  - (double robust/TMLE, SL)
- Time-varying covariates
  - (MSM)
- IV / collider

**Regularized regression versus the high-dimensional propensity score for confounding adjustment in secondary database analyses**

JM Franklin, W Eddings, RJ Glynn... - American journal of ..., 2015 - academic.oup.com

... We present a simulation study that compares the **high-dimensional propensity score** algorithm for variable selection with approaches that utilize direct adjustment for all potential confounders via regularized regression, including **ridge** regression and lasso regression ...

[PDF] **Scalable collaborative targeted learning for large scale and high-dimensional data**

C Ju, S Gruber, SD Lendle, JM Franklin... - UC Berkeley Division of ..., 2016 - core.ac.uk

... **dimensional propensity score** (hdPS) algorithm is a method to extract information from electronic medical claims data that produces hundreds or even thousands of candidate covariates, increasing the dimension of the data dramatically. [16] In order to apply C- **TMLE** to large ...

**High-dimensional propensity score** algorithm in comparative effectiveness research with time-varying interventions

R Neugebauer, JA Schmittdiel, Z Zhu... - Statistics in ..., 2015 - Wiley Online Library

... The **high-dimensional propensity score** (hdPS) algorithm was proposed 1 for automation of confounding adjustment in problems ... of confounders 'by hand' is not practical because of the **high** dimensionality of ... t. At each time point , expert-selected covariates (listed in **Table** 1) are ...

# Extensions of hdPS

- Only ~ <mark>30%</mark> of the selected hdPS covariates were common.
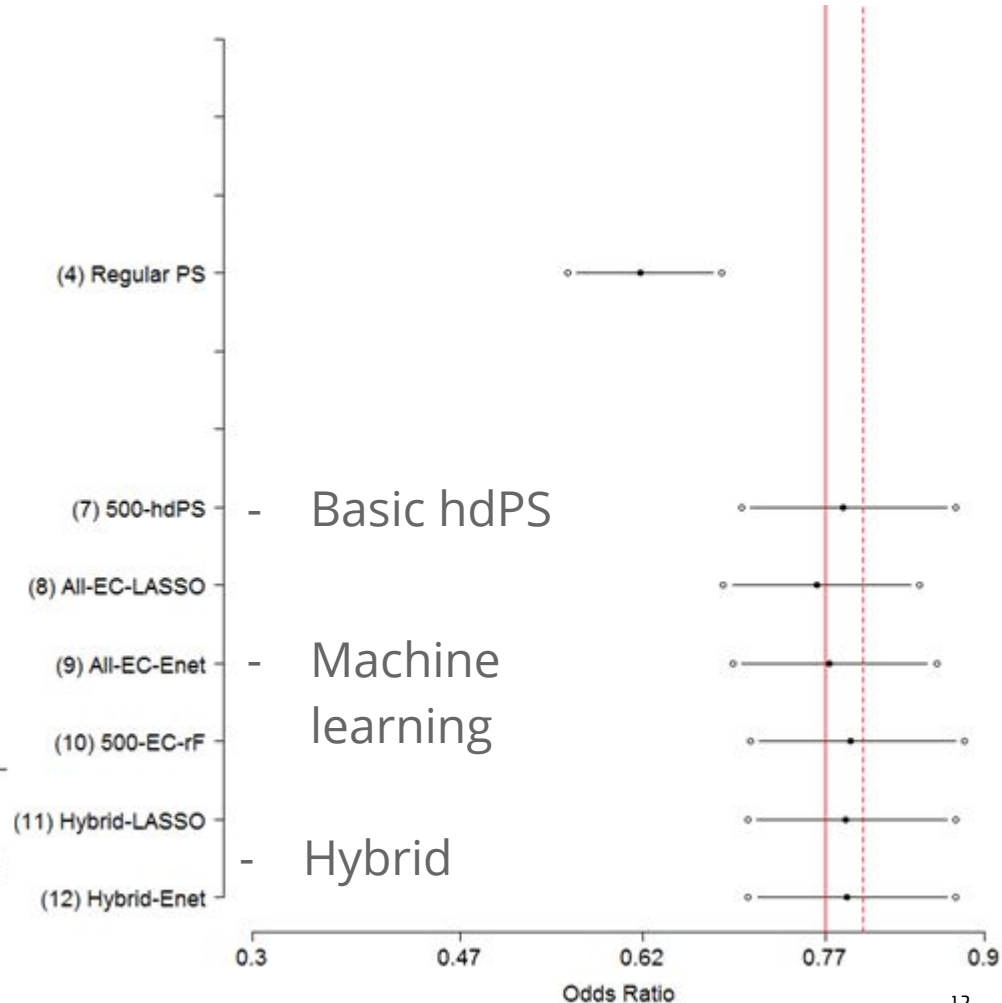- <mark>Statistical inefficiency</mark>

ORIGINAL ARTICLE

Can We Train Machine Learning Methods to Outperform the High-dimensional Propensity Score Algorithm?

Mohammad Ehsanul Karim,[a,b] Menglan Pang,[c,d] and Robert W. Platt[c,e,f]

- Basic hdPS

- Machine learning

- Hybrid

12

# Inflated SE

- "… overfitting of propensity score models can lead to inflated variance of effect estimates and therefore to estimation inaccuracy in situations where relatively many covariates are included in the propensity score model" (# of exposed vs. # of covariates)
- hdPS context

# Application of hdPS

**JAMA study (2017)**: **Serotonergic Antidepressant Use** during pregnancy vs. **Autism Spectrum Disorder** in Children

- Unadjusted:                    HR, 2.16 [95% CI, 1.64-2.86]
- Multivariable adjusted: HR, 1.59 [95% CI, 1.17-2.17]
- IPTW hdPS:                    HR, 1.61 [95% CI, 0.997-2.59] ⟶ "not associated"!!
- 1:1 hdPS matching:      HR, 1.64 [95% CI, 1.07-2.53] (sensitivity analysis 1)
- Pre-pregnancy data:      HR, 1.85 [95% CI, 1.37-2.51] (sensitivity analysis 2)

"Adjusting for too many pre-exposure covariates will lead to **collinearity and statistical inefficiency** …."

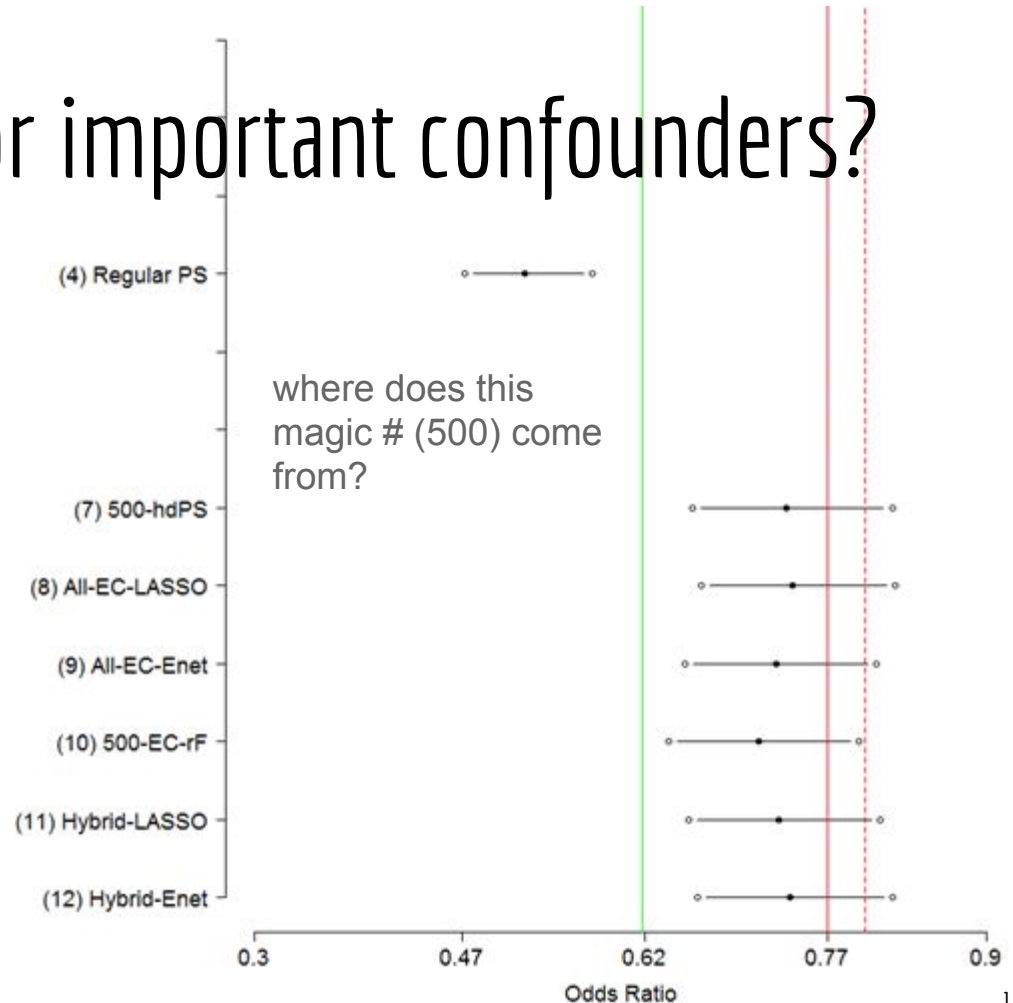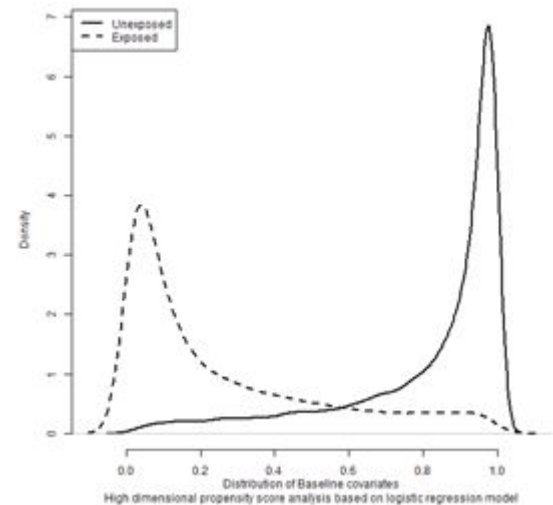14

# Collective substitute for important confounders?

- use of a hdPS "to adjust for ==500 covariates== that might **collectively contribute to confounding**"
- "association … may not be causal" (JAMA editorial)

- Most simulations based on ==plasmode==

where does this magic # (500) come from?



15

# Current practices and Open questions

Balance diagnostics

- PS analysis not reported, hdPS being main analysis!
- <mark>Deviation from PS</mark>
  - design vs. analysis stage; selective inference?
- <mark>Balance diagnostics</mark> in high-dimension
  - balance in p?
- <mark>Trimming</mark>:
  - practical/near positivity assumption violation
  - target population? bias-variance trade-off

# Thank you!

ehsan.karim@ubc.ca