

Reproducing the results in *Optimal transport weights for causal inference*

Eric Dunipace

2/23/2022

Basics

This R Markdown document will reproduce the tables and figures in the paper *Optimal transport weights for causal inference*. If this file is compiled in this directory, the tables and figures will be generated from the source files found in the `code` folder.

To generate tables and figures

By default, this document will re-compile the tables and figures from the simulation studies using the already run simulation data from the original paper. This document will also input the tables and figures from the case study. The case study can be re-run by setting the corresponding chunk to `eval = TRUE`.

To re-run simulations

This will be a lengthy process if not on a cluster. The original files used to run the simulations on the cluster are found in the subfolder `code\original_sim` for both the convergence analysis and Hainmueller simulations. How to run these on a cluster is included below.

Required R packages

To re-run these simulations and/or recompile tables and figures, there are several libraries needed.

For just compiling tables and figures, here's all of the libraries that should be loaded:

```
library(causalOT)
library(dplyr)
library(ggplot2)
library(ggsci)
library(xtable)
library(cowplot)
library(tidyr)
```

For re-running simulations, the following packages are needed

```
library(causalOT)
library(doRNG)
library(Rmosek)
```

as well as the Python libraries `numpy`, `scipy`, `pykeops`, and `geomloss`.

Simulations studies

Study of bias/RMSE using Hainmueller (2012)

To generate the figures and tables for the simulations, we can use the following code

```
source("code/Hainmueller.R")
```

This will take the simulations for the Hainmueller setting and calculate the summary statistics. Then it generates the Latex table in Table 1. Some of the references in the table will not work since this is not in the larger paper with the .bib file.

Note this will download the simulation data if it's not already found in the data folder of this workflow.

Re-running on a cluster

To re-run this data, I recommend using a cluster. There are several other files in `code/original_sim` to discuss.

1. `hain_setting_array.R` generates a setting array that the cluster run will refer to in order to setup the simulation settings.
2. `generate_seeds.R` generates a list of seeds to be used by the analysis (and also for the convergence simulations)
3. `combine_sim_res.R` will combine the raw simulation results into one folder.

The files in `code/original_sim/Hainmueller` will re-run the raw analysis. The command

```
sbatch --array=1-1500 hain.sh
```

should re-run the results as desired.

Output

At the end we get the following table, which should match the paper.

design	overlap	method	constraint	Bias			RMSE		
				Hajek	DR	WOLS	Hajek	DR	WOLS
A	high	GLM	none	0.01	0.00	0.00	0.13	0.10	0.10
		CBPS	means	0.16	0.00	0.00	0.21	0.10	0.10
		SBW	means	0.00	0.00	0.00	0.09	0.09	0.09
		SCM	none	0.09	0.00	0.00	0.17	0.14	0.14
		COT	none	0.00	0.00	0.00	0.10	0.10	0.10
			means	0.00	0.00	0.00	0.28	0.28	0.28
	medium	GLM	none	0.04	0.00	0.00	0.19	0.11	0.11
		CBPS	means	0.24	0.00	0.00	0.28	0.11	0.11
		SBW	means	0.00	0.00	0.00	0.10	0.10	0.10
		SCM	none	0.14	0.00	0.00	0.21	0.15	0.15
		COT	none	0.00	0.00	0.00	0.12	0.12	0.12
			means	0.02	0.02	0.02	0.28	0.28	0.28
	low	GLM	none	0.09	0.00	0.00	0.28	0.12	0.12
		CBPS	means	0.26	0.00	0.00	0.32	0.12	0.12
		SBW	means	0.00	0.00	0.00	0.11	0.11	0.11
		SCM	none	0.22	0.01	0.01	0.28	0.17	0.17
		COT	none	0.00	0.00	0.00	0.14	0.14	0.14
			means	-0.01	-0.01	-0.01	0.29	0.29	0.29
B	high	GLM	none	-0.01	-0.01	-0.02	1.18	1.14	1.14
		CBPS	means	0.24	-0.01	-0.02	1.12	1.11	1.09
		SBW	means	-0.01	-0.01	-0.01	1.00	1.00	1.00
		SCM	none	0.36	0.27	0.28	1.63	1.57	1.55
		COT	none	0.01	0.01	0.01	0.61	0.61	0.61
			means	0.01	0.01	0.01	0.42	0.42	0.42
	medium	GLM	none	1.12	1.10	1.04	1.72	1.69	1.70
		CBPS	means	1.20	1.06	0.95	1.72	1.64	1.56
		SBW	means	0.63	0.63	0.63	1.20	1.20	1.20
		SCM	none	1.19	1.12	1.10	2.05	1.97	1.95
		COT	none	0.23	0.23	0.23	0.74	0.74	0.74
			means	-0.03	-0.03	-0.03	0.43	0.43	0.43
	low	GLM	none	0.19	0.06	0.02	1.72	1.49	1.51
		CBPS	means	0.45	0.06	0.01	1.42	1.46	1.42
		SBW	means	0.03	0.03	0.03	1.03	1.03	1.03
		SCM	none	0.64	0.42	0.43	1.75	1.69	1.65
		COT	none	0.05	0.05	0.05	0.85	0.85	0.85
			means	0.00	0.00	0.00	0.41	0.41	0.41

Table 1: Performance of various weighting methods under the simulation settings of [?]. Bold values are the values with the lowest bias or root mean-squared error (RMSE) of the methods under the same conditions. GLM refers to weighting by the inverse of the propensity score as calculated from a logistic regression model, CBPS is the covariate balancing propensity score, SBW is the stable balancing weights, SCM is the synthetic control method, and COT is the optimal transport formulation proposed in this paper. The estimators are Hajek weights (Hajek), doubly-robust augmented IPW (DR), and weighted least squares (WOLS). All weights are normalized to sum to 1. Constraints refer to balancing constraints and are one of “none” for no constraints or “mean” for mean constraints.

Convergence and confidence intervals

Similarly, the convergence and confidence interval simulations can be compiled with the following code chunk

```
source("code/convergence.R")
```

Note this will download the simulation data if it’s not already found in the data folder of this workflow.

Re-running on a cluster

To re-run this data, I again **recommend** using a cluster. There are several other files in `code/original_sim` as mentioned previously:

1. `generate_seeds.R` generates a list of seeds to be used by the analysis (and also for the convergence

simulations)

2. `combine_sim_res.R` will combine the raw simulation results into one folder.

Then the files in `code/original_sim/Convergence` will re-run the analysis. On a Slurm based cluster running

```
sbatch --array=1-1000 conv.sh
```

should be sufficient.

Outputs for convergence

Then we can look at the plots demonstrating convergence in terms of Sinkhorn divergence and L_2 norm.

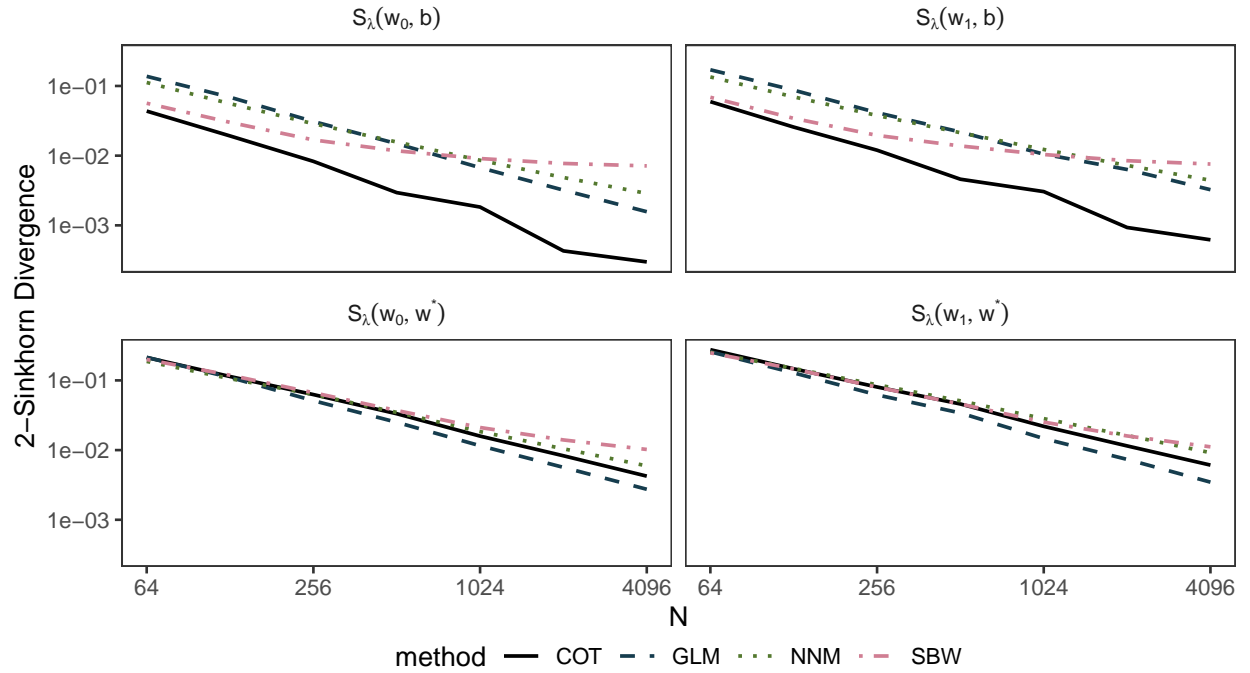


Figure 1: Convergence of the weights to the distributions specified by the empirical distributions (top) and the distributions specified by the true propensity score/Radon-Nikodym derivatives (bottom). Weights are a Causal Optimal Transport (COT), Nearest Neighbor Matching (NNM), a Probit model (GLM), and Stable Balancing Weights (SBW). Lines denote means across 1000 simulations. Both axes are on the log scale.

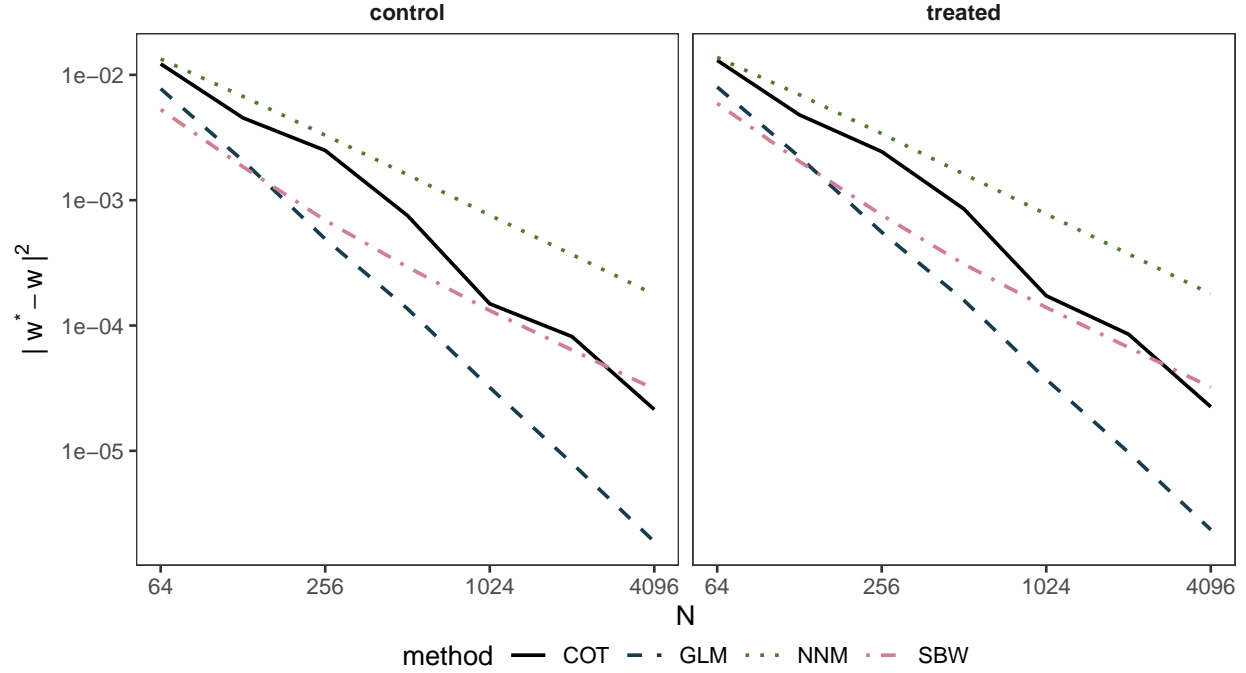


Figure 2: Convergence of the estimated weights to the values of the true inverse propensity score in terms of the L_2 norm. Weights are a Causal Optimal Transport (COT), Nearest Neighbor Matching (NNM), a Probit model (GLM), and Stable Balancing Weights (SBW). Lines denote means across 1000 simulations. Both axes are on the log scale.

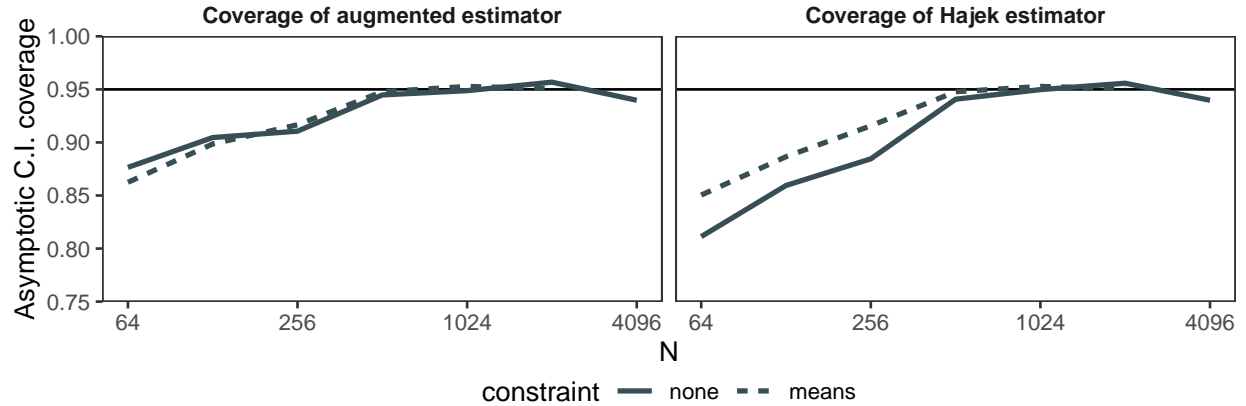


Figure 3: Coverage of the true treatment effect

Outputs for confidence intervals

This analysis also gets the following figures for confidence interval coverage.

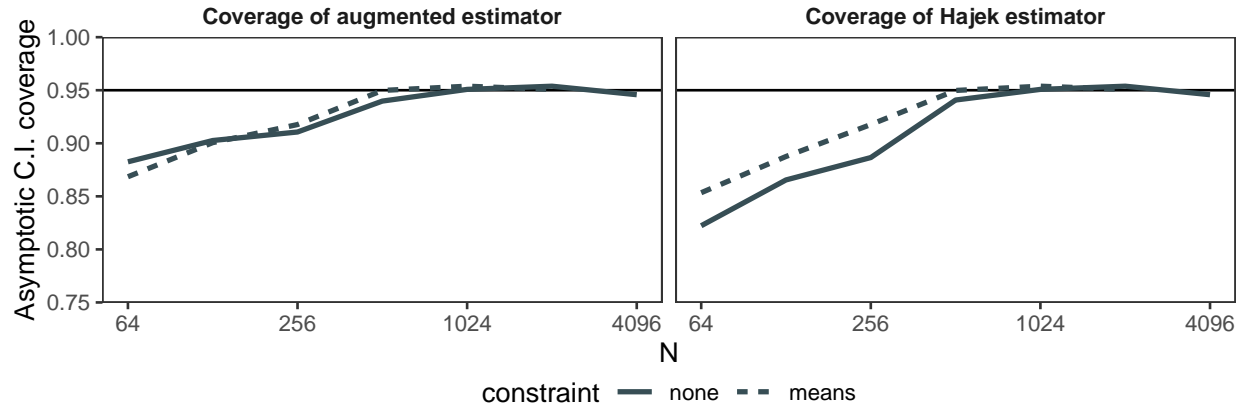


Figure 4: Coverage of the estimated average treatment effect

Case study

The data

The data come from a study by Blum et al. looking at the effect of misoprostol vs. oxytocin at stopping post-partum hemorrhage. The data is described in the documents for the `causalOT` package which can be accessed by running `?causalOT::pph`. A more detailed description is found in the paper *Optimal transport weights for causal inference*.

Running the analysis

The analysis can be run by running the following chunk with `eval = TRUE`

```
source("code/miso_egypt.R")
```

Figures/Tables

We now turn to the tables and figures generated by this analysis. The first is the mean balance of the covariates before and after weighting.

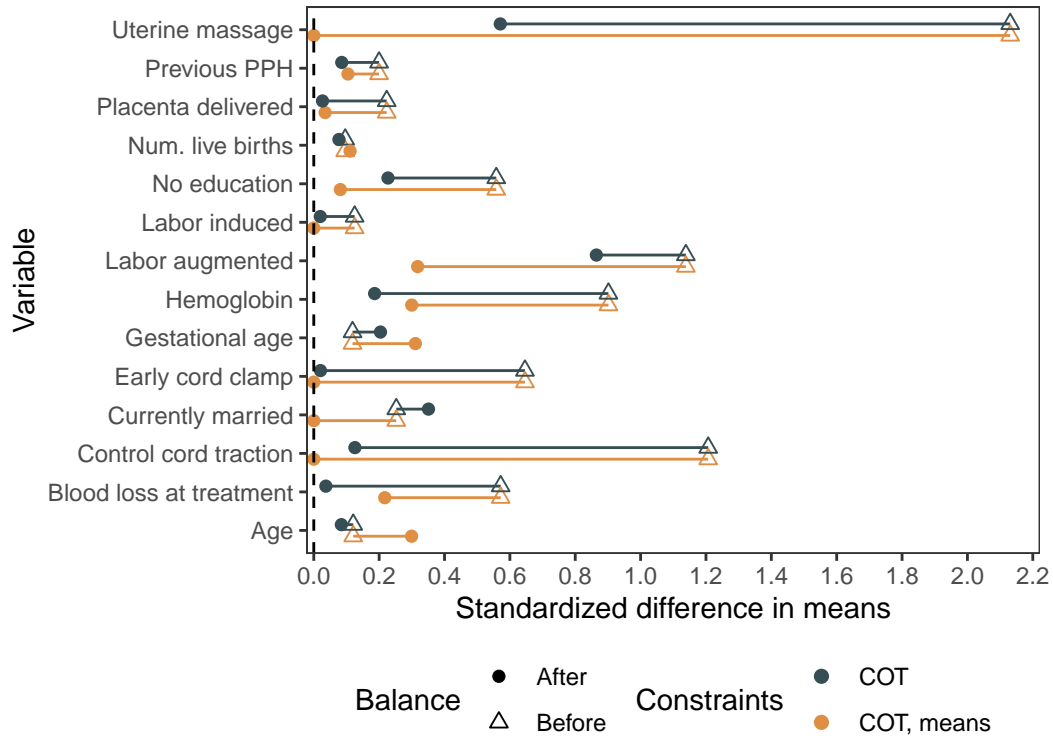


Figure 5: Change in the standardized difference in means between the two groups before and after weighting. An examination in the change in balance before and after utilizing the optimal transport methods with the listed constraints for the misoprostol receiving participants in Egypt versus the oxytocin receiving participants at the other sites. “COT” corresponds to no constraints and “COT, means” corresponds to constraints on mean balance.

The next figure looks at the changes in 2-Sinkhorn divergence before and after weighting.

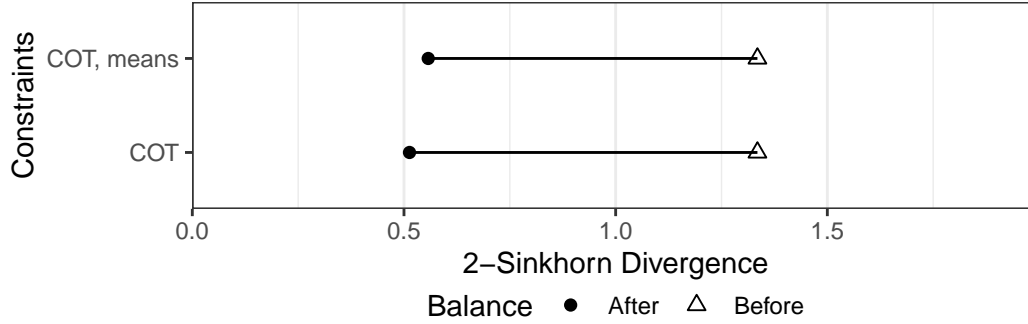


Figure 6: Change in the 2-Sinkhorn divergence between the two groups before and after weighting. An examination in the change in balance before and after utilizing the optimal transport methods with the listed constraints for the misoprostol receiving participants in Egypt versus the oxytocin receiving participants at the other sites. “COT” corresponds to no constraints and “COT, means” corresponds to constraints on the mean balance.

We can see that balance is improved though not perfect.

However, some of the effect estimates look real good!

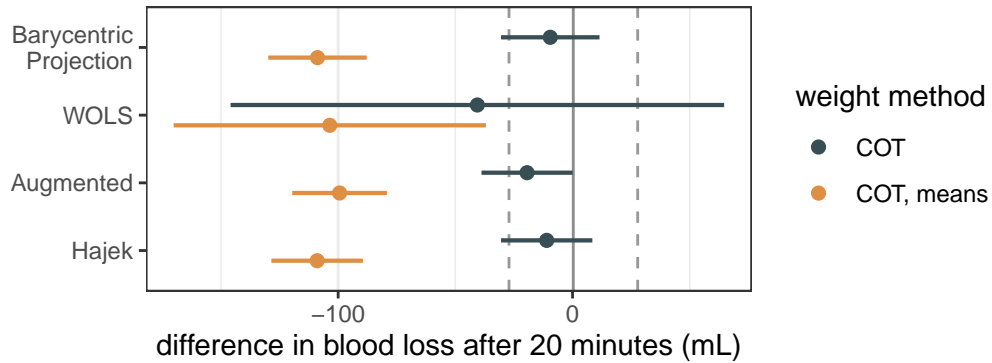


Figure 7: Results for treatment effect estimation between the treated misoprostol group in Egypt versus the control oxytocin group at four other sites. We see that methods without mean constraints are able to recover the effect estimate. Note that “COT” corresponds to no constraints and “COT, means” corresponds to constraints on mean balance.

We can see how mean balancing leads to very poor estimates.

This figure is also available in table form.

Method	Hajek		Augmented		Weighted OLS		Barycentric projection	
	Est.	C.I.	Est.	C.I.	Est.	C.I.	Est.	C.I.
COT	-11.0	(-30.5, 8.5)	-19.4	(-38.8, 0.0)	-40.6	(-145.9, 64.7)	-9.5	(-30.5, 11.5)
COT, means	-108.9	(-128.5, -89.4)	-99.4	(-119.6, -79.2)	-103.6	(-170.2, -36.9)	-108.8	(-129.8, -87.8)

Table 2: Estimates and confidence intervals for optimal transport weighting methods applied to a modification of the data in [?]. We have constructed a control group for the misoprostol receiving patients at Egyptian site using the controls from the four other sites. The original treatment effect at the Egypt site was 0.369 mL with a 95% C.I. of $(-27.1, 27.8)$. The augmented method uses a gaussian process as the outcome model.

This allows us to see better the overlap in confidence intervals.

Conclusions

This document reproduces the results for the paper *Optimal transport weights for causal inference* by Eric Dunipace. Any questions or concerns can be addressed by filing an issue at <https://github.com/ericdunipace/causalOT/issues>.