

The Counterfactual χ -GAN

Amelia J Averitt, MPH MA MPhil

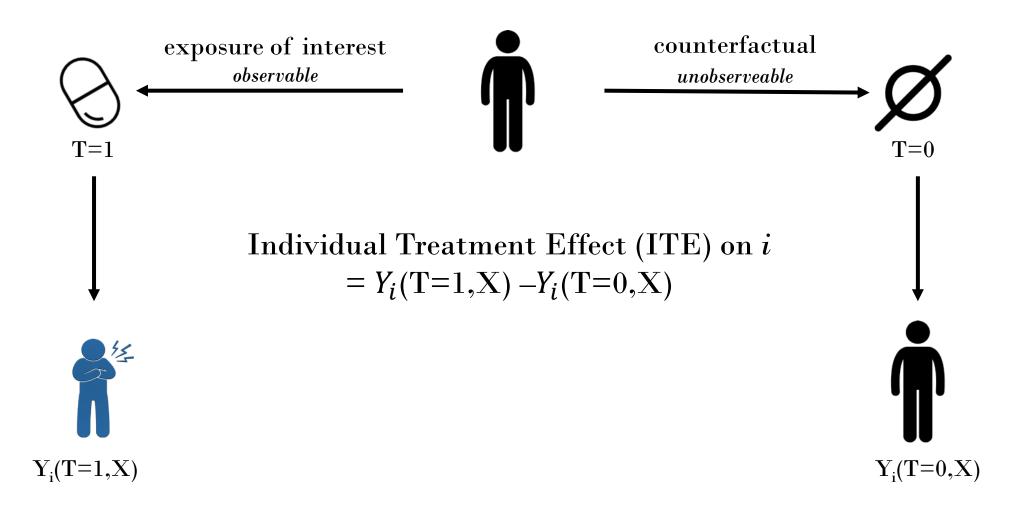


COLUMBIA UNIVERSITY BIOMEDICAL INFORMATICS

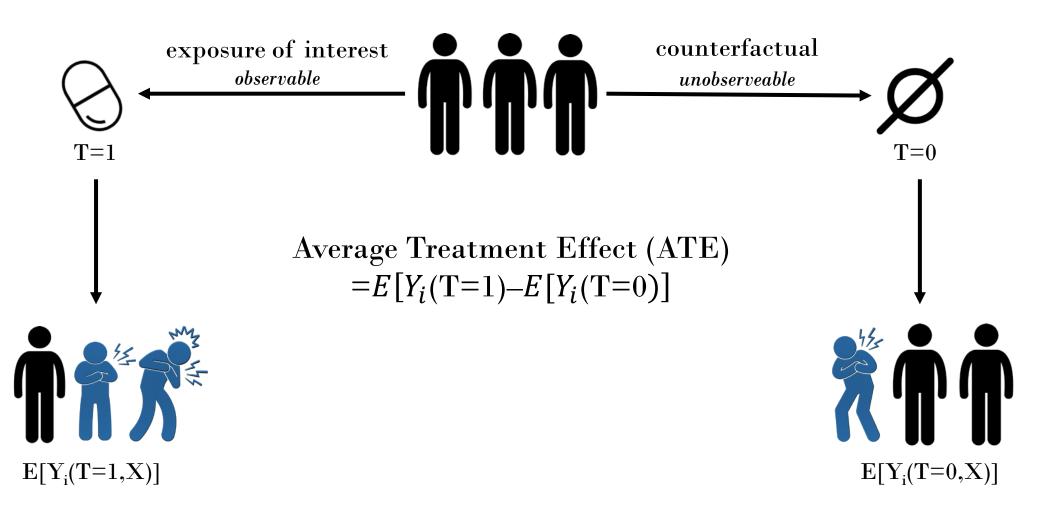


"Causal inference refers to the process of drawing a conclusion about cause and effect relationships" Vogt 2011

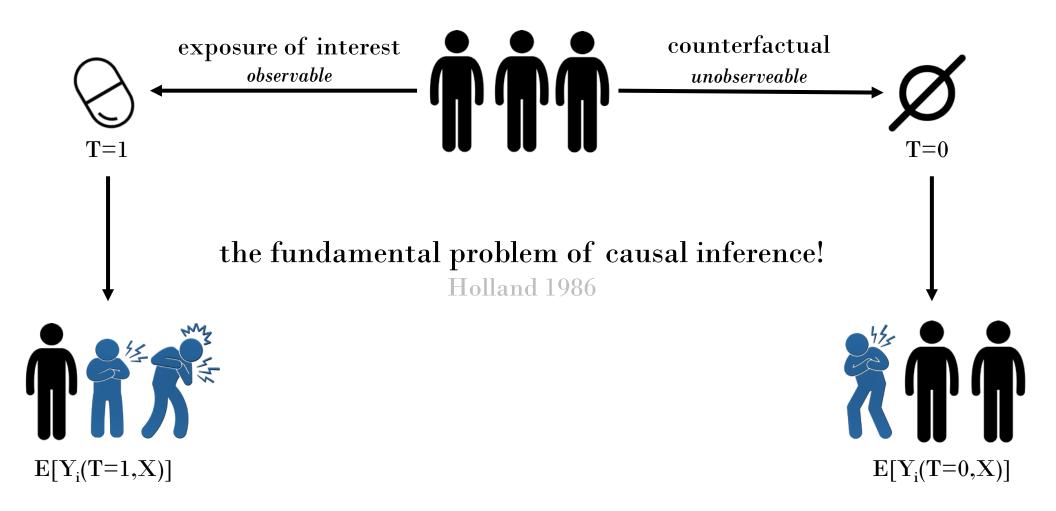




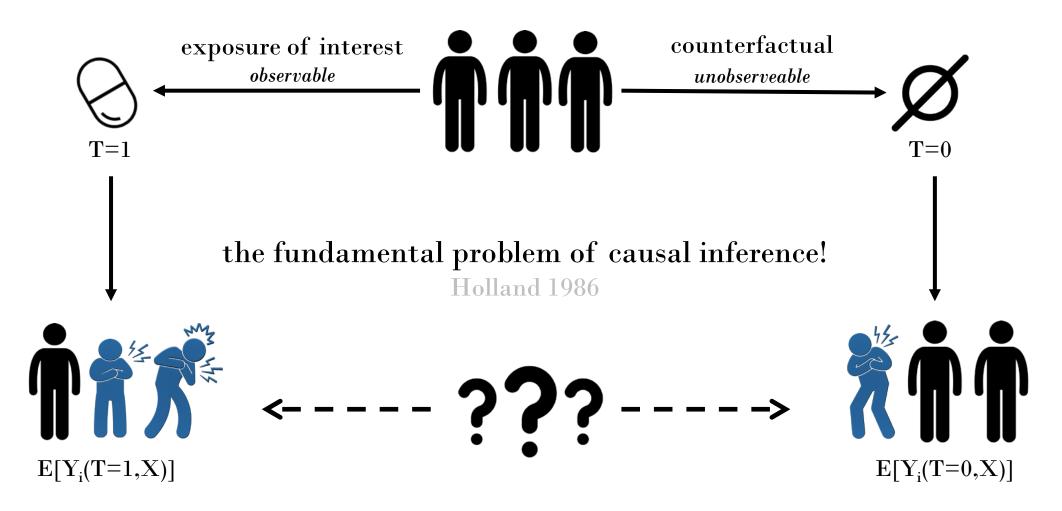














In an cohort of treatment and comparator units, $Y_i(T=1)$ and $Y_i(T=0)$ are potential outcomes in that either of these two outcomes can be potentially observed.

Contrast the mean reported outcome in each arm

$$\widehat{E}[Y|T=1] = \frac{1}{n} \sum_{i=1}^{n} Y_i(T=1, X=X_i)$$



$$\widehat{E}[Y|T=0] = \frac{1}{n} \sum_{i=1}^{n} Y_i(T=0, X=X_i)$$



Approximate the counterfactual treatment effect when two assumptions are met.





Approximate the counterfactual treatment effect when two assumptions are met.

- 1. Stable Unit Treatment Value Assumption (SUTVA) Cox 1958; Rubin 1986
 - the [potential outcome of] one unit should be unaffected by the particular assignment of treatments to the other units. Laffers 2016



Approximate the counterfactual treatment effect when two assumptions are met.

- 1. Stable Unit Treatment Value Assumption (SUTVA) Cox 1958; Rubin 1986
 - the [potential outcome of] one unit should be unaffected by the particular assignment of treatments to the other units. Laffers 2016

- 2. Strong Ignorability/ Exchangeability Rosenbaum 1983a
 - $(Y_i(1), Y_i(0)) \perp T_i \mid X_i = x \quad \forall x$



violations of these assumptions bias causal estimates

enforcing or correcting these assumptions improves causal estimates







by enforcing feature balance

Imbens 2009; Morgan 2014; Ho 2007

$$\tilde{F}(X|T=1) = \tilde{F}(X|T=0)$$
 where $\tilde{F}(\cdot)$ is the empirical distribution



by enforcing feature balance

Imbens 2009; Morgan 2014; Ho 2007

$$\tilde{F}(X|T=1) = \tilde{F}(X|T=0)$$
 where $\tilde{F}(\cdot)$ is the empirical distribution

experimental data from an RCT!



by enforcing feature balance

Imbens 2009; Morgan 2014; Ho 2007

$$\tilde{F}(X|T=1) = \tilde{F}(X|T=0)$$
 where $\tilde{F}(\cdot)$ is the empirical distribution

experimental data from an RCT!





by enforcing feature balance

Imbens 2009; Morgan 2014; Ho 2007

$$\tilde{F}(X|T=1) = \tilde{F}(X|T=0)$$
 where $\tilde{F}(\cdot)$ is the empirical distribution

experimental data from an RCT!

expensive, unethical, not representative, poor generalizability, & narrow scope

World Medical Association 1997; Rothman 2000; DiMasi 2014; Gabler 2009; Longford 1999; Kravitz 2004; Lachin 1988; Kernan 1999





or we could use observational data

data that is passively collected without any engineering adjustments Czitrom 1997

• electronic health record (EHR) data Murdoch 2013



or we could use observational data

data that is passively collected without any engineering adjustments Czitrom 1997

• electronic health record (EHR) data Murdoch 2013

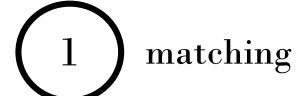
suitable for studying rare outcomes there's a lot of it Imai 2009 representative Concato 2004; Thadani 2006; Kleinberg 2011 inferences are more externally valid Steckler 2008; Rothwell 2006

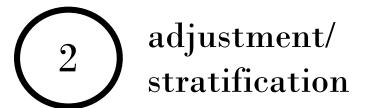


observational data is nonrandomized and requires manipulations to enforce strong ignorability



manipulations to enforce strong ignorability with observational data









manipulations to enforce strong ignorability with observational data









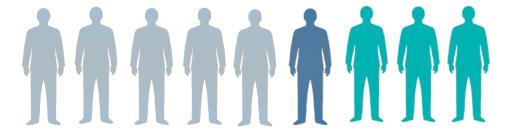
what is weighting?

when the sample is not representative of the population, we can disproportionally consider units to make the sample look more like the population.

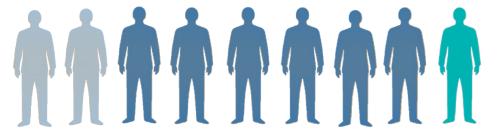


weighting and the counterfactual

population 1



population 2





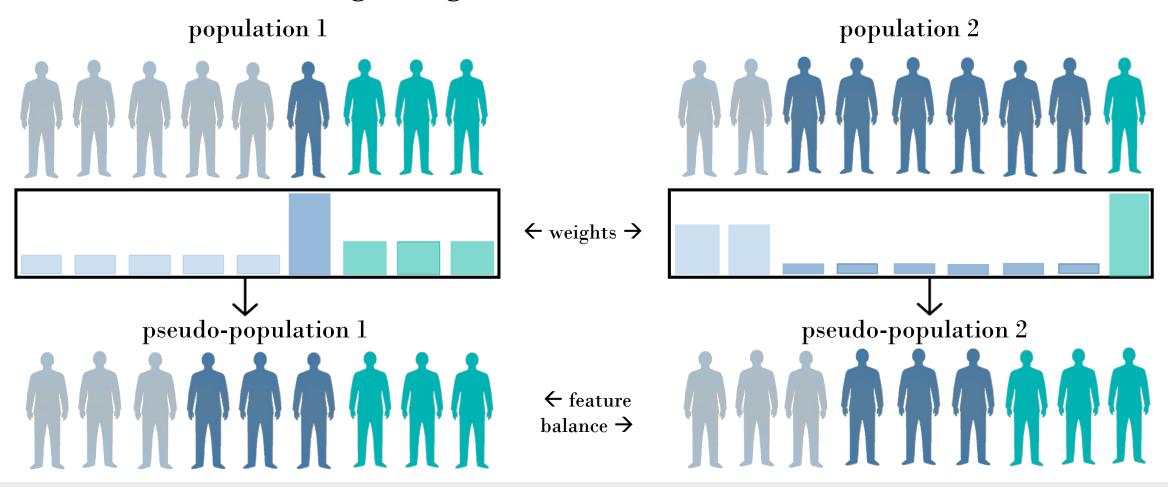
weighting and the counterfactual

population 1

Population 2



weighting and the counterfactual

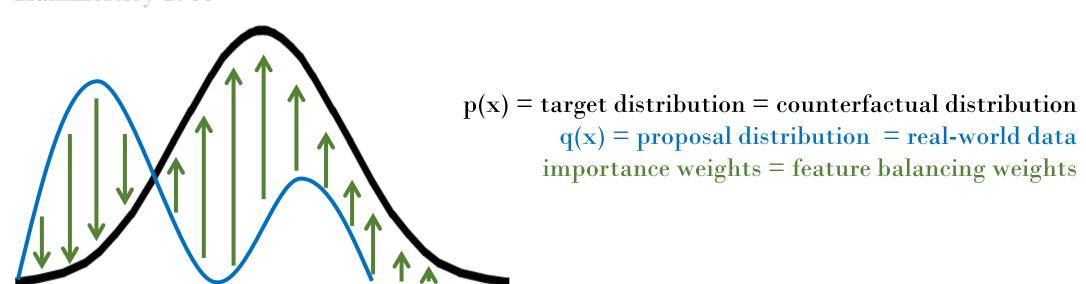






this is related to importance sampling

Importance sampling draws samples from a *proposal distribution* and re-weights the distribution using importance weights so that the weighted distribution represents your target distribution. Hammersley 1966







a common method of weighting



a common method of weighting

Inverse Probability of Treatment Weighting (IPW). Units are weighted according to inverse of their probability of being assigned to the treatment conditional on their measured, baseline features, a metric often called the propensity score. Rubin 2015; Rosenblatt 1965; Rosenbaum 1984; Austin 2011

$$w = \frac{T}{P(T=1|X)} + \frac{1-T}{P(T=1|X)}$$



a common method of weighting

Inverse Probability of Treatment Weighting (IPW). Units are weighted according to inverse of their probability of being assigned to the treatment conditional on their measured, baseline features, a metric often called the propensity score. Rubin 2015; Rosenblatt 1965; Rosenbaum 1984; Austin 2011

$$w = \frac{T}{P(T=1|X)} + \frac{1-T}{P(T=1|X)}$$



• flexible and robust causal modeling under selection on observables Imai 2013



- model dependent!
- unstable weights/feature imbalance/bias if propensity scores very close to 0 or 1 King 2016



how can we learn balancing weights for causal inference such that the weights are more stable and model agnostic?





how can we learn balancing weights for causal inference such that the weights are more stable and model agnostic?

use a Generative Adversarial Network (GAN)!





why GANs?



implicit generative models only specify a stochastic procedure with which to generate data Mohamed 2017

- full distributional matching on feature
- less prone to instability originating from model specification



why GANs?



implicit generative models only specify a stochastic procedure with which to generate data Mohamed 2017

- full distributional matching on feature
- less prone to instability originating from model specification

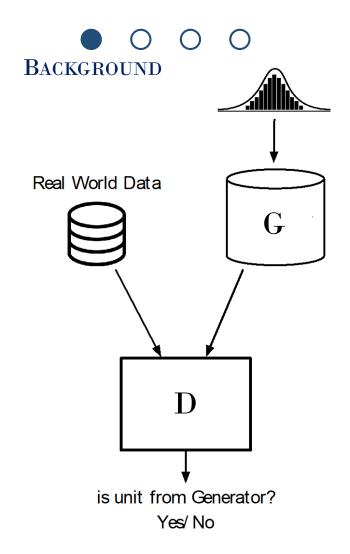
prescribed models provide an explicit parametric specification for a distribution Diggle 1984

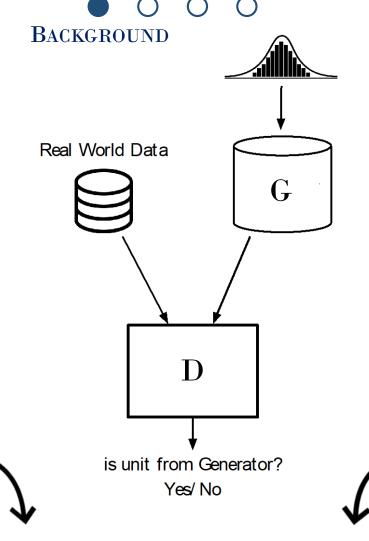
- often used to model propensity scores, etc that are used in weighting
- model dependence

 $\overline{\mathbf{VS}}$









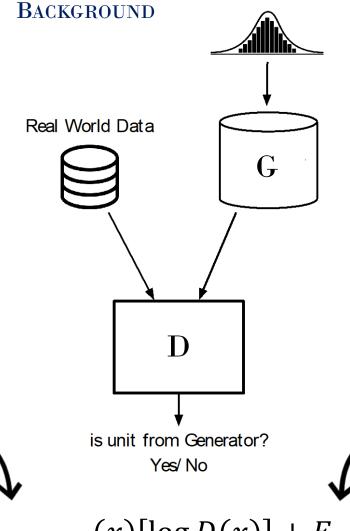
log probability of D predicting that "real" data is genuine

log probability of D predicting that "fake" data is not genuine

 $min_{G}max_{D}V(D,G) = E_{x \sim p_{data}}(x)[\log D(x)] + E_{z \sim p_{z}}(z)[\log(1 - D(G(z)))]$

the optimal solution to this expression minimizes the Jensen-Shannon divergence.

log probability of D predicting that "real" data is genuine



log probability of D predicting that "fake" data is not genuine

 $min_{G}max_{D}V(D,G) = E_{x \sim p_{data}}(x)[\log D(x)] + E_{z \sim p_{z}}(z)[\log(1 - D(G(z)))]$

vanilla GAN

the optimal solution to this expression minimizes the Jensen-Shannon divergence. BACKGROUND Real World Data G D is unit from Generator? Yes/ No

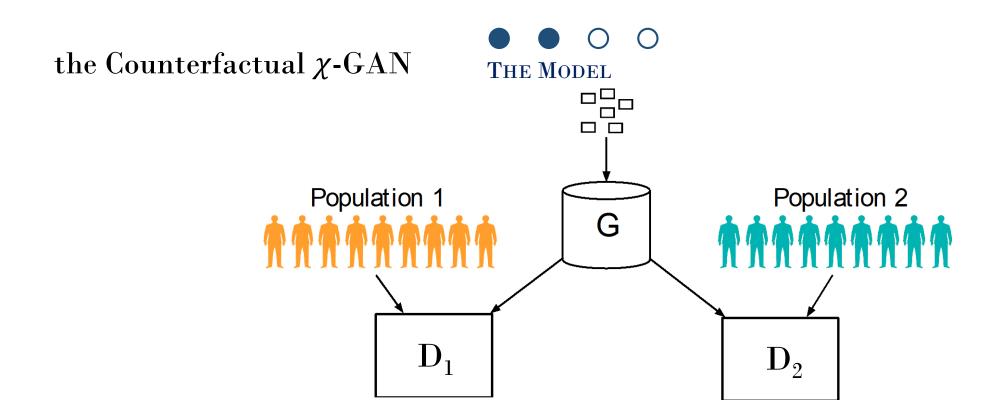
this won't suffice for causal inference. we need something new!

log probability of D predicting that "real" data is genuine



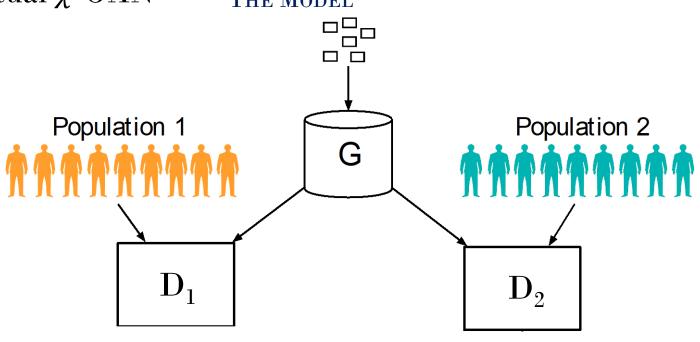
log probability of D predicting that "fake" data is not genuine

 $min_{G}max_{D}V(D,G) = E_{x \sim p_{data}}(x)[\log D(x)] + E_{z \sim p_{z}}(z)[\log(1 - D(G(z)))]$



• • O O
THE MODEL

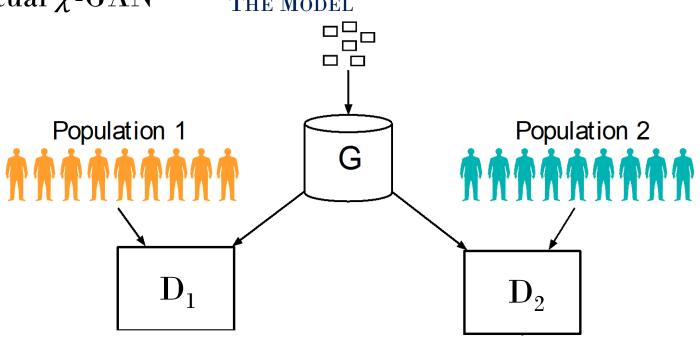
difference #1: two GANs joined at the generator



THE MODEL G

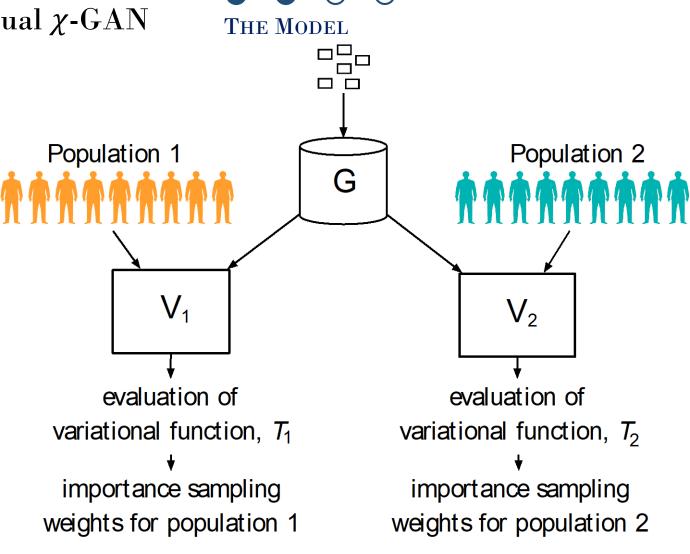
difference #1: two GANs joined at the generator

difference #2: minimize the χ divergence



difference #1: two GANs joined at the generator

difference #2: minimize the χ divergence Nowozin 2016



THE MODEL

difference #1: two GANs joined at the generator

difference #2: minimize the χ divergence

Population 1 Population 2 G V_2 evaluation of evaluation of variational function, T_1 variational function, T_2 importance sampling importance sampling weights for population 1 weights for population 2

learns featurebalancing weights through an adversarial training process.



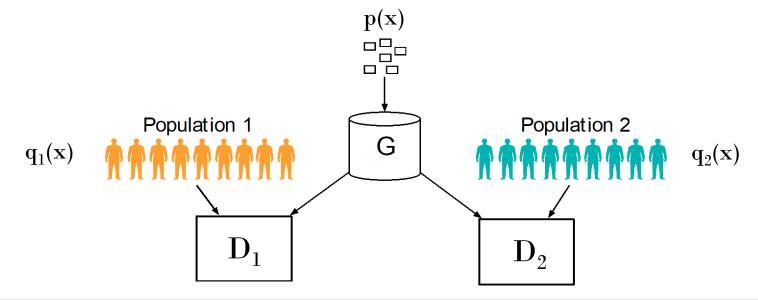
Nowozin 2016



Under importance sampling, those ranges with in which the ratio of p(x)/q(x)

- high \rightarrow will have high importance weights, contribute more to expectations
- low \rightarrow will have very small importance weights, contribute negligibly to expectations

the target distribution, p(x) functions as the generator, which embodies the *overlapping* portions of the empirical distributions, q(x) of the treatment arms.





$$\chi = \int q(x) \left[\frac{p(x)^2}{q(x)^2} - 1 \right] dx$$
Dieng 2016



$$\chi = \int q(x) \left[\frac{p(x)^2}{q(x)^2} - 1 \right] dx$$
Dieng 2016

$$\sigma_q^s = \frac{\mu_q^2}{n} \left[\int q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$



$$\chi = \int_{\text{Dieng 2016}} q(x) \left[\frac{p(x)^2}{q(x)^2} - 1 \right] dx$$
 change the objective function to minimize this divergence.
$$\sigma_q^s = \frac{\mu_q^2}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 change the objective function to minimize this divergence. equivalent to minimizing the variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 change the objective function to minimize this divergence.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.
$$\sigma_q^s = \frac{1}{n} \left[\int_{\text{C}} q(x) \left(\frac{p(x)}{q(x)} \right)^2 - 1 \right] dx$$
 variance.

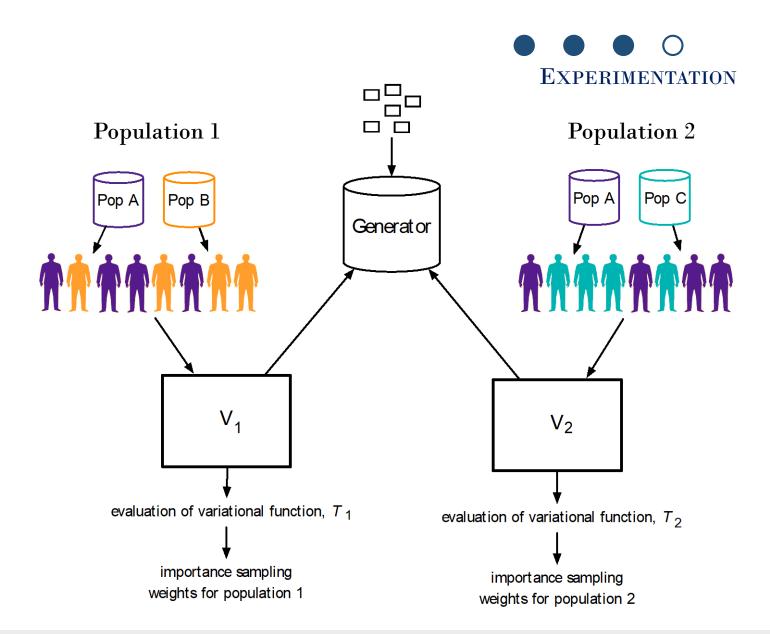


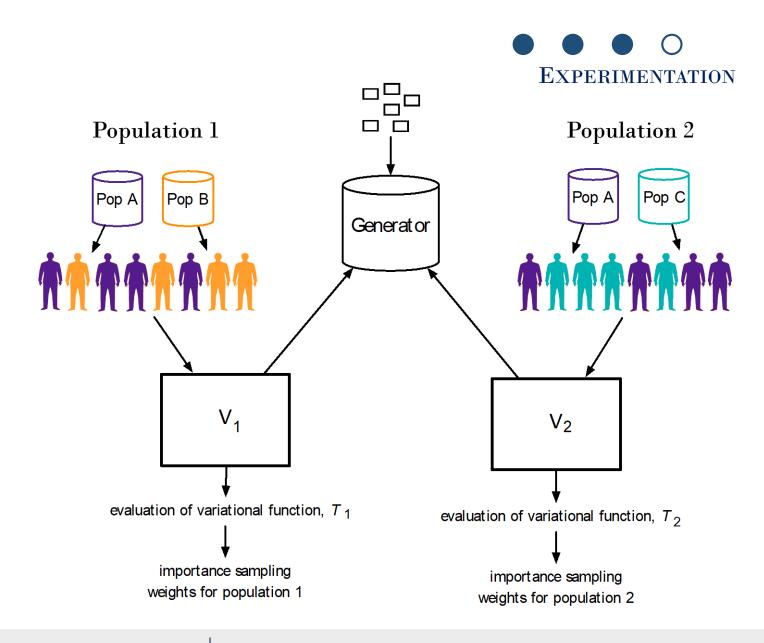
- 1. simulation
- 2. application to clinical data



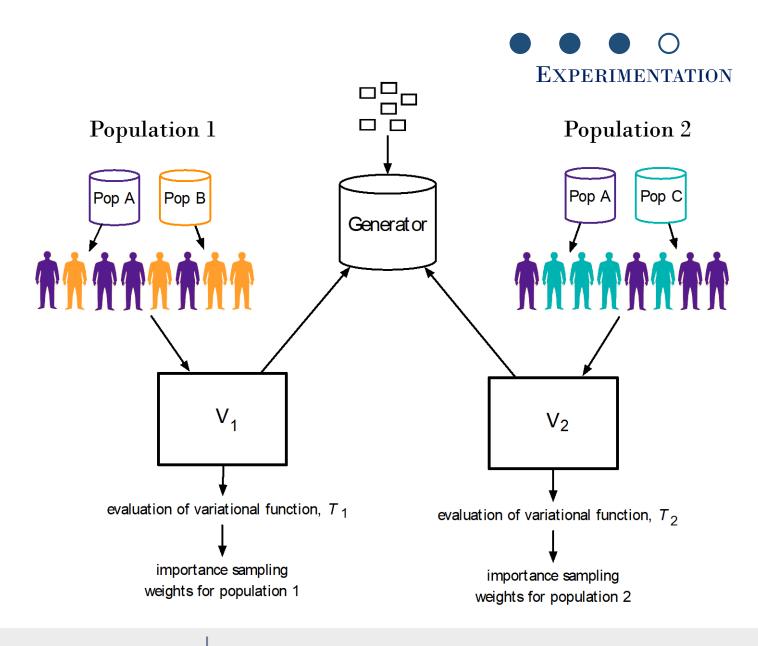
1. simulation

2. application to clinical data





- **3 subpopulations** A, B, C were drawn from randomly generated multivariate normal distributions.
- 5 discrete / 5 continuous features
- **2 populations** (Pop1 & Pop2) are mixtures of subpopulations
- Pop 1 = subpop A & subpop B
- Pop 2 = subpop A & subpop C
- 4000 per arm, 2000 per subpopulation



3 subpopulations A, B, C were drawn from randomly generated multivariate normal distributions.

• 5 discrete / 5 continuous features

2 populations (Pop1 & Pop2) are mixtures of subpopulations

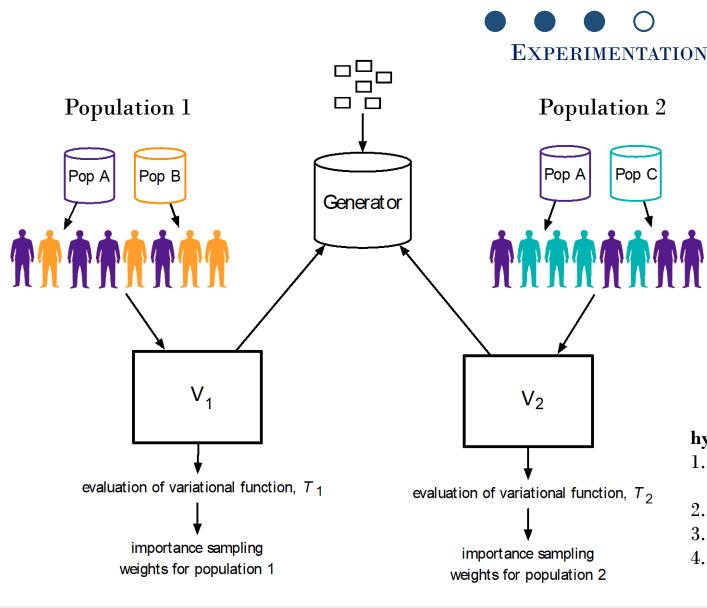
- Pop 1 = subpop A & subpop B
- Pop 2 = subpop A & subpop C
- 4000 per arm, 2000 per subpopulation

1 outcome. conditional on subpop

- Pop $1A \sim Gaussian (60, 1)$
- Pop 1B \sim Gaussian (40, 1)
- Pop $2A \sim Gaussian(-10, 1)$
- Pop $2C \sim Gaussian (10, 1)$

$$ATE_{mixture} = Pop1 - Pop2 = 50$$

$$ATE_{overlap} = Pop 1A - Pop2A = 70$$



3 subpopulations A, B, C were drawn from randomly generated multivariate normal distributions.

• 5 discrete / 5 continuous features

2 populations (Pop1 & Pop2) are mixtures of subpopulations

- Pop 1 = subpop A & subpop B
- Pop 2 = subpop A & subpop C
- 4000 per arm, 2000 per subpopulation

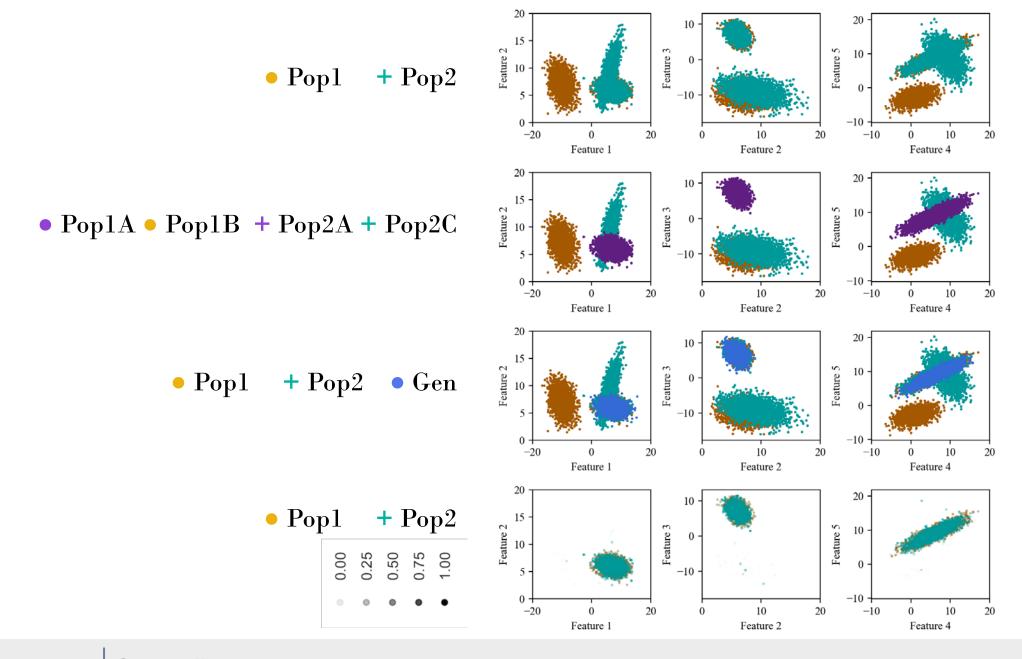
1 outcome. conditional on subpop

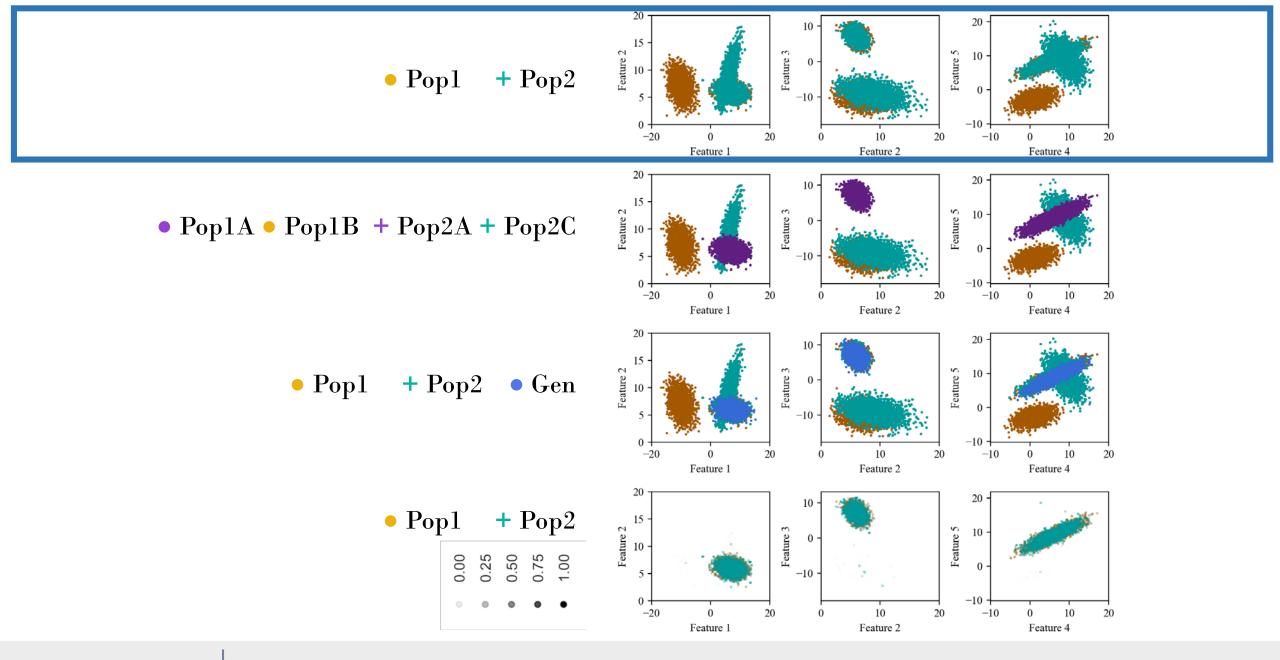
- Pop $1A \sim Gaussian (60, 1)$
- Pop 1B \sim Gaussian (40, 1)
- Pop $2A \sim Gaussian(-10, 1)$
- Pop $2C \sim Gaussian (10, 1)$

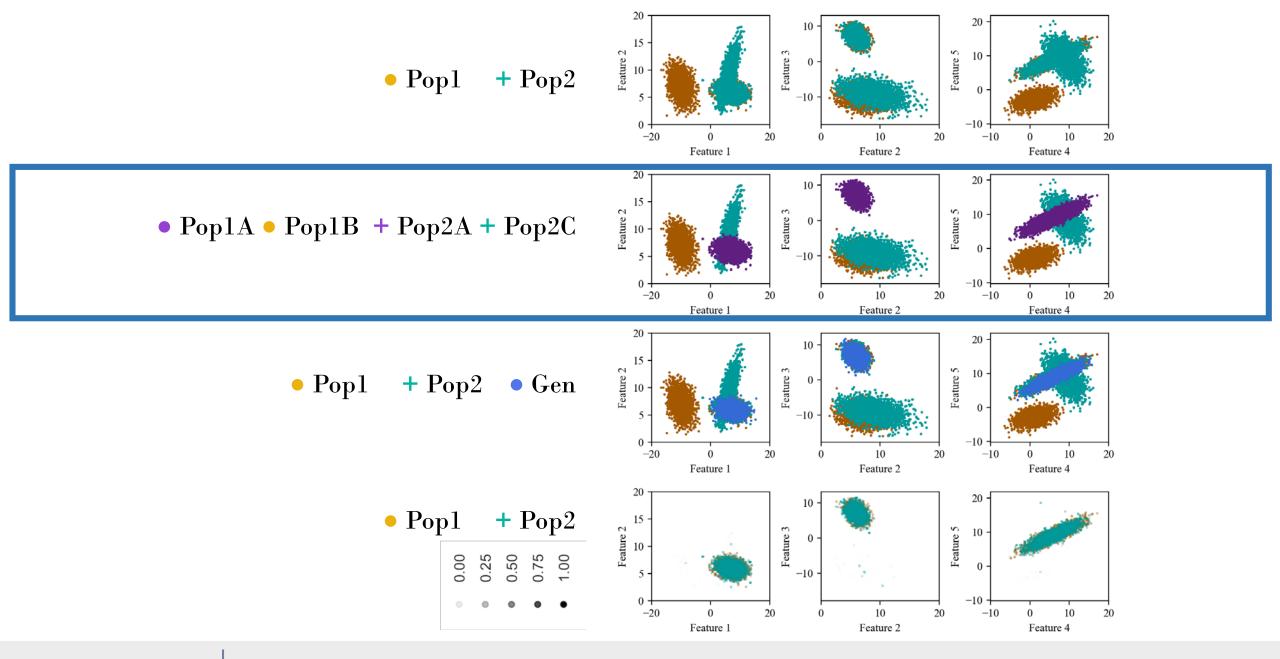
$$\begin{aligned} & \text{ATE}_{\text{mixture}} = \text{Pop1} - \text{Pop2} = 50 \\ & \text{ATE}_{\text{overlap}} = \text{Pop } 1\text{A} - \text{Pop2A} = 70 \end{aligned}$$

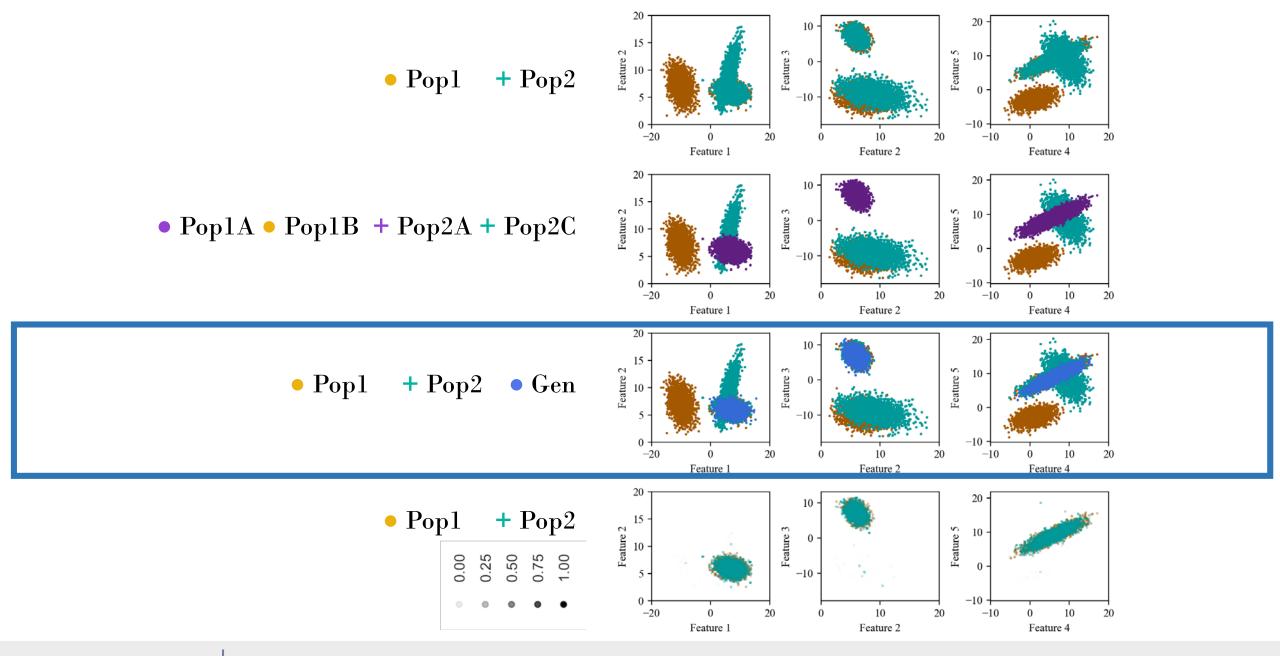
hypotheses

- 1. weights from overlapping population (Pop 1A/Pop 2A) will be high; weights from Pop 1B/Pop 2C will be low
- 2. weighting will make features more similar between Pops
- 3. cGAN-weighted ATE will be less biased than comparators
- 4. effective sample size of cGAN will be more reasonable than comparators

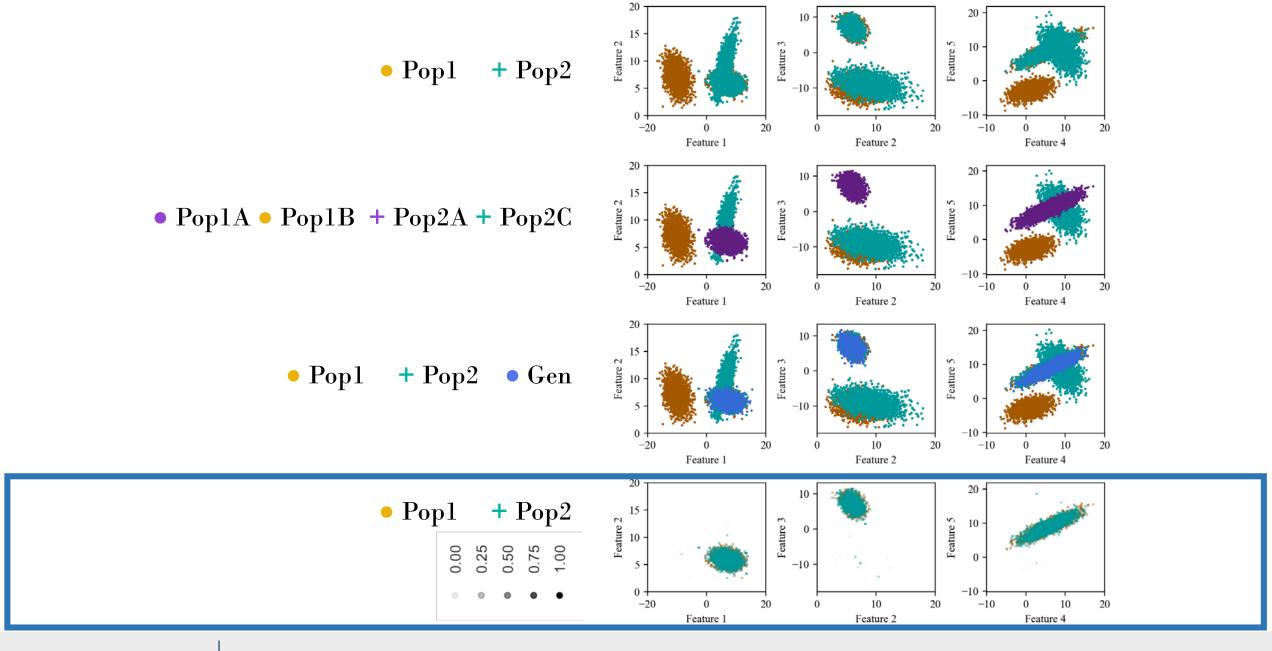




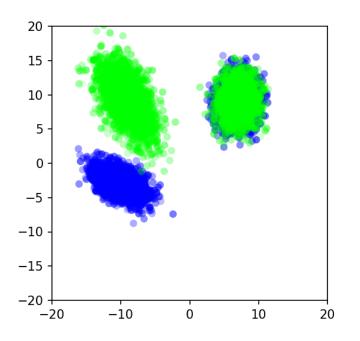


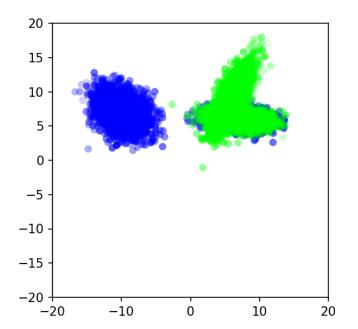


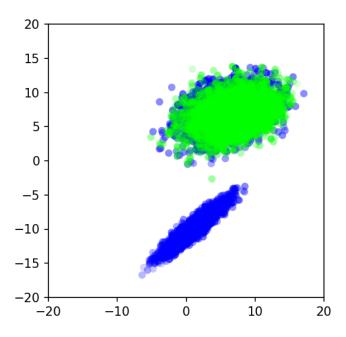




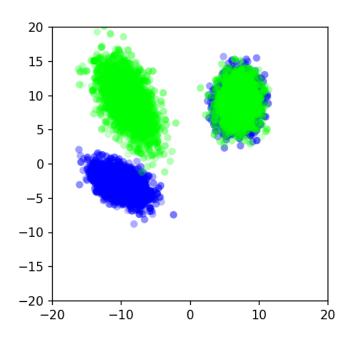


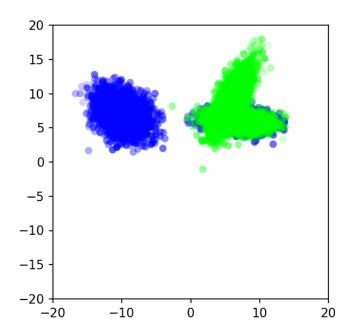


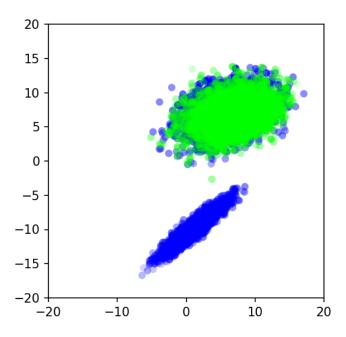






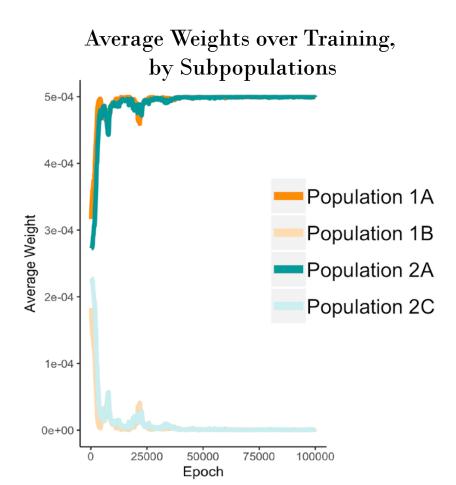








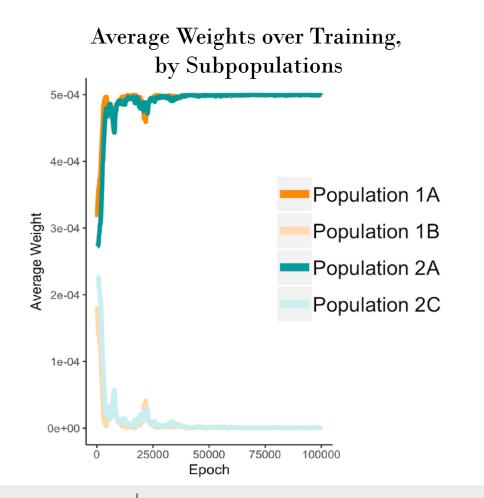
HYPOTHESIS 1: weights from Pop 1A/Pop 2A will be high; weights from Pop 1B/Pop 2C will be low

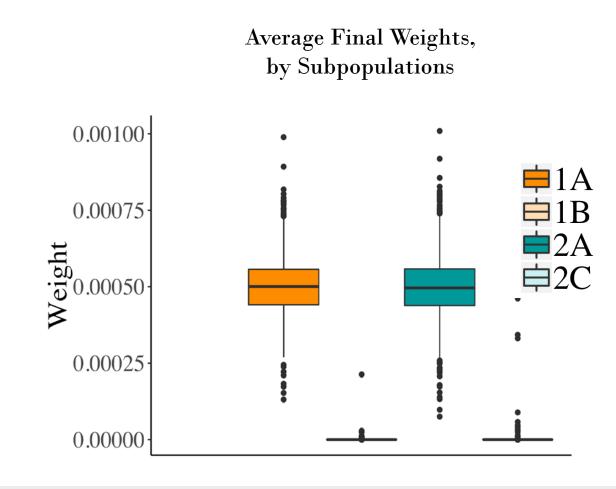






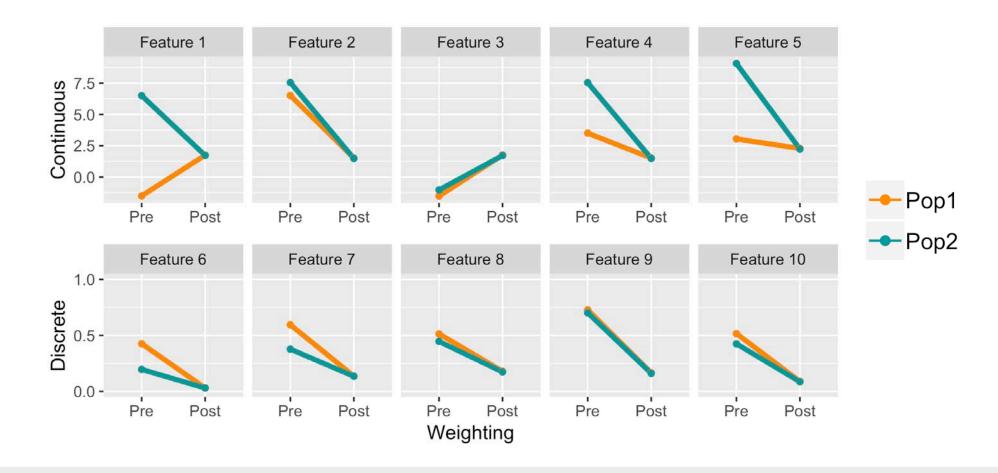
HYPOTHESIS 1: weights from Pop 1A/Pop 2A will be high; weights from Pop 1B/Pop 2C will be low







HYPOTHESIS 2: weighting will make features more similar between populations





 $ATE_{mixture} = 50$ $ATE_{overlap} = 70$

HYPOTHESIS 3: cGAN-weighted ATE will be less biased than comparators

Weighting Method
unweighted
cGAN
Inverse probability of treatment (IPTW)
Clipped IPTW
Binary regression propensity score
generalized boosted modeling of propensity scores McCaffrey 2004
covariate-balancing propensity scores Imai 2014
non-parametric covariate-balancing propensity scores Fong 2018
entropy balancing weights Hainmueller 2012
empirical balancing calibration weights Chan 2016
optimization-based weights Zubizarreta 2015



 $\begin{aligned} & \text{ATE}_{\text{mixture}} = 50 \\ & \text{ATE}_{\text{overlap}} = 70 \end{aligned}$

HYPOTHESIS 3: cGAN-weighted ATE will be less biased than comparators

Weighting Method
unweighted
cGAN
Inverse probability of treatment (IPTW)
Clipped IPTW
Binary regression propensity score
generalized boosted modeling of propensity scores McCaffrey 2004
covariate-balancing propensity scores Imai 2014
non-parametric covariate-balancing propensity scores Fong 2018
entropy balancing weights Hainmueller 2012
empirical balancing calibration weights Chan 2016
optimization-based weights Zubizarreta 2015

ATE
50.03
70.01
92.00
87.24
92.00
84.51
91.83
37.65
104.13
52.06
52.07



 $\begin{aligned} \text{ATE}_{\text{mixture}} &= 50\\ \text{ATE}_{\text{overlap}} &= 70 \end{aligned}$

HYPOTHESIS 4: effective sample size of cGAN will be more reasonable than comparators

Weighting Method
unweighted
cGAN
Inverse probability of treatment (IPTW)
Clipped IPTW
Binary regression propensity score
generalized boosted modeling of propensity scores McCaffrey 2004
covariate-balancing propensity scores Imai 2014
non-parametric covariate-balancing propensity scores Fong 2018
entropy balancing weights Hainmueller 2012
empirical balancing calibration weights Chan 2016
optimization-based weights Zubizarreta 2015

ATE
50.03
70.01
92.00
87.24
92.00
84.51
91.83
37.65
104.13
52.06
52.07

ESS Kish 1965
8000
3870
6551
6997
6551
7207
6686
11
65
65
114



- 1. simulation
- 2. application to clinical data





NewYork-Presbyterian OHDS



Efficacy and Tolerability of Sitagliptin Compared with Glimepiride in Elderly Patients with Type 2 Diabetes Mellitus and Inadequate Glycemic Control: A Randomized, **Double-Blind, Non-Inferiority Trial**

Paul Hartley¹ · Yue Shentu² · Patricia Betz-Schiff² · Gregory T. Golm² · Christine McCrary Sisk² · Samuel S. Engel² · R. Ravi Shankar²



clinicaltrials.gov/ct2/show/NCT01189890

Hartley 2015







Efficacy and Tolerability of Sitagliptin Compared with Glimepiride in Elderly Patients with Type 2 Diabetes Mellitus and Inadequate Glycemic Control: A Randomized, Double-Blind, Non-Inferiority Trial

Paul Hartley¹ · Yue Shentu² · Patricia Betz-Schiff² · Gregory T. Golm² · Christine McCrary Sisk² · Samuel S. Engel² · R. Ravi Shankar²



clinicaltrials.gov/ct2/show/NCT01189890

Hartley 2015

COLUMBIA UNIVERSITY DEPARTMENT OF BIOMEDICAL INFORMATICS

eligible patients

- diagnosis of Type II Diabetes
 Mellitus
- prescription to sitagliptin or glimepiride
- aged 65-80.





Efficacy and Tolerability of Sitagliptin Compared with Glimepiride in Elderly Patients with Type 2 Diabetes Mellitus and Inadequate Glycemic Control: A Randomized, Double-Blind, Non-Inferiority Trial

Paul Hartley¹ · Yue Shentu² · Patricia Betz-Schiff² · Gregory T. Golm² · Christine McCrary Sisk² · Samuel S. Engel² · R. Ravi Shankar²



clinicaltrials.gov/ct2/show/NCT01189890

eligible patients

- diagnosis of Type II Diabetes
 Mellitus
- prescription to sitagliptin or glimepiride
- aged 65-80.

a **sub-sample** of sitagliptin users was taken to match the count of the glimepiride arm (N=608 vs N=144).

not necessary

Hartley 2015







Efficacy and Tolerability of Sitagliptin Compared with Glimepiride in Elderly Patients with Type 2 Diabetes Mellitus and Inadequate Glycemic Control: A Randomized, Double-Blind, Non-Inferiority Trial

Paul Hartley¹ · Yue Shentu² · Patricia Betz-Schiff² · Gregory T. Golm² · Christine McCrary Sisk² · Samuel S. Engel² · R. Ravi Shankar²



clinicaltrials.gov/ct2/show/NCT01189890

eligible patients

- diagnosis of Type II Diabetes
 Mellitus
- prescription to sitagliptin or glimepiride
- aged 65-80.

a sub-sample of situaliptin users was taken to match the count of the glimepiride arm (N=608 vs N=144).

• not necessary

37 features

- repeated measurements: the most recent result was selected.
- missing data: values were imputed

Hartley 2015







Efficacy and Tolerability of Sitagliptin Compared with Glimepiride in Elderly Patients with Type 2 Diabetes Mellitus and Inadequate Glycemic Control: A Randomized, Double-Blind, Non-Inferiority Trial

Paul Hartley¹ · Yue Shentu² · Patricia Betz-Schiff² · Gregory T. Golm² · Christine McCrary Sisk² · Samuel S. Engel² · R. Ravi Shankar²



clinicaltrials.gov/ct2/show/NCT01189890

Hartley 2015



- diagnosis of Type II Diabetes
 Mellitus
- prescription to sitagliptin or glimepiride
- aged 65-80.

a **sub-sample** of sitagliptin users was taken to match the count of the glimepiride arm (N=608 vs N=144).

not necessary

37 features

- repeated measurements: the most recent result was selected.
- missing data: values were imputed

hypothesis

1. cGAN will improve feature balance over comparator methods



HYPOTHESIS 1: cGAN will improve feature balance over comparator methods

ASDM is a common metric of feature balance Austin 2011. A lower ASDM is indicative of feature balance

$$ASDM = |\frac{\bar{x}_{treatment} - \bar{x}_{control}}{\sqrt{\frac{s_{treatment}^2 + s_{control}^2}{2}}}|$$

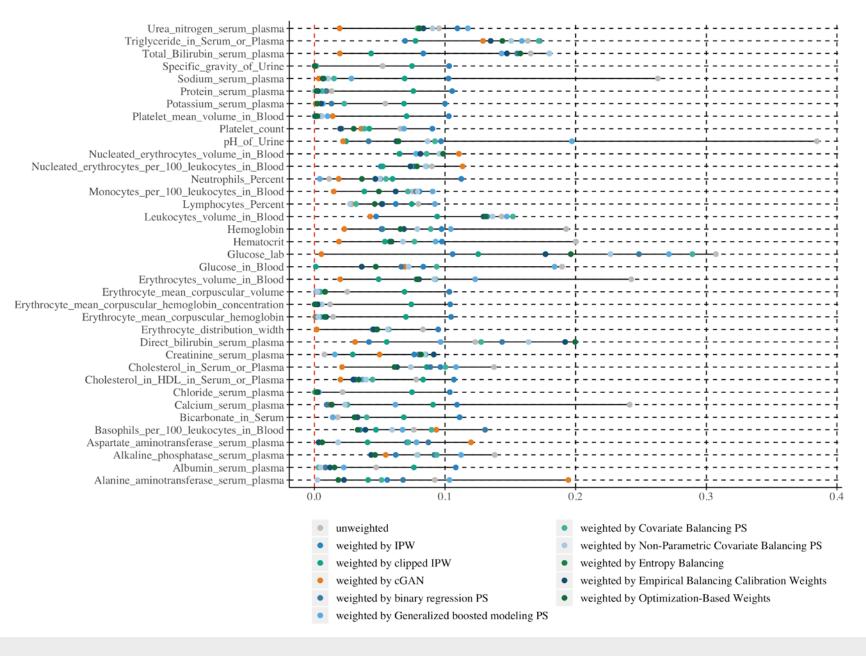
same comparators as simulation

HYPOTHESIS 1: cGAN will improve feature balance over comparator methods

ASDM is a common metric of feature balance Austin 2011. A lower ASDM is indicative of feature balance

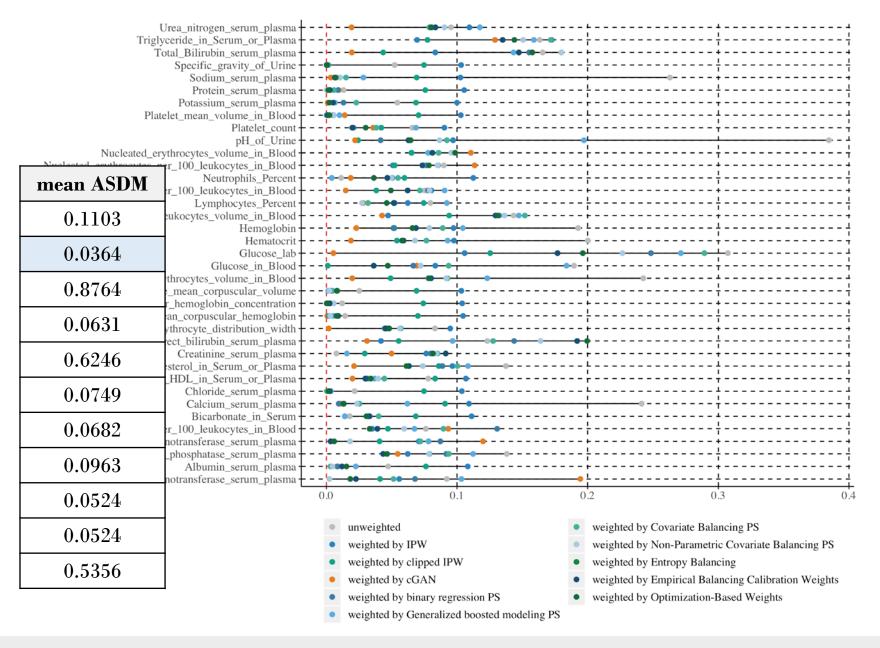
$$ASDM = |\frac{\bar{x}_{treatment} - \bar{x}_{control}}{\sqrt{\frac{s_{treatment}^2 + s_{control}^2}{2}}}$$

same comparators as simulation



HYPOTHESIS 1: cGAN will improve feature balance over comparator methods

Weighting Method
unweighted
cGAN
IPTW
clipped IPTW
binary regression
generalized boosted modeling
covariate-balancing
non-parametric covariate-balancing
entropy balancing weights
empirical balancing calibration
optimization-based weights







the experiments suggest that Counterfactual χ -GAN is an effective method of learning feature balancing weights to support counterfactual inference!



the experiments suggest that Counterfactual χ -GAN is an effective method of learning feature balancing weights to support counterfactual inference!

the Counterfactual χ -GAN could provide an alternative means to causal inference from observational data.





the experiments suggest that Counterfactual χ -GAN is an effective method of learning feature balancing weights to support counterfactual inference!

the Counterfactual χ -GAN could provide an alternative means to causal inference from observational data.

furthermore, if we assume that all potentially confounding variables are observed and included as features, average treatment effect estimates from Counterfactual χ -GAN weighted models may be less biased.



Limitations

GANs are unstable

Parameter tuning is hard

What is the best way to assess convergence?

discrete data - gradients are unbiased, but high variance





Limitations

GANs are unstable

Parameter tuning is hard

What is the best way to assess convergence?

discrete data - gradients are unbiased, but high variance

Future Directions

application to clinical data.

compare to RCT. need multisite

collaborators

assessing variance of outcome.
this requires a more complex
simulation





References





Thank you



Natnicha Vanitchanant



Rajesh Ranganath



Adler J. Perotte



aja2149@cumc.columbia.edu



http://people.dbmi.columbia.edu/AJA/



https://www.linkedin.com/in/ameliajaveritt



@AJAveritt

back up



learning for discrete data

We leverage a **score function estimator**. This score function-based estimator exchanges a gradient of an expectation for an expectation of a gradient which we can make an unbiased Monte Carlo estimate and incorporate into a modified stochastic backpropagation procedure.

$$\nabla_{\pi_a} \mathbb{E}_{\hat{q}_t(x;\pi_a)}[f(x)] = \mathbb{E}_{\hat{q}_t(x;\pi_a)}[f(x)\nabla_{\pi_a} \log \hat{q}_t(x;\pi_a)]$$

Glasserman 2003; Fu 2006; Schulman 2015