

OTClean: Data Cleaning for Conditional Independence Violations using Optimal Transport

Alireza Pirhadi
Mohammad Hossein Moslemi
Mostafa Milani



Babak Salimi
Alexander Cloninger



Content

- **Motivation:** Conditional Independence (CI) Constraints
 - Algorithmic Fairness and Data Cleaning
- **Problem Formulation**
 - Data Repair w.r.t. CI Constraints using Optimal Transport (OT)
- **Solutions**
 - Quadratic Constrained Linear Program (QCLP) Formulation
 - FastOTClean and Optimization
- **Experimental Results**

CI and Algorithmic Fairness

Dissecting racial bias in an algorithm used to manage the health of populations

Ziad Obermeyer^{1,2,*}, Brian Powers³, Christine Vogeli⁴, Sendhil Mullainathan^{5,*†}

† See all authors and affiliations

Science 25 Oct 2019;
Vol. 366, Issue 6464, pp. 447-453
DOI: 10.1126/science.aax2342

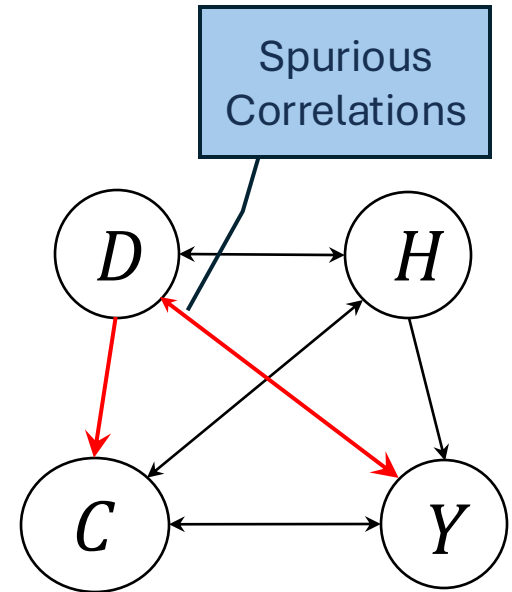


Health systems rely on commercial prediction algorithms to identify and help patients with complex health needs. We show that a widely used algorithm, typical of this industry-wide approach and **affecting millions of patients**, exhibits significant **racial bias: At a given risk score, Black patients are considerably sicker than White patients, as evidenced by signs of uncontrolled illnesses ... despite health care cost appearing to be an effective proxy for health by some measures of predictive accuracy, large racial biases arise ...**

CI and Algorithmic Fairness

- **Demographic Info (D):** Contains Race, Gender, and...
- **Health Metrics (H):** Encompasses Medical history and...
- **Cost Data (C):** Consists of Financial records related to...
- **Label (Y):** Healthcare Needs

$$(D \perp C | H)$$

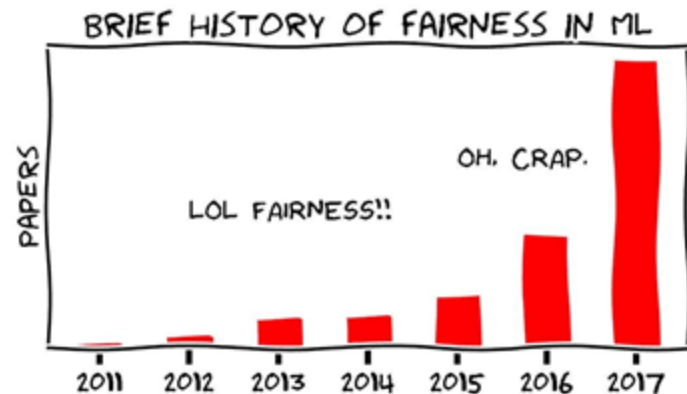


CI and Algorithmic Fairness (Another Example)

Priors (prior arrests or convictions), Age, Race, and Recidivism (label)



- **Statistical Parity:**
Predictions \perp Race
- **Equality of Odds:**
Predictions \perp Race | Recidivism
- **Conditional Statistical Parity:**
Predictions \perp Race | Prior

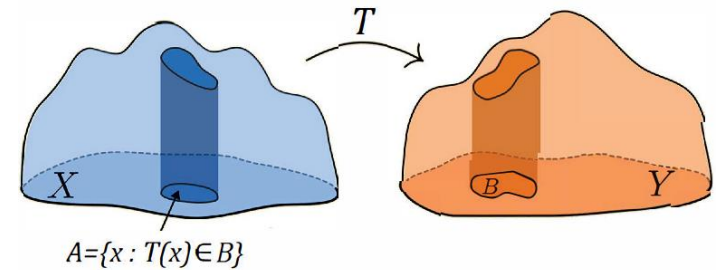


CI and Data Cleaning

- Sporadic correlations in training data due to
 - Erroneous values and dirty data
 - Biases (sampling, measurement or other biases)
- Downstream ML models underperform since these sporadic correlations do not exist at usage time

Optimal Transport (Monge Formulation)

- The most efficient way of transferring mass from a probability distribution P to another distribution Q



- Monge formulation:

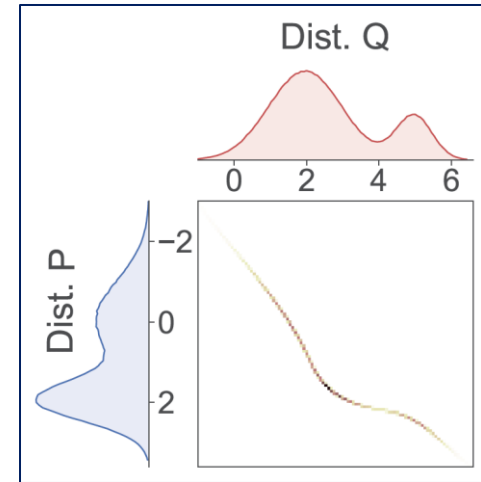
$$OT_{Monge}(P, Q) = \operatorname{argmin}_{T: \mathcal{X} \rightarrow \mathcal{Y}} \sum_{\mathbf{x}_i \in \mathcal{X}} c(\mathbf{x}_i, T(\mathbf{x}_i))$$

Cost function

Transport Map

Optimal Transport (Kantorovich)

- Kantorovich formulation of OT uses a probabilistic map or “plan”
- Kantorovich formulation:



$$OT(P, Q) = \operatorname{argmin}_{\pi \in \Pi(P, Q)} \sum_{x_i \in \mathcal{X}, y_i \in \mathcal{Y}} c(x_i, y_j) \pi(x_i, y_j)$$

Transport Plan

Problem Formulation

- A CI Constraint σ : $X \perp\!\!\!\perp Y \mid Z$

$$P_{X,Y|Z}(\bar{x}, \bar{y}, \bar{z}) = P_{X|Z}(\bar{x}, \bar{z}) \times P_{Y|Z}(\bar{y}, \bar{z})$$

- CI Data cleaner



Problem Formulation

- CI data cleaner T^* :

$$T^* = \operatorname{argmin}_{T: \mathcal{V} \rightarrow \mathcal{V}} \sum_{v_i \in D} c(v_i, T(v_i)) \quad s.t. \quad T(D) \models \sigma$$

- Probabilistic optimal data cleaner π^* :

$$\pi^* = \operatorname{argmin}_{\pi} \sum_{v_i, v'_j \in \mathcal{V}} c(v_i, v'_j) \pi(v_i, v'_j)$$
$$s.t. \quad \pi(v) = P^D, \pi(v') \models \sigma$$

First Solution: QCLP

- Quadratically Constrained Linear Program (QCLP)

- **Objective function:**

$$\min_{\hat{\pi}} \sum_{i,j \in [1, d_V]} c(v_i, v_j) \times \hat{\pi}_{i,j}$$

- **Constraints:**

1. Validity constraints:

$$\hat{\pi}_{i,j} \geq 0 \quad \forall i \in [1, d_V], j \in [1, d_V]$$

2. Marginal constraints:

$$\sum_{j \in [1, d_V]} \hat{\pi}_{i,j} = P^D(v_i) \quad \forall i \in [1, d_V]$$

3. Independence constraints: ensures satisfaction of CI constraint σ

Second Solution: Fast Approximation

- Relaxed OT with Entropic Regularize:

The entropic regularization parameter

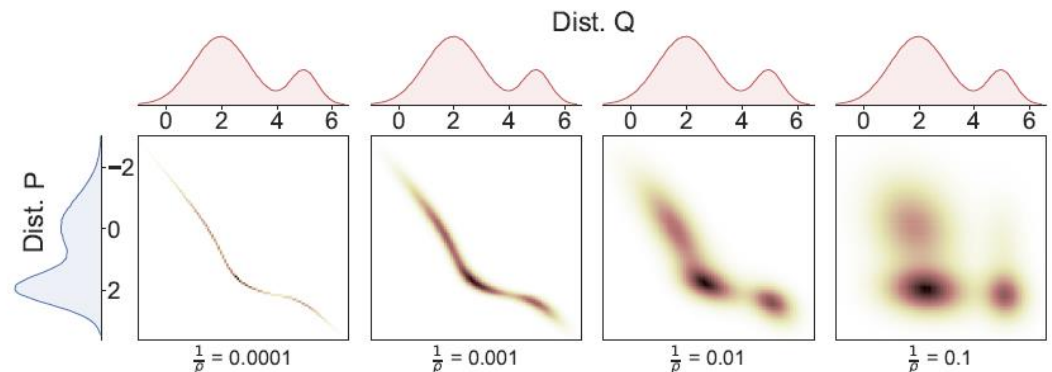
$$\operatorname{argmin}_{\pi} \sum_{x_i \in \mathcal{X}, y_i \in \mathcal{Y}} c(x_i, y_j) \pi(x_i, y_j) - \frac{1}{\rho} H(\pi)$$

$$+ \lambda \left(D_{KL}(\pi(v'), Q) + D_{KL}(\pi(v), P^D) \right) + \mu \delta_{\sigma}(Q)$$

The relaxation regularization coefficient

The CI constraint term

Entropic regularization allows
Sinkhorn algorithm
[M. Cuturi, NIPS'13]



Second Solution: FastOTClean

Algorithm 2: FASTOTCLEAN: Fast Computation of Probabilistic Data Cleaner for Conditional Independence

Input: Database D , cost function c , and CI constraint

$$\sigma : X \perp\!\!\!\perp Y \mid Z$$

Output: Transport plan (probabilistic data cleaner) π

```

1  $\mathbf{p} := \text{vector}(P^D)$ ;  $\mathbf{C} := \text{matrix}(c)$ ;
2 Randomly initialize  $\mathbf{q}$  ▷ An initial guess for  $Q$ 
3  $\mathbf{u} := \mathbb{1}_{d_X}$ ;  $\mathbf{v} := \mathbb{1}_{d_Y}$ ;  $\mathbf{K} := e^{-\frac{c}{\rho}}$ ; ▷ Sinkhorn Initialization
4 while  $\mathbf{q}$  is not converged do ▷ Sinkhorn iterations
5   while  $\mathbf{u}$  and  $\mathbf{v}$  are not converged do
6      $\mathbf{u} := (\mathbf{p} \oslash (\mathbf{K} \cdot \mathbf{v}))^{\frac{\rho\lambda}{\rho\lambda+1}}$ ,  $\mathbf{v} := (\mathbf{q} \oslash (\mathbf{K} \cdot \mathbf{u}))^{\frac{\rho\lambda}{\rho\lambda+1}}$ ;
7      $\pi = \text{diag}(\mathbf{u}) \cdot \mathbf{K} \cdot \text{diag}(\mathbf{v})$ ;
8   for each  $z \in \mathcal{Z}$  do
9     Initialize  $\mathbf{W}_z, \mathbf{H}_z$  randomly.
10    while  $\mathbf{W}_z$  and  $\mathbf{H}_z$  are not converged do
11      Update  $\mathbf{W}_z$  to minimize
12         $D_{\text{KL}}(\pi(X', Y', Z' = z) \mid \mathbf{W}_z \cdot \mathbf{H}_z^T)$  with  $\mathbf{H}_z$  fixed
13      Update  $\mathbf{H}_z$  to minimize
14         $D_{\text{KL}}(\pi(X', Y', Z' = z) \mid \mathbf{W}_z \cdot \mathbf{H}_z^T)$  with  $\mathbf{W}_z$  fixed
15    Construct  $\mathbf{q}$  using  $\mathbf{W}_z$ s and  $\mathbf{H}_z$ s computed in the previous step
16 return  $\pi$ ;
```

Initializations

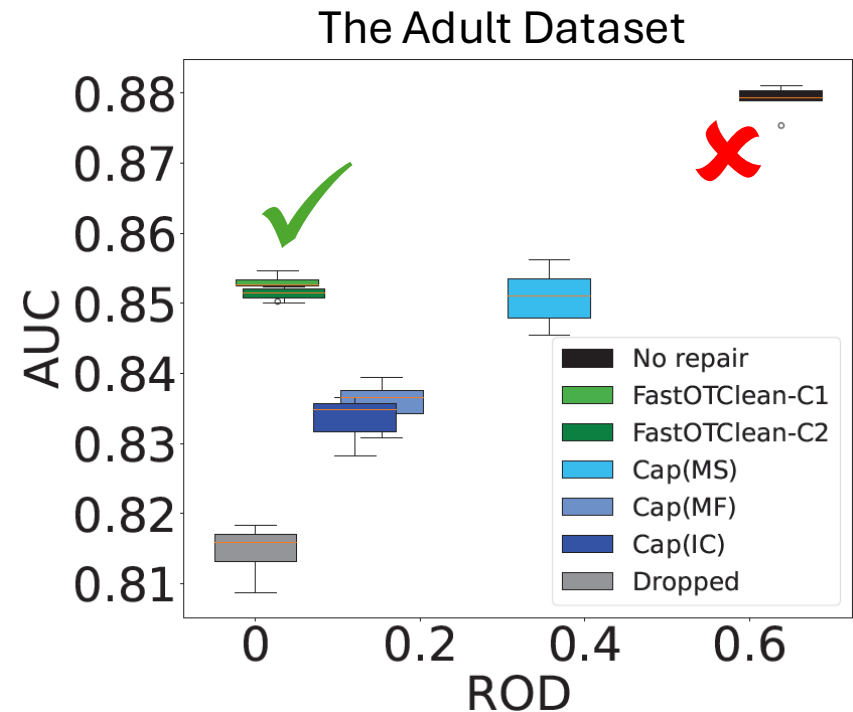
Sinkhorn iterations

CI constraints and alternating updates

Experimental Results (Fairness)

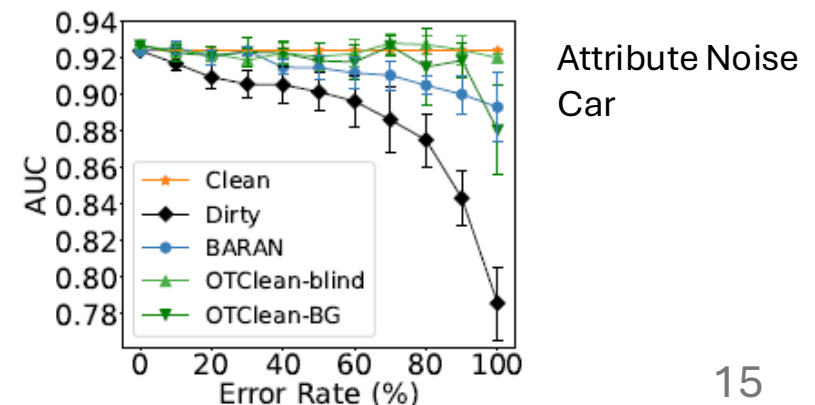
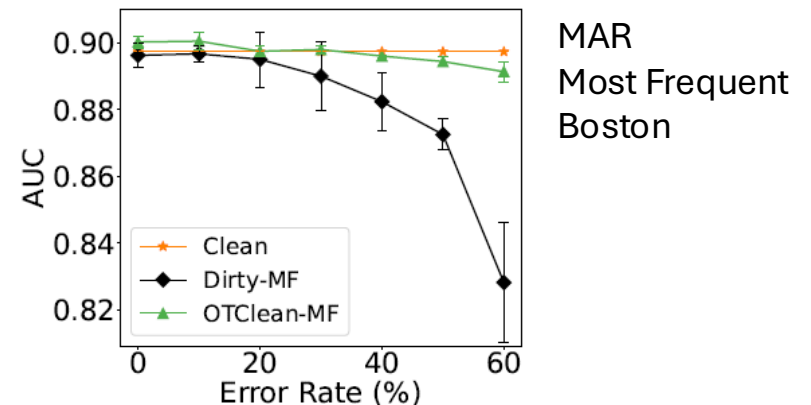
- **Datasets:** Adult and COMPAS
- OTClean with two cost functions

- **Baselines:**
 - Capuchin (Cap): three variations
 - [Salimi et al., SIGMOD'19]
 - No repair
 - Dropped
- **Performance Measures:**
 - AUC for accuracy
 - ROD (Ratio of Observation Discrimination) for fairness



Experimental Results (Data Cleaning)

- Datasets: Car and Boston
- Attribute noise and missing values (MAR and MNAR)
- Performance Measures:
 - AUC for accuracy
- Baselines
 - Clean data
 - Dirty data
 - BARAN (attribute noise)
 - BARAN (missing values)
 - Knn, MF, etc. (missing values)



References

- Z. Obermeyer, B. Powers, C. Vogeli, and S. Mullainathan, Dissecting racial bias in an algorithm used to manage the health of populations, Science, 2019
- S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing and machine-learning applications. IEEE Signal Processing Magazine, 34(4):43–59, 2017
- M. Cuturi, Sinkhorn distances: Lightspeed computation of optimal transport, NIPS, 2013
- B Salimi, A. L. Rodriguez, B. Howe, D. Suciu, Interventional Fairness: Causal Database Repair for Algorithmic Fairness, SIGMOD, 2019