

---

# Optimal Transport under Group Fairness Constraints

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Ensuring fairness in matching algorithms is a key challenge in allocating scarce  
2 resources and positions. Focusing on Optimal Transport (OT), we introduce a novel  
3 notion of group fairness requiring that the probability of matching two entities from  
4 any two given groups in the OT plan follows a user-specified target. We develop two  
5 relaxation strategies to solve this constrained problem. The first one involves solv-  
6 ing a penalized OT problem, for which we derive novel finite-sample complexity  
7 guarantees. Our second strategy leverages bi-level optimization to learn a ground  
8 cost that produces a fair OT solution, which can be reused to match new samples.

## 9 1 Introduction

10 Algorithmic matching mechanisms play an increasing role in modern societies, handling the distri-  
11 bution of rare goods by connecting individuals or firms through algorithmic decisions rather than  
12 through price driven markets. Examples of such mechanisms include online job recommendations,  
13 college admissions systems, and kidney allocation circuits. An increasing concern is the fairness  
14 of such mechanisms. For instance, school assignment systems have been criticized for excessively  
15 matching students from privileged backgrounds with elite institutions regardless of their academic  
16 potential [25, 67]. These concerns stem from the fact that the matching decision made by the central  
17 algorithm is not independent of a group-defining attribute such as social background. Although group  
18 fairness has been extensively explored in supervised and unsupervised learning settings [4], extending  
19 these notions to matching problems requires novel formal definitions and algorithmic tools.

20 In this work, we study group fairness in matching through optimal transport (OT), which has been a  
21 long-standing tool to model matching problems in economics and social sciences since the pioneering  
22 work of Kantorovitch [32, 48, 9, 54, 63, 39]—see Appendix A for a detailed survey of related work.  
23 Our contribution is twofold. Our first contribution is a formal framework for fair optimal transport  
24 under group fairness constraints, featuring two new fairness definitions designed for this setting.  
25 Secondly, we design two algorithms that rely on penalized optimal transport and bi-level optimization,  
26 respectively, to achieve group fairness in optimal transport. For the first approach, we derive a novel  
27 sample complexity bound, which could be of independent interest for general penalized entropic  
28 optimal transport problems. All proofs are given in the Appendix.

## 29 2 Formalizing Group Fairness in Optimal Transport

### 30 2.1 Optimal Transport

31 **Entropic Optimal Transport.** We consider the entropic regularized optimal transport problem  
32 [61, 34, 27, 62, 44] between two probability distributions  $\mu$  and  $\eta$  with ground cost  $c$ :

$$\mathbf{W}_c^\varepsilon(\mu, \eta) := \min_{\pi \in \Pi(\mu, \eta)} \int_{\mathbb{R}^d \times \mathbb{R}^d} c(x, y) d\pi(x, y) + \varepsilon \mathbf{KL}(\pi | \mu \otimes \eta) ,$$

33 where the minimum is taken over couplings  $\Pi(\mu, \eta)$  of  $\mu$  and  $\eta$  and where  $\mathbf{KL}$  denotes the Kullback-  
 34 Leibler (KL) divergence (see [19]). If  $\mu$  and  $\nu$  are sums of Dirac measures, *i.e.*, if  $\mu = 1/n \sum_{i=1}^n \delta_{\mathbf{x}_i}$   
 35 and  $\eta = 1/m \sum_{j=1}^m \delta_{\mathbf{y}_j}$ , the problem above reduces to

$$\mathbf{W}_C^\varepsilon := \min_{\Pi \in \Pi} \sum_{i=1}^n \sum_{j=1}^m \Pi_{ij} C_{ij} + \varepsilon \mathbf{KL}(\Pi) ,$$

36 where  $\Pi := \{ \Pi \in \mathbb{R}_+^{n \times m} \mid \Pi \mathbf{1}_m = 1/n \mathbf{1}_n, \Pi^\top \mathbf{1}_n = 1/m \mathbf{1}_m \}$ ,  $\mathbf{KL}(\Pi) := \sum_{ij} \Pi_{ij} \log \Pi_{ij}$  and  
 37  $\mathbf{C} \in \mathbb{R}^{n \times m}$  with  $C_{ij} := c(x_i, y_j)$ . The KL divergence is strictly convex, and consequently, we  
 38 denote the corresponding unique minimizer by  $\Pi_\varepsilon(\mathbf{C})$ .

## 39 2.2 Fair Optimal Transport

40 We consider two distributions  $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{S})$  and  $\eta \in \mathcal{P}(\mathcal{Y} \times \mathcal{W})$  to be matched, where  $\mathcal{X}$  and  $\mathcal{Y}$   
 41 represent feature spaces while  $\mathcal{S}$  and  $\mathcal{W}$  correspond to sensitive attributes (e.g., gender, ethnicity or  
 42 age) defining groups of entities in  $\mu$  and  $\eta$  respectively. We restrict ourselves to the case where  $\mathcal{X}$  and  
 43  $\mathcal{Y}$  are compact subspaces of  $\mathbb{R}^k$  for the rest of this paper. Additionally,  $\mathcal{W}$  and  $\mathcal{S}$  are assumed to be  
 44 finite sets which we identify, respectively, with  $\{1, \dots, K_w\}$  and  $\{1, \dots, K_s\}$ . We let  $p_i = \mathbb{P}(S = i)$   
 45 and  $q_j = \mathbb{P}(W = j)$  the probability of each group  $i$  in  $\mu$  and  $j$  in  $\eta$ , respectively, and denote by  $\mathbf{p}$   
 46 and  $\mathbf{q}$  the resulting distributions over  $\mathcal{W}$  and  $\mathcal{S}$ .

47 **Fairness targets.** Optimal transport imposes structure on the coupling of two probability distribu-  
 48 tions through the cost function, making it more likely to match “nearby” points.

49 **Example 2.1.** Consider a school assignment system, in which students are characterized by their  
 50 geographic position  $X \in \mathbb{R}^2$  and a social status  $S \in \{\text{high}, \text{low}\}$ . Similarly, schools are charac-  
 51 terized by geographic positions  $Y \in \mathbb{R}^2$  and a prestige level  $W \in \{\text{elite}, \text{non-elite}\}$ . Suppose that  
 52 students with  $S = \text{high}$  are more likely to live near elite schools, which are often situated in privileged  
 53 neighborhoods, while students with low social status tend to live near non-elite schools. In this  
 54 setting, assigning students to schools using optimal transport with Euclidean cost—minimizing the  
 55 distance between students and schools—will result in highly segregated schools.

56 In this example, the block-sparse structure of the optimal transport plan can be seen as a source of  
 57 unfairness, as the matching will be highly correlated with the social status of the students and the  
 58 elitist nature of the schools. To define a fair transport plan, we assume that we are given a *fairness*  
 59 *target*  $\mathbf{F}$ , which is a  $K_s \times K_w$  matrix that specifies, for each pair of groups  $(s, w) \in \mathcal{S} \times \mathcal{W}$ , the  
 60 desired probability of matching members of group  $s$  with group  $w$ . To be valid, the matrix  $\mathbf{F}$  should  
 61 itself be a coupling between  $\mathbf{p}$  and  $\mathbf{q}$ : it should be non-negative and satisfy the constraints

$$\sum_{w=1}^{K_w} \mathbf{F}_{sw} = p_s \text{ and } \sum_{s=1}^{K_s} \mathbf{F}_{sw} = q_w$$

62 for all  $(s, w)$ . We hence write  $\Pi(\mathbf{p}, \mathbf{q})$  the set of fairness targets.

63 **Example 2.2.** Consider the segregated schooling system introduced in Example 2.1. Let  $p := \mathbb{P}(S =$   
 64  $\text{low})$  and  $q := \mathbb{P}(W = \text{non-elite})$ . A newly appointed city administrator wishes to strongly limit  
 65 social homogamy within public schools, and hence requires 60% of student with low social status to  
 66 be matched to elite schools. The fairness target can be written as

$$\mathbf{F} = \begin{bmatrix} 0.4 \times p & 0.6 \times p \\ q - 0.4 \times p & 1 - q - 0.6 \times p \end{bmatrix}. \quad (1)$$

67 **Cost-Insentive Fairness.** We can now define a notion of *fair coupling* with respect to a given  
 68 fairness target  $\mathbf{F}$ .

69 **Definition 2.3 (F-Fair coupling).** Let  $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{S})$  and  $\nu \in \mathcal{P}(\mathcal{Y} \times \mathcal{W})$ . A coupling  $\pi \in \Pi(\mu, \nu)$   
 70 is said to be **F-fair** if for all  $(s, w) \in \mathcal{S} \times \mathcal{W}$  we have

$$\pi_{SW}(s, w) = \mathbf{F}_{sw} ,$$

71 where  $\pi_{SW}$  is the coupling on  $\mathcal{S} \times \mathcal{W}$  that is obtained from  $\pi$  by marginalizing over  $x$  and  $y$ , that is

$$\pi_{SW}(s, w) := \pi(\mathcal{X} \times \{s\} \times \mathcal{Y} \times \{w\}) = \int_{\mathcal{X} \times \{s\} \times \mathcal{Y} \times \{w\}} d\pi(x, u, y, t) .$$

72 In other words, a coupling is fair if the amount of mass moved from group  $s$  to group  $w$  is equal to  
 73  $\mathbf{F}_{sw}$ . We denote by  $\Pi_{\text{fair}}^{\mathbf{F}}(\mu, \eta)$  the set of  $\mathbf{F}$ -fair couplings and define the corresponding fair optimal  
 74 transport problem as the problem of finding

$$\Pi_{\varepsilon, \text{fair}}^* \in \arg \min_{\pi \in \Pi_{\text{fair}}^{\mathbf{F}}(\mu, \eta)} \int c(x, y) d\pi(x, y) + \varepsilon \mathbf{KL}(\pi | \mu \otimes \eta) . \quad (2)$$

75 **Proposition 2.4.** *Assume that  $\mu$  and  $\eta$  have compact support. Then, for any  $\mathbf{F} \in \Pi(\mathbf{p}, \mathbf{q})$ , there*  
 76 *exists a unique fair optimal transport plan.*

77 **Cost-Sensitive Fairness.** In some scenarios, it is desirable to enforce fairness while taking the  
 78 ground cost  $c$  into account. Instead of requiring that a fraction  $\mathbf{F}_{sw}$  of all matches occur between  
 79 individuals from groups  $s$  and  $w$ , we may require that these groups collectively bear a share  $\mathbf{F}_{sw}$  of  
 80 the total transport cost. This leads to the following definition.

81 **Definition 2.5 (Cost Fair OT).** *Let  $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{S})$  and  $\nu \in \mathcal{P}(\mathcal{Y} \times \mathcal{W})$ . A coupling  $\pi \in \Pi(\mu, \nu)$  is*  
 82 *said to be **F-cost fair** if for all  $(s, w) \in \mathcal{S} \times \mathcal{W}$  we have*

$$\int_{\mathcal{X} \times \{s\} \times \mathcal{Y} \times \{w\}} c(x, y) d\pi(x, u, y, t) = \mathbf{F}_{sw} \bar{c} ,$$

83 where  $\bar{c} := \mathbb{E}_{\mu \otimes \eta}[c(X, Y)]$  is the total cost.

84 **Example 2.6.** *Consider again the segregated school system introduced in previous examples. En-*  
 85 *forcing cost-sensitive fairness in this context involves weighting matches according to the distance*  
 86 *between students and schools, thereby discouraging matches that are overly costly.*

87 **Finite-Sample Fair Optimal Transport.** Consider two datasets  $(\mathbf{x}_i, \mathbf{s}_i)_{i=1}^n \in \mathcal{X} \times \mathcal{S}$  and  
 88  $(\mathbf{y}_i, \mathbf{w}_i)_{i=1}^m \in \mathcal{Y} \times \mathcal{W}$  drawn i.i.d. from  $\mu$  and  $\nu$  respectively. The finite sample version of the  
 89 fair optimal transport (2) corresponds to replacing  $\mu$  and  $\nu$  by their empirical counterparts, that is,  
 90  $\mu_n = 1/n \sum_i \delta_{(x_i, s_i)}$  and  $\nu_n = 1/m \sum_j \delta_{(y_j, w_j)}$ . Let  $\mathbf{B}_{sw} := (\mathbf{1}_{s_i=s} \mathbf{1}_{w_j=w})_{i,j} \in \{0, 1\}^{n \times m}$  and  
 91 the sample problem is thus given by

$$\min_{\Pi \in \Pi} \text{Tr}[\Pi^\top \mathbf{C}] + \varepsilon \mathbf{KL}(\Pi) \quad (3)$$

$$\text{s.t. } \forall (s, w) \in \mathcal{S} \times \mathcal{W}, \text{Tr}[\Pi^\top \mathbf{B}_{sw}] = \sum_{i|\mathbf{s}_i=s} \sum_{j|\mathbf{w}_j=w} \Pi_{ij} = \mathbf{F}_{sw} . \quad (4)$$

92 For cost-sensitive fairness, (4) is replaced by  $\text{Tr}[\mathbf{B}_{sw} \Pi^\top \mathbf{C}] = \bar{c} \mathbf{F}_{sw}$  with  $\bar{c} := (nm)^{-1} \sum \mathbf{C}_{ij}$ .  
 93 This problem can be solved using a modified Sinkhorn algorithm—see Appendix B for details.

### 94 3 Two Strategies for Fair OT

95 **Relaxing Fairness Constraints.** Enforcing exact fairness can lead to transport plans with pro-  
 96 hibitively high cost. In practice, it is often preferable to allow some tolerance to balance the trade-off  
 97 between fairness and overall efficiency. In our setting, we aim to find a transport plan that is fair  
 98 while remaining reasonably close to the original plan. To achieve this, we adopt a  $\rho$ -relaxed fairness  
 99 approach, requiring that for all  $(s, w) \in \mathcal{S} \times \mathcal{W}$ , the fairness violation satisfies  $\mathcal{L}_{\mathbf{F}}(\Pi) \leq \rho$ , where

$$\mathcal{L}_{\mathbf{F}}(\Pi) := \sum_{(s,w) \in \mathcal{S} \times \mathcal{W}} \left( \text{Tr}[\Pi^\top \mathbf{B}_{sw}] - \mathbf{F}_{sw} \right)^2, \quad (\rho\text{-relaxed cost-insensitive fairness}) \quad (5)$$

or

$$\mathcal{L}_{\mathbf{F}}(\Pi) := \sum_{(s,w) \in \mathcal{S} \times \mathcal{W}} \left( \text{Tr}[\mathbf{B}_{sw} \Pi^\top \mathbf{C}] - \bar{c} \mathbf{F}_{sw} \right)^2, \quad (\rho\text{-relaxed cost-sensitive fairness}) \quad (6)$$

100 and  $\rho \geq 0$  is a tolerance level that leads to a relaxation of (3)-(4) where (4) is replaced by  $\mathcal{L}_{\mathbf{F}}(\Pi) \leq \rho$ .  
 101 To find optimal transport plans under relaxed fairness constraints, we propose two strategies.

102 **3.1 Fairness-Penalized OT**

103 Our first strategy is a direct penalization strategy. We solve a penalized version of optimal transport  
 104 problem under  $\rho$ -relaxed fairness by introducing a Lagrange multiplier associated to the constraint  
 105  $\mathcal{L}_{\mathbf{F}}(\Pi) \leq \rho$ :

$$\text{OT}_{\mathbf{F}}^{\varepsilon} := \min_{\Pi \in \Pi} \text{Tr}[\Pi^{\top} \mathbf{C}] + \varepsilon \mathbf{KL}(\Pi) + \lambda \mathcal{L}_{\mathbf{F}}(\Pi) , \quad (7)$$

106 where  $\lambda, \varepsilon > 0$  control regularization strength and  $\mathcal{L}_{\mathbf{F}}(\Pi)$  is one of the two fairness penalties intro-  
 107 duced in Equations (5)-(6), which may depend on the cost. This is an instance of penalized entropic opti-  
 108 mal transport which can be solved efficiently using a generalized conditional gradient algorithm [64].

109 **Sample Complexity of Penalized OT.** One of the main contributions of our work is to establish  
 110 a sample complexity bound for the fairness penalized entropic optimal transport cost, by building  
 111 upon results from entropic optimal transport [34, 55]. To formalize this, note that the optimization  
 112 problem in (7) can be defined for arbitrary measures  $\alpha \in \mathcal{P}(\mathcal{X} \times 0, 1)$  and  $\beta \in \mathcal{P}(\mathcal{Y} \times 0, 1)$  as  
 113 (see Appendix D.2 for a formal definition)

$$m^*(\alpha, \beta) := \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathbb{R}^d \times \mathbb{R}^d} c(x, y) d\pi(x, y) + \varepsilon \mathbf{KL}(\pi | \alpha \otimes \beta) + \lambda \mathcal{L}_{\mathbf{F}}(\pi) . \quad (8)$$

114 For simplicity, we restrict to the simple setting where which  $K_s = K_w = 2$  and the fairness penalty  
 115 is chosen to be cost-insensitive. We note, however, that our proof extends naturally to the more  
 116 general case. Similar to [65] and [34], our proof relies on the fact that the measures are compactly  
 117 supported. We also require the following assumption on the ground cost function.

118 **Assumption 1.** The cost function  $c : (\mathcal{X} \times \{0, 1\}) \times (\mathcal{Y} \times \{0, 1\}) \rightarrow \mathbb{R}_+$  is Lipschitz and  $\mathcal{C}^{\infty}$ .

**Theorem 3.1.** Assume for simplicity that  $m = n$ . Under Assumption 1, we have

$$\mathbb{E} |m^*(\mu_n, \eta_n) - m^*(\mu, \eta)| \lesssim \frac{\log(n)}{\sqrt{n}} .$$

119 Comparing our result with existing sample complexity bounds for the *unpenalized* entropic OT  
 120 problem [34, 55], we achieve similar bounds up to a logarithmic factor.

121 **Remark 1.** In the formulation above, the cost function is defined on inputs from  $\mathcal{X} \times 0, 1$  and  $\mathcal{Y} \times 0, 1$   
 122 for full generality. By setting the cost function  $c$  to be independent of the sensitive attribute, we  
 123 recover our original setting.

124 **Proof sketch.** Our proof works by finding adequate random variables  $Y_n, Z_n$  such that

$$Y_n - m^*(\mu, \eta) \leq m^*(\mu_n, \eta_n) - m^*(\mu, \eta) \leq Z_n - m^*(\mu, \eta) . \quad (9)$$

125 *Lower bound.* We obtain  $Y_n$  through linearization of the convex penalty, which is guaranteed to yield  
 126 a proper lower bound through the connection between the optimality conditions of our penalized  
 127 problem and its linearized counterpart [64]. This linearized problem may then be cast as an entropic  
 128 optimal transport problem with modified cost, allowing us to leverage standard results on the sample  
 129 complexity of entropic optimal transport [33, 55].

130 *Upper bound.* To obtain  $Z_n$ , we proceed in two steps. First, we evaluate the optimal population  
 131 coupling  $\pi_*$  on the dataset, that is, we consider the random coupling

$$\widehat{\pi}_n^* := \frac{1}{n^2} \sum_{ij} \frac{d\pi^*}{d\mu \otimes \nu}((x_i, s_i), (y_j, w_j)) \delta_{((x_i, s_i), (y_j, w_j))} \quad (10)$$

132 and let  $\widehat{Z}_n := \int c d\widehat{\pi}_n^* + \varepsilon \mathbf{KL}(\widehat{\pi}_n^* | \mu_n \otimes \nu_n) + \lambda \mathcal{L}_{\mathbf{F}}(\widehat{\pi}_n^*)$ . Observe that  $\widehat{\pi}_n^*$  does not necessarily  
 133 satisfy the marginal constraints that define  $\Pi(\mu_n, \eta_n)$ , and, consequently,  $\widehat{Z}_n$  is not necessarily  
 134 an upper bound of  $m_n^*$ . We thus secondly project  $\widehat{\pi}_n^*$  onto the constraint set. We use the round

135 algorithm proposed by Altschuler et al. [1, Algorithm 2] (see Appendix D.2.2). The rounded coupling  
 136  $\bar{\pi}_n^* = \text{round}(\hat{\pi}_n^*)$  now is a valid coupling of the empirical distributions, and we hence use

$$Z_n := \int c d\bar{\pi}_n^* + \varepsilon \mathbf{KL}(\bar{\pi}_n^* || \mu_n \otimes \nu_n) + \mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) . \quad (11)$$

137 We conclude the proof by using  $|Z_n - m_\infty^*| \leq |Z_n - \hat{Z}_n| + |\hat{Z}_n - m^*(\mu, \nu)|$ . The first term can be  
 138 bounded using regularity arguments and concentration of the projected coupling, and the second term  
 139 using  $\mathbb{E}[\hat{Z}_n] = m^*(\mu, \nu)$ , independence of samples, and compactness arguments.

### 140 3.2 Bilevel Optimization Formulation

141 A second strategy builds on cost-learning in optimal transport [49] and consists in solving the bilevel  
 142 optimization problem

$$\min_{\theta \in \Theta} \mathcal{L}_{\mathbf{F}}(\mathbf{\Pi}(c_\theta), \mathbf{C}, \Gamma) + \lambda^{-1} \mathcal{D}(c_\theta, c_{\text{base}}) \text{ s.t. } \mathbf{\Pi}(c_\theta) = \arg \min_{\mathbf{\Pi}} \text{Tr}[\mathbf{\Pi}^\top \mathbf{C}_\theta] + \varepsilon \mathbf{KL}(\mathbf{\Pi}) \quad (12)$$

143 where  $\mathbf{C}_\theta := [c_\theta(\mathbf{x}_i, \mathbf{y}_j)]_{ij}$  denotes a parameterized family of cost matrices,  $c_{\text{base}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}+$   
 144 is a user-specified baseline cost function, and  $\mathcal{D}$  is a discrepancy measure between cost functions that  
 145 encourages the learned cost to remain close to the baseline. Problem (12) can be efficiently solved  
 146 using gradient-based methods, as its inner component is an entropically regularized optimal transport  
 147 problem that is strongly convex over the probability simplex. Typical approaches rely on iterative  
 148 differentiation [7, 30, 53, 59] or (approximate) implicit differentiation [20, 21, 24, 35, 47, 60], and  
 149 can accommodate a wide variety of parameterized costs. Relevant examples include Mahalanobis  
 150 distances,  $c_{\mathbf{M}}(x, y) = (x - y)^\top \mathbf{M}(x - y)$  where  $\mathbf{M}$  is a PSD matrix, and their nonlinear counterpart,  
 151  $c_\theta(x, y) = \|\phi_\theta(x) - \phi_\theta(y)\|_2^2$ , where  $\phi_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$  is a neural network with parameters  $\theta$ .

152 Our bilevel cost-learning approach allows for high flexibility and modularity and yields an inter-  
 153 pretable cost function, which can be readily reused on new samples. It however requires to solve an  
 154 intricate non-convex problem and comes with no theoretical guarantees, in contrast to the penalized  
 155 approach. We provide a thorough comparison of both approaches in Appendix C.

## 156 4 Numerical Illustrations

157 We illustrate our algorithms using simulated mixture of Gaussian data, focusing on cost-sensitive  
 158 fairness. Figure 1 displays our results for the penalized approach. The left panel shows how the  
 159 fairness loss evolves as the penalty increases. We observe that our matching gets increasingly fair  
 160 with larger penalties. The center and right panels show the fairness loss as a function function of the  
 161 cost difference and KL divergence, respectively, between the fair plan and the unfair plan (that is the  
 162 plan computed with  $\lambda = 0$ ). One notices that lower fairness loss corresponds to higher bias with  
 163 respect to the unfair optimal transport plan. We conduct the same experiment with the cost learning  
 164 approach. The corresponding results—reported in Appendix E—yield similar conclusions.

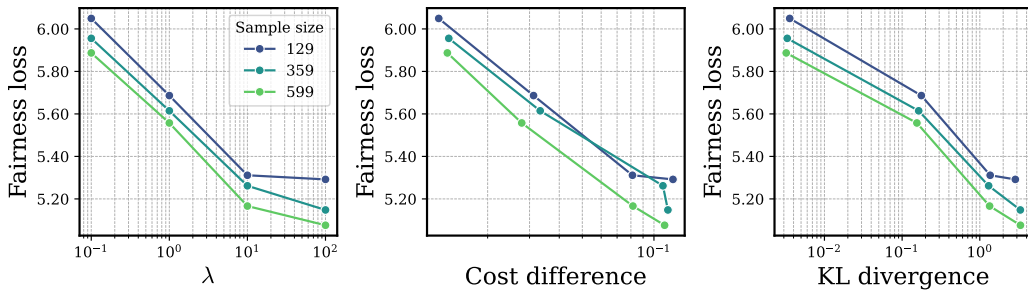


Figure 1: Fairness vs. penalty (**left**), fairness-cost difference (**center**) and fairness-KL divergence (**right**) trade-offs (w.r.t. the non-penalized problem) with varying sample sizes for the **penalized OT problem** with cost-sensitive fairness loss.

## References

- 165
- 166 [1] J. Altschuler, J. Niles-Weed, and P. Rigollet. Near-linear time approximation algorithms for  
167 optimal transport via sinkhorn iteration. *Advances in neural information processing systems*,  
168 30, 2017. [5](#), [15](#), [16](#)
- 169 [2] I. Ashlagi and A. E. Roth. Free riding and participation in large scale, multi-hospital kidney  
170 exchange. *Theoretical Economics*, 9(3):817–863, 2014. [10](#)
- 171 [3] S. Baharlouei, M. Nouiehed, A. Beirami, and M. Razaviyayn. R\`enyi fair inference. *arXiv*  
172 *preprint arXiv:1906.12005*, 2019. [10](#)
- 173 [4] S. Barocas, M. Hardt, and A. Narayanan. *Fairness and Machine Learning: Limitations and*  
174 *Opportunities*. MIT Press, 2023. [1](#)
- 175 [5] J.-D. Benamou, G. Carlier, M. Cuturi, L. Nenna, and G. Peyré. Iterative bregman projections  
176 for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–  
177 A1138, 2015. [12](#)
- 178 [6] M. Blondel, V. Seguy, and A. Rolet. Smooth and sparse optimal transport. In *International*  
179 *conference on artificial intelligence and statistics*, pages 880–889. PMLR, 2018. [10](#)
- 180 [7] J. Bolte, E. Pauwels, and S. Vaiter. One-step differentiation of iterative algorithms. In *Advances*  
181 *in Neural Information Processing Systems (NeurIPS)*, 2023. [5](#)
- 182 [8] T. Calders, F. Kamiran, and M. Pechenizkiy. Building classifiers with independency constraints.  
183 In *2009 IEEE international conference on data mining workshops*, pages 13–18. IEEE, 2009.  
184 [10](#)
- 185 [9] G. Carlier, A. Dupuy, A. Galichon, and Y. Sun. Sista: learning optimal transport costs under  
186 sparsity constraints. *Communications on Pure and Applied Mathematics*, 76(9):1659–1677,  
187 2023. [1](#)
- 188 [10] R. Castera. *Efficiency and Fairness in Matching Problems*. PhD thesis, UGA (Université  
189 Grenoble Alpes), 2024. [10](#)
- 190 [11] R. Castera, P. Loiseau, and B. S. Pradelski. Statistical discrimination in stable matchings. In  
191 *Proceedings of the 23rd ACM Conference on Economics and Computation*, pages 373–374,  
192 2022. [10](#)
- 193 [12] R. Castera, F. Garrido-Lucero, M. Molina, S. Murras, P. Loiseau, and V. Perchet. The price of  
194 fairness in bipartite matching. *arXiv preprint arXiv:2403.00397*, 2024. [10](#)
- 195 [13] Y.-K. Che and O. Tercieux. Efficiency and stability in large matching markets. *Journal of*  
196 *Political Economy*, 127(5):2301–2342, 2019. [10](#)
- 197 [14] J. Cho, G. Hwang, and C. Suh. A fair classifier using mutual information. In *2020 IEEE*  
198 *international symposium on information theory (ISIT)*, pages 2521–2526. IEEE, 2020. [10](#)
- 199 [15] S. Chowdhary, G. De Pasquale, N. Lanzetti, A.-A. Stoica, and F. Dorfler. Fairness in social  
200 influence maximization via optimal transport. *Advances in Neural Information Processing*  
201 *Systems*, 37:10380–10413, 2024. [10](#)
- 202 [16] E. Chzhen, C. Denis, M. Hebiri, L. Oneto, and M. Pontil. Fair regression with wasserstein  
203 barycenters. *Advances in Neural Information Processing Systems*, 33:7321–7331, 2020. [10](#)
- 204 [17] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy. Optimal transport for domain  
205 adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865,  
206 2016. [10](#)
- 207 [18] M.-P. Couto, F. Bugeja-Bloch, and L. Frouillou. Parcoursup: les prémices d’un accroissement  
208 de la stratification sociale et scolaire des formations du supérieur. *Agora débats/jeunesses*, 89  
209 (3):23–38, 2021. [10](#)

- 210 [19] I. Csiszár.  $i$ -divergence geometry of probability distributions and minimization problems.  
 211 *The Annals of Probability*, 3(1):146–158, 1975. ISSN 00911798, 2168894X. URL <http://www.jstor.org/stable/2959270>. 2, 18  
 212
- 213 [20] M. Dagr  ou, P. Ablin, S. Vaiter, and T. Moreau. A framework for bilevel optimization that  
 214 enables stochastic and global variance reduction algorithms. In *Advances in Neural Information*  
 215 *Processing Systems (NeurIPS)*, 2022. 5
- 216 [21] M. Dagr  ou, T. Moreau, S. Vaiter, and P. Ablin. A Lower Bound and a Near-Optimal Algorithm  
 217 for Bilevel Empirical Risk Minimization. In *International Conference on Artificial Intelligence*  
 218 *and Statistics (AISTATS)*, 2024. 5
- 219 [22] S. Devic, D. Kempe, V. Sharan, and A. Korolova. Fairness in matching under uncertainty. In  
 220 *International Conference on Machine Learning*, pages 7775–7794. PMLR, 2023. 10
- 221 [23] V. Divol and S. Gaucher. Demographic parity in regression and classification within the  
 222 unawareness framework. *arXiv preprint arXiv:2409.02471*, 2024. 10
- 223 [24] M. Eisenberger, A. Toker, L. Leal-Taixe, F. Bernard, and D. Cremers. A Unified Framework  
 224 for Implicit Sinkhorn Differentiation. In *2022 IEEE/CVF Conference on Computer Vision and*  
 225 *Pattern Recognition (CVPR)*, pages 499–508, New Orleans, LA, USA, 2022. IEEE. 5
- 226 [25] G. Fack, J. Grenet, and A. Benhenda. L’impact des proc  dures de sectorisation et d’affectation  
 227 sur la mixit   sociale et scolaire dans les lyc  es d’  le-de-france. 2014. 1
- 228 [26] E. Fehr, A. Klein, and K. M. Schmidt. Fairness and contract design. *Econometrica*, 75(1):  
 229 121–154, 2007. 10
- 230 [27] J. Feydy, T. S  journ  , F.-X. Vialard, S.-i. Amari, A. Trouv  , and G. Peyr  . Interpolating between  
 231 optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference*  
 232 *on Artificial Intelligence and Statistics*, pages 2681–2690. PMLR, 2019. 1
- 233 [28] R. Flamary, N. Courty, A. Gramfort, M. Z. Alaya, A. Boisbunon, S. Chambon, L. Chapel,  
 234 A. Corenflos, K. Fatras, N. Fournier, L. Gautheron, N. T. Gayraud, H. Janati, A. Rakotomamonjy,  
 235 I. Redko, A. Rolet, A. Schutz, V. Seguy, D. J. Sutherland, R. Tavenard, A. Tong, and T. Vayer.  
 236 Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021. URL  
 237 <http://jmlr.org/papers/v22/20-451.html>. 20
- 238 [29] R. Flamary, C. Vincent-Cuaz, N. Courty, A. Gramfort, O. Kachaiev, H. Quang Tran, L. David,  
 239 C. Bonet, N. Cassereau, T. Gnassounou, E. Tanguy, J. Delon, A. Collas, S. Mazelet, L. Chapel,  
 240 T. Kerdoncuff, X. Yu, M. Feickert, P. Krzakala, T. Liu, and E. Fernandes Montesuma. Pot  
 241 python optimal transport (version 0.9.5), 2024. URL <https://github.com/PythonOT/POT>.  
 242 20
- 243 [30] L. Franceschi, M. Donini, P. Frasconi, and M. Pontil. Forward and Reverse Gradient-Based  
 244 Hyperparameter Optimization. In *International Conference on Machine Learning (ICML)*, 2017.  
 245 5
- 246 [31] D. Gale and L. S. Shapley. College admissions and the stability of marriage. *The American*  
 247 *mathematical monthly*, 69(1):9–15, 1962. 10
- 248 [32] A. Galichon. *Optimal transport methods in economics*. Princeton University Press, 2018. 1
- 249 [33] A. Genevay. *Entropy-regularized optimal transport for machine learning*. PhD thesis, Universit    
 250 Paris sciences et lettres, 2019. 4, 15
- 251 [34] A. Genevay, L. Chizat, F. Bach, M. Cuturi, and G. Peyr  . Sample complexity of sinkhorn  
 252 divergences. In *The 22nd international conference on artificial intelligence and statistics*, pages  
 253 1574–1583. PMLR, 2019. 1, 4
- 254 [35] S. Ghadimi and M. Wang. Approximation Methods for Bilevel Programming. *arXiv preprint*  
 255 *arXiv:1802.02246*, 2018. 5

- 256 [36] P. Gordaliza, E. Del Barrio, G. Fabrice, and J.-M. Loubes. Obtaining fairness using optimal  
257 transport theory. In *International conference on machine learning*, pages 2357–2365. PMLR,  
258 2019. 10
- 259 [37] T. L. Gouic, J.-M. Loubes, and P. Rigollet. Projection to fairness in statistical learning. *arXiv*  
260 *preprint arXiv:2005.11720*, 2020. 10
- 261 [38] M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. *Advances in*  
262 *neural information processing systems*, 29, 2016. 10
- 263 [39] Y. Hazard and T. Kitagawa. Who with whom? learning optimal matching policies. *arXiv*  
264 *preprint arXiv:2507.13567*, 2025. 1
- 265 [40] B. T. Hirsch and E. J. Schumacher. Match bias in wage gap estimates due to earnings imputation.  
266 *Journal of labor economics*, 22(3):689–722, 2004. 10
- 267 [41] F. Hu, P. Ratz, and A. Charpentier. Fairness in multi-task learning via wasserstein barycenters.  
268 In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*,  
269 pages 295–312. Springer, 2023. 10
- 270 [42] Y. Kamada and F. Kojima. Fair matching under constraints: Theory and applications. *Review of*  
271 *Economic Studies*, 91(2):1162–1199, 2024. 10
- 272 [43] G. Karni, G. N. Rothblum, and G. Yona. On fairness and stability in two-sided matchings. *arXiv*  
273 *preprint arXiv:2111.10885*, 2021. 10
- 274 [44] N. Keriven. Entropic optimal transport in random graphs. *arXiv preprint arXiv:2201.03949*,  
275 2022. 1
- 276 [45] T. König, L. Mechtenberg, D. Kübler, and R. Schmacker. Fairness in matching markets:  
277 Experimental evidence. Technical report, WZB Discussion Paper, 2023. 10
- 278 [46] J. Korman and R. McCann. Optimal transportation with capacity constraints. *Transactions of*  
279 *the American Mathematical Society*, 367(3):1501–1521, 2015. 10
- 280 [47] J. Larsen, L. Hansen, C. Svarer, and M. Ohlsson. Design and regularization of neural networks:  
281 The optimal use of a validation set. In *Neural Networks for Signal Processing VI. Proceedings*  
282 *of the 1996 IEEE Signal Processing Society Workshop*, pages 62–71, Kyoto, Japan, 1996. IEEE.  
283 5
- 284 [48] D. Leite and C. De Bacco. Urban transportation networks and optimal transport-based infras-  
285 tructures: similarity and economy of scale. *arXiv preprint arXiv:2209.06751*, 2022. 1
- 286 [49] R. Li, X. Ye, H. Zhou, and H. Zha. Learning to match via inverse optimal transport. *Journal*  
287 *of Machine Learning Research*, 20(80):1–37, 2019. URL [http://jmlr.org/papers/v20/](http://jmlr.org/papers/v20/li18-700.html)  
288 [18-700.html](http://jmlr.org/papers/v20/li18-700.html). 5
- 289 [50] T. Liu, J. Puigcerver, and M. Blondel. Sparsity-constrained optimal transport. *arXiv preprint*  
290 *arXiv:2209.15466*, 2022. 10
- 291 [51] E. Lobo, J. Payan, C. Cousins, and Y. Zick. Fair and welfare-efficient constrained multi-  
292 matchings under uncertainty. *Advances in Neural Information Processing Systems*, 37:74579–  
293 74616, 2025. 10
- 294 [52] A. Lowy, S. Baharlouei, R. Pavan, M. Razaviyayn, and A. Beirami. A stochastic optimization  
295 framework for fair risk minimization. *arXiv preprint arXiv:2102.12586*, 2021. 10
- 296 [53] D. Maclaurin, D. Duvenaud, and R. P. Adams. Gradient-based Hyperparameter Optimization  
297 through Reversible Learning. In *International Conference on Machine Learning (ICML)*, 2015.  
298 5
- 299 [54] R. Mastrandrea, P. Pagnottoni, N. Pecora, and A. Spelta. An optimal transport approach to  
300 model the community structure of the international trade network. *Social Networks*, 82:111–133,  
301 2025. 1

- 302 [55] G. Mena and J. Niles-Weed. Statistical bounds for entropic optimal transport: sample complexity  
303 and the central limit theorem. *Advances in neural information processing systems*, 32, 2019. 4
- 304 [56] K. Nguyen, H. Nguyen, and N. Ho. Towards marginal fairness sliced wasserstein barycenter.  
305 *arXiv preprint arXiv:2405.07482*, 2024. 10
- 306 [57] M. Nutz. Introduction to entropic optimal transport. *Lecture notes, Columbia University*, 2021.  
307 19
- 308 [58] F.-P. Paty and M. Cuturi. Subspace robust wasserstein distances. In *International conference on*  
309 *machine learning*, pages 5072–5081. PMLR, 2019. 10
- 310 [59] E. Pauwels and S. Vaiter. The Derivatives of Sinkhorn–Knopp Converge. *SIAM Journal on*  
311 *Optimization*, 33(3):1494–1517, 2023. ISSN 1052-6234, 1095-7189. 5
- 312 [60] F. Pedregosa. Hyperparameter optimization with approximate gradient. In *International*  
313 *Conference on Machine Learning (ICML)*, 2016. 5
- 314 [61] G. Peyré, M. Cuturi, et al. Computational optimal transport: With applications to data science.  
315 *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019. 1
- 316 [62] A.-A. Pooladian and J. Niles-Weed. Entropic estimation of optimal transport maps. *arXiv*  
317 *preprint arXiv:2109.12004*, 2021. 1
- 318 [63] Z. Qu, A. Galichon, W. Gao, and J. Ugander. On sinkhorn’s algorithm and choice modeling.  
319 *Operations Research*, 2025. 1
- 320 [64] A. Rakotomamonjy, R. Flamary, and N. Courty. Generalized conditional gradient: analysis of  
321 convergence and applications. *arXiv preprint arXiv:1510.06567*, 2015. 4, 10, 17
- 322 [65] P. Rigollet and A. J. Stromme. On the sample complexity of entropic optimal transport. *arXiv*  
323 *preprint arXiv:2206.13472*, 2022. 4, 16, 20
- 324 [66] F. Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94,  
325 2015. 13
- 326 [67] M. Simioni and P. Steiner. *Comment ça matche: Une sociologie de l’appariement*. Presses de  
327 Sciences Po, 2022. 1
- 328 [68] T. Sühr, A. J. Biega, M. Zehlike, K. P. Gummedi, and A. Chakraborty. Two-sided fairness  
329 for repeated matchings in two-sided markets: A case study of a ride-hailing platform. In  
330 *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery &*  
331 *data mining*, pages 3082–3092, 2019. 10
- 332 [69] Z. Xiong, N. Dalmaso, S. Sharma, F. Lecue, D. Magazzeni, V. Potluru, T. Balch, and M. Veloso.  
333 Fair wasserstein coresets. *Advances in Neural Information Processing Systems*, 37:132–168,  
334 2024. 10
- 335 [70] M. Zhang, X. Dai, and L. Li. Fairness-aware organ exchange and kidney paired donation. *arXiv*  
336 *preprint arXiv:2503.06431*, 2025. 10

337 **A Related Work**

338 **Fairness in Supervised Learning.** To model fairness issues, one usually considers in supervised  
339 learning datapoints of the form  $(X, Y, S) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{S}$ , where  $X$  represents an individuals' features,  
340  $Y$  is an outcome and  $S$  is a sensitive attribute — such as gender, age or income. In its most common  
341 sens, statistical fairness is defined through independence between the sensitive attribute  $S$  and the  
342 output of an algorithm  $f$ , that is an algorithm is fair if  $f(X) = \widehat{Y} \perp\!\!\!\perp S$ . To solve this problem, the  
343 literature on fair supervised learning has considered mutual information based penalties to enforce  
344 independence between the prediction and the sensitive attribute [3, 14, 52], and has also devised  
345 various criteria to be minimized as fairness objectives such as demographic parity (DP) [8, 23] or  
346 equality of odds (EO) [38].

347 **Fairness in Matching Mechanisms.** Since the pioneering work of Gale and Shapley [31] on  
348 the stable marriage problem, a series of works have provided theoretical studies of the fairness of  
349 individual preference-based matching mechanisms [2, 43, 11, 22, 10, 12, 51, 70]. In the same time,  
350 a growing body of work has provided empirical assessment of fairness in matching mechanisms  
351 [40, 26, 68, 13, 18, 45, 42].

352 **Fairness and Optimal Transport.** Recent work has drawn many fruitful connection between  
353 fairness and optimal transport, mostly by leveraging OT as a tool for obtaining or characterizing  
354 fair algorithms [36, 37, 41, 15, 69, 23] — one of the key insights of this line of research being that  
355 the problem of finding a fair predictor can be framed as a Wasserstein barycenter problem [37, 16].  
356 Closest to our work is the recent article by Nguyen et al. [56], who consider the constrained problem  
357 of finding a Wasserstein barycenter between multiple distributions while controlling the pairwise  
358 differences in distances to the barycenter between distributions. Crucially, their setup does not involve  
359 any sensitive attribute. We address a different problem, namely obtaining transport plans under mass  
360 constraints between groups defined through sensitive attributes. To the best of our knowledge, this  
361 problem is still unstudied.

362 **Constrained Optimal Transport.** Our work borrows from recent developments in constrained  
363 optimal transport, which seeks to enforce structural properties on the transportation plan [17, 6, 58,  
364 50]; for instance, Korman and McCann [46] analyze a variant of optimal transport in which the  
365 amount of mass that can be transported between two units is upper bounded. In particular, we build  
366 on Rakotomamonjy et al. [64] to obtain finite sample guarantees for a penalized optimal transport  
367 problem.

368 **B A Modified Sinkhorn Algorithm for Fair Optimal Transport**

369 The following is an immediate consequence of first order conditions on the Lagrangian of Problem 3.

**Proposition B.1.** There exists  $\mathbf{f} \in \mathbb{R}^n$ ,  $\mathbf{g} \in \mathbb{R}^m$  and  $\mathbf{h} = [h_{sw}]_{sw} \in \mathbb{R}^{K_s \times K_w}$  such that the solution to Problem 3 has the form

$$\mathbf{\Pi} = \text{diag}(e^{\mathbf{f}/\varepsilon})(\mathbf{K} \odot \mathbf{H})\text{diag}(e^{\mathbf{g}/\varepsilon})$$

where

$$\mathbf{K} := [e^{-\mathbf{C}_{ij}/\varepsilon}]_{ij} \text{ and } \mathbf{H} := \sum_{sw} e^{h_{sw}/\varepsilon} \mathbf{B}_{sw}.$$

370 **Remark 2.** A similar formulation holds for the cost sensitive problem.

371 *Proof.* Introducing dual variables  $\mathbf{f} \in \mathbb{R}^n$ ,  $\mathbf{g} \in \mathbb{R}^m$  and  $\mathbf{H} = (h_{sw})_{sw} \in \mathbb{R}^{K_s \times K_w}$ , the Lagrangian  
372 writes

$$\begin{aligned} \mathcal{E}(\mathbf{\Pi}, \mathbf{f}, \mathbf{g}, \mathbf{h}) &= \text{Tr}[\mathbf{\Pi}^\top \mathbf{C}] + \varepsilon \mathbf{KL}(\mathbf{\Pi}) - \mathbf{f}^\top (\mathbf{\Pi} \mathbf{1}_m - \mathbf{a}) - \mathbf{g}^\top (\mathbf{\Pi}^\top \mathbf{1}_n - \mathbf{b}) \\ &\quad - \sum_{s,w} h_{sw} [\text{Tr}[\mathbf{\Pi}^\top \mathbf{B}_{sw}] - p_s q_w]. \end{aligned}$$

373 First order conditions yield for every  $i, j$

$$\frac{\partial \mathcal{E}(\mathbf{\Pi}, \mathbf{f}, \mathbf{g}, \mathbf{h})}{\partial \mathbf{\Pi}_{ij}} = \mathbf{C}_{ij} + \varepsilon \log(\mathbf{\Pi}_{ij}) + \varepsilon - \mathbf{f}_i - \mathbf{g}_j - \sum_{sw} h_{sw} [\mathbf{B}_{sw}]_{ij} = 0$$

374 which we may rewrite as

$$\begin{aligned} \mathbf{\Pi}_{ij} &= \exp(\mathbf{f}_i/\varepsilon) \exp\left(-\varepsilon^{-1} \mathbf{C}_{ij} + 1 \varepsilon^{-1} \sum_{sw} h_{sw} [\mathbf{B}_{sw}]_{ij} - 1\right) \exp(\mathbf{g}_j/\varepsilon) \\ &= \exp(\mathbf{f}_i/\varepsilon) \exp(-\mathbf{C}_{ij}/\varepsilon - 1) \prod_{s'=1}^{K_s} \prod_{w'=1}^{K_w} \exp(h_{s'w'} [\mathbf{B}_{s'w'}]_{ij}/\varepsilon) \exp(\mathbf{g}_j/\varepsilon). \end{aligned}$$

375 Now, remark that for any  $(i, j) \in \{1, \dots, n\}^2$ , there is only one  $(s, w) \in \{1, \dots, K_s\} \times \{1, \dots, K_w\}$   
376 such that  $[\mathbf{B}_{sw}]_{i,j} \neq 0$ . Therefore, in the product  $\prod_{s'=1}^{K_s} \prod_{w'=1}^{K_w} \exp(h_{s'w'} [\mathbf{B}_{s'w'}]_{ij})$ , there is only  
377 one term distinct from 1. As a consequence

$$\begin{aligned} \prod_{s'=1}^{K_s} \prod_{w'=1}^{K_w} \exp(h_{s'w'} [\mathbf{B}_{s'w'}]_{ij}/\varepsilon) &= \prod_{s'=1}^{K_s} \prod_{w'=1}^{K_w} \exp(h_{s'w'} \mathbf{1}_{\mathbf{s}_i=s'} \mathbf{1}_{\mathbf{w}_j=w'}/\varepsilon) \\ &= \sum_{s'=1}^{K_s} \sum_{w'=1}^{K_w} \exp(h_{s'w'}/\varepsilon) \mathbf{1}_{\mathbf{s}_i=s'} \mathbf{1}_{\mathbf{w}_j=w'} \\ &= \sum_{s'=1}^{K_s} \sum_{w'=1}^{K_w} \exp(h_{s'w'}/\varepsilon) [\mathbf{B}_{s'w'}]_{ij}. \end{aligned}$$

378 This leads to the matrix form

$$\mathbf{\Pi} = \text{diag}(e^{\mathbf{f}/\varepsilon})(\mathbf{K} \odot \mathbf{H})\text{diag}(e^{\mathbf{g}/\varepsilon})$$

379 where

$$\begin{aligned} \mathbf{K} &:= e^{-\mathbf{C}/\varepsilon - 1} \\ \mathbf{H} &:= \sum_{s=1}^{K_s} \sum_{w=1}^{K_w} e^{h_{sw}/\varepsilon} \mathbf{B}_{sw}. \end{aligned}$$

380

□

|                                  | Penalized OT | Fair Cost Learning |
|----------------------------------|--------------|--------------------|
| Convex ?                         | ✓            | ✗                  |
| Finite sample bounds ?           | ✓            | ✗                  |
| Reusable ?                       | ✗            | ✓                  |
| Interpretability and flexibility | ✗            | ✓                  |

Figure 2: Comparison of both approaches.

A fair optimal transport plan can hence be obtained using a modified Sinkhorn algorithm which we call FairSinkhorn. Define the map

$$\Phi(\mathbf{u}, \mathbf{v}) := \left( \text{Tr} \left[ \left( [\mathbf{u}_i \mathbf{v}_j]_{ij} \odot \mathbf{B}_{sw} \right) \left( \mathbf{K} \odot \mathbf{B}_{sw} \right)^\top \right] \right)_{sw}.$$

381 The algorithm FairSinkhorn combines projections on the set of marginal constraints with a [projec-](#)  
382 [tion on the set of fairness constraints, highlighted in blue.](#)

---

**Algorithm 1** FairSinkhorn Algorithm

---

- 1: **Inputs:** Cost  $\mathbf{C} \in \mathbb{R}^{n \times m}$ ,  $\varepsilon > 0$ , marginals  $(\mathbf{a}, \mathbf{b}) \in \mathbb{R}_+^n \times \mathbb{R}_+^m$ , fairness target  $\mathbf{F} = (\mathbf{F}_{sw})_{sw} \in [0, 1]^{K_s \times K_w}$ , max iter.  $T$ .
- 2:  $\mathbf{K} \leftarrow e^{-\mathbf{C}/\varepsilon-1}$
- 3: **for**  $t = 1, \dots, T$  **do**
- 4:  $\mathbf{u}^{(t+1)} \leftarrow \mathbf{a} \oslash \left( (\mathbf{K} \odot \mathbf{T}^{(t)}) \mathbf{v}^{(t)} \right)$
- 5:  $\mathbf{v}^{(t+1)} \leftarrow \mathbf{b} \oslash \left( (\mathbf{K} \odot \mathbf{T}^{(t)})^\top \mathbf{u}^{(t+1)} \right)$
- 6:  $\mathbf{L}^{(t+1)} \leftarrow \mathbf{F} \oslash \Phi(\mathbf{u}^{(t+1)}, \mathbf{v}^{(t+1)})$
- 7:  $\mathbf{T}^{(t+1)} \leftarrow \sum_{sw} \ell_{sw}^{(t+1)} \mathbf{B}_{sw}$  w.  $\mathbf{L}^{(t+1)} = (\ell_{sw}^{(t+1)})_{sw}$
- 8: **end for**
- 9: **return**

$$\mathbf{\Pi} = \text{diag}(\mathbf{u}^{(T+1)}) (\mathbf{K} \odot \mathbf{T}^{(T+1)}) \text{diag}(\mathbf{v}^{(T+1)})$$


---

383 **Remark 3.** Since FairSinkhorn alternates between projections on sets defined by linear constraints,  
384 it is guaranteed to converge [5].

## 385 C Which approach should one choose ?

386 **Convexity.** While the penalized OT problem is convex, the fair cost learning approach requires  
387 bi-level optimization — hence, the solution is not guaranteed to be unique and algorithms might not  
388 converge easily.

389 **Theoretical Guarantees.** Our first approach comes with strong theoretical guarantees, which ensure  
390 a convergence towards the value of the population problem at rate  $n^{-1/2}$ . The rate of convergence of  
391 our second approach remains an open problem.

392 **Reusability.** When using the first approach, one need to solve a new OT problem for every new  
393 batch of samples. Interestingly, once a cost function that enforces fair matchings is learned, it can be  
394 reuse on any set of points. This also ensures lower variance of fairness metrics (see experiments).

395 **Interpretability and flexibility.** Our second approach learns a parametrized cost function which  
 396 can be interpreted and used for subsequent downstream tasks. It also offers greater flexibility, since  
 397 the baseline cost function and the discrepancy  $\mathcal{D}$  can be set by the user.

## 398 D Proofs

### 399 D.1 Proof of Proposition 2.4

400 **Lemma D.1.** Let  $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{S})$  and  $\nu \in \mathcal{P}(\mathcal{Y} \times \mathcal{W})$ . Let  $p \in \mathcal{P}(\mathcal{S})$  and  $q \in \mathcal{P}(\mathcal{W})$  be obtained  
 401 from  $\mu$  and  $\nu$  by marginalizing, respectively,  $x$  and  $y$ , that is,

$$p(S = s) = \mu(\mathcal{X} \times \{s\}) \quad (13)$$

$$q(W = w) = \nu(\mathcal{Y} \times \{w\}). \quad (14)$$

402 Finally let  $F \in \Pi(p, q)$ .

403 There exists  $\pi \in \Pi(\mu, \nu)$  such that  $\pi$  is  $F$ -fair, that is

$$\pi(\mathcal{X} \times \{s\} \times \mathcal{Y} \times \{w\}) = F(S = s, W = w).$$

404 *Proof.* Given measurable sets  $A \subseteq \mathcal{X}$  and  $B \subseteq \mathcal{Y}$  let  $\pi$  be defined as

$$\pi(A \times \{s\} \times B \times \{w\}) := \frac{F(S = s, W = w)\mu(A \times \{s\})\nu(B \times \{w\})}{p(S = s)q(W = w)}$$

405 and let's check that **i)**  $\pi \in \Pi(\mu, \nu)$  and **ii)**  $\pi$  is  $F$ -fair:

406 **i)** We want to check that  $\pi(A \times \{s\} \times \mathcal{Y} \times \mathcal{W}) = \mu(A \times \{s\})$  and similarly for the other  
 407 marginal. To this end note that

$$\begin{aligned} \pi(A \times \{s\} \times \mathcal{Y} \times \mathcal{W}) &= \sum_{w \in \mathcal{W}} \frac{F(S = s, W = w)\mu(A \times \{s\})\nu(\mathcal{Y} \times \{w\})}{p(S = s)q(W = w)} \\ &= \mu(A \times \{s\}) \sum_{w \in \mathcal{W}} \frac{F(S = s, W = w)}{p(S = s)}, \end{aligned}$$

408 where the second equality follows from (14). To finish note that from  $F \in \Pi(p, q)$  follows  
 409 that  $\sum_{w \in \mathcal{W}} F(S = s, W = w) = p(S = s)$ . A similar argument shows that  $\pi(\mathcal{X} \times \mathcal{S} \times$   
 410  $B \times \{w\}) = \nu(B \times \{w\})$ .

411 **ii)** To see that  $\pi$  is  $F$ -fair note that it is immediate from (13) and (14) that

$$\begin{aligned} \pi(\mathcal{X} \times \{s\} \times \mathcal{Y} \times \{w\}) &= \frac{F(S = s, W = w)\mu(\mathcal{X} \times \{s\})\nu(\mathcal{Y} \times \{w\})}{p(S = s)q(W = w)} \\ &= F(S = s, W = w). \end{aligned}$$

412 □

413 **Proposition D.1.** Assume that  $\mu$  and  $\eta$  have bounded support. Let  $F$  be a coupling of the marginals  
 414 of  $\mu$  and  $\nu$  denoted by  $p$  and  $q$  defined via (13)-(14). Then there exists a unique  $F$ -fair transport plan.

415 *Proof.* Assume  $\mathcal{X}$  and  $\mathcal{Y}$  to be compact. Similar to the proof of Theorem 1.4. in [66] we can prove  
 416 that

$$\Lambda := \{\pi \in \mathcal{P}(\mathcal{X} \times \{0, 1\} \times \mathcal{Y} \times \{0, 1\}) : \pi_{SW} = F\}$$

417 is compact with respect to the weak topology: Let  $\pi_n$  be a sequence in  $\Lambda$ . They are probability  
 418 measures, so that their mass is 1, and hence they are bounded in the dual of  $C(\mathcal{X} \times \{0, 1\} \times \mathcal{Y} \times \{0, 1\})$ .  
 419 Hence usual weak- $\ast$  compactness in dual spaces guarantees the existence of a subsequence  $\pi_n \rightarrow \pi$   
 420 converging to a probability  $\pi$ . We just need to check  $\pi \in \Lambda$ . This may be done by fixing  $\phi \in$   
 421  $C(\{0, 1\} \times \{0, 1\})$  and from  $\pi_n \in \Lambda$  it follows that

$$\int \phi d\pi_n = \int \phi d[(\pi_n)_{SW}] = \int \phi dF = \sum_{s,w} \phi(s, w)F(S = s, W = w).$$

Now pass to the limit to obtain

$$\int \phi d\pi_{SW} = \int \phi d\pi = \sum_{s,w} \phi(s,w) F(S=s, W=W)$$

422 which shows that  $\pi \in \Lambda$ .

423 To finish, just note that  $\Pi_{\text{fair}}^{\mathbf{F}} = \Pi(\mu, \nu) \cap \Lambda$  is the intersection of two compact sets (with respect to  
424 weak topology) and Lemma D.1 establishes that it is non-empty. Continuity of the map defining  
425 the transport problem is enough to conclude the existence of a minimizer. Uniqueness is then a  
426 consequence of the strict convexity of KL.  $\square$

## 427 D.2 Proof of Theorem 3.1

428 **Notations.** Let  $\mathbf{F} \in \mathbb{R}^{K_s \times K_w}$  a given a fairness target. In what follows, we focus on  $K_s = K_w = 2$ ,  
429 and adopt the following notational conventions.

$$\begin{aligned} \xi_{11}(u, t) &:= ut - \mathbf{F}_{11} & \xi_{10}(u, t) &:= u(1-t) - \mathbf{F}_{10}, \\ \xi_{01}(u, t) &:= (1-u)t - \mathbf{F}_{01}, & \xi_{00}(u, t) &:= (1-u)(1-t) - \mathbf{F}_{00}. \end{aligned}$$

430 Given a measure  $\pi \in \mathcal{M}(\mathcal{X} \times \{0, 1\} \times \mathcal{Y} \times \{0, 1\})$ , define

$$\mathcal{L}_{\mathbf{F}}(\pi) := \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi \rangle^2 = \sum_{s,w \in \{0,1\}^2} \left( \int \xi_{sw}(z, w) d\pi(x, z, y, w) \right)^2.$$

431 For any distribution  $\alpha, \beta$  over a set  $\mathcal{X}$  such that  $\alpha$  is absolutely continuous with respect to  $\beta$  (i.e.,  
432  $\alpha \ll \beta$ ), we denote by  $\frac{d\alpha}{d\beta}$  the Radon-Nikodym derivative of  $\alpha$  with respect to  $\beta$ , and

$$\mathbf{KL}(\alpha || \beta) := \int_{\mathcal{X}} \log \left( \frac{d\alpha}{d\beta}(x) \right) d\alpha(x).$$

433 Given any two measures  $\alpha \in \mathcal{M}(\mathcal{X} \times \{0, 1\})$  and  $\beta \in \mathcal{M}(\mathcal{Y} \times \{0, 1\})$  define

$$m^*(\alpha, \beta) := \min_{\pi \in \Pi(\alpha, \beta)} \langle c, \pi \rangle + \varepsilon \mathbf{KL}(\pi || \alpha \otimes \beta) + \lambda \mathcal{L}_{\mathbf{F}}(\pi). \quad (15)$$

434 The goal is to prove the following result

435 **Theorem D.2.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be compact subsets of  $\mathbb{R}^d$  and let  $\mu$  and  $\eta$  be probability measures on  
436  $\mathcal{X} \times \{0, 1\}$  and  $\mathcal{Y} \times \{0, 1\}$ , respectively. Let  $(x_i, s_i)_{i=1}^n$  be  $n$  independent and identically distributed  
437 (i.i.d.) samples from  $\mu$  and let  $(y_j, w_j)_{j=1}^n$  be  $n$  i.i.d. samples from  $\eta$ ; assume further that the samples  
438 from  $\mu$  are independent from those from  $\eta$ . Finally, suppose that the cost  $c$  is  $L$ -Lipschitz and  $\mathcal{C}^\infty$ .  
439 Then

$$\mathbb{E}_{\mu \otimes \eta} |m^*(\mu_n, \eta_n) - m^*(\mu, \eta)| \leq \mathcal{O}(\log(n)/\sqrt{n^{-1}}), \quad (16)$$

440 *Proof.* Let  $m_\infty^* := m^*(\mu, \eta)$  and  $m_n^* := m^*(\mu_n, \eta_n)$ , and  $\pi_\infty^*$  be the minimizer attaining the  $m_\infty^*$   
441 that is a measure on  $\mathcal{X} \times \{0, 1\} \times \mathcal{Y} \times \{0, 1\}$  with marginals  $\mu$  and  $\eta$  such that

$$m_\infty^* = \langle c, \pi_\infty^* \rangle + \varepsilon \mathbf{KL}(\pi_\infty^* || \mu \otimes \eta) + \lambda \mathcal{L}_{\mathbf{F}}(\pi_\infty^*).$$

442 Let  $p_\infty^*$  denote the Radon-Nikodym density of  $\pi_\infty^*$  with respect to  $\mu \otimes \eta$  (note that the Kull-  
443 back-Leibler term forces  $\pi_\infty^*$  to be absolutely continuous with respect to  $\mu \otimes \eta$ ).

444 The idea is to “sandwich” the random variable  $m_n^* - m_\infty^*$  between two random variables with  
445 expectation upper bounded by  $\mathcal{O}(\log(n)/\sqrt{n})$ , that is, to find random variables  $Y_n, Z_n$  such that

$$Y_n - m_\infty^* \leq m_n^* - m_\infty^* \leq Z_n - m_\infty^* \quad (17)$$

446 and

$$\mathbb{E}[|Y_n - m_\infty^*|] \leq \mathcal{O}(\log(n)/\sqrt{n^{-1}}) \quad \text{and} \quad \mathbb{E}[|Z_n - m_\infty^*|] \leq \mathcal{O}(\log(n)/\sqrt{n^{-1}}).$$

447 To see that this is enough to establish (16), observe that (17) implies that

$$|m_n^* - m_\infty^*| \leq \max(|Y_n - m_\infty^*|, |Z_n - m_\infty^*|) \leq |Y_n - m_\infty^*| + |Z_n - m_\infty^*|.$$

448 **Lower bound.** To define  $Y_n$  start by observing that by linearizing the fairness penalization, we can  
 449 incorporate this linearization into the transport cost, so that the problem is now a proper entropic  
 450 optimal transport one. This allows us to leverage existing results on the sample complexity of optimal  
 451 transport. Moreover, the convexity of  $\mathcal{L}_{\mathbf{F}}$  implies that this linearization does yield a lower bound. To  
 452 this end note that for any measure  $\pi \in \mathcal{M}(\mathcal{X} \times \{0, 1\} \times \mathcal{Y} \times \{0, 1\})$  and any  $s, w \in \{0, 1\}^2$ , the  
 453 inequality  $0 \leq (\langle \xi_{sw}, \pi - \pi_{\infty}^* \rangle)^2 = (\langle \xi_{sw}, \pi \rangle)^2 + (\langle \xi_{sw}, \pi_{\infty}^* \rangle)^2 - 2\langle \xi_{sw}, \pi \rangle \langle \xi_{sw}, \pi_{\infty}^* \rangle$ , implies that

$$\mathcal{L}_{\mathbf{F}}(\pi) \geq -\mathcal{L}_{\mathbf{F}}(\pi_{\infty}^*) + 2 \sum_{(s,w) \in \{0,1\}^2} \langle \langle \xi_{sw}, \pi_{\infty}^* \rangle \xi_{sw}, \pi \rangle \quad (18)$$

$$= \mathcal{L}_{\mathbf{F}}(\pi_{\infty}^*) + \langle 2 \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_{\infty}^* \rangle \xi_{sw}, \pi - \pi_{\infty}^* \rangle. \quad (19)$$

454 This leads to a lower bound on  $m_n^*$  given by

$$m_n^* \geq -\lambda \mathcal{L}_{\mathbf{F}}(\pi_{\infty}^*) + \underbrace{\min_{\substack{\pi_1 = \mu_n \\ \pi_2 = \eta_n}} \langle c + 2\lambda \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_{\infty}^* \rangle \xi_{sw}, \pi \rangle + \varepsilon \mathbf{KL}(\pi \| \mu_n \otimes \eta_n)}_{:= \hat{m}_n^*} := Y_n \quad (20)$$

455 where  $\hat{m}_n^*$  is an entropic optimal transport problem with cost  $\hat{c}$ . From the sample complexity of  
 456 optimal transport (see Theorem 18 in [33]) it follows that

$$\mathbb{E}_{\mu \otimes \eta} |\hat{m}_n^* - \hat{m}_{\infty}^*| \leq \mathcal{O}(\sqrt{n}^{-1}), \quad (21)$$

457 with  $\hat{m}_{\infty}^*$  being the population version of  $\hat{m}_n^*$ , i.e.,

$$\hat{m}_{\infty}^* := \min_{\substack{\pi_1 = \mu \\ \pi_2 = \eta}} \langle c + 2\lambda \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_{\infty}^* \rangle \xi_{sw}, \pi \rangle + \varepsilon \mathbf{KL}(\pi \| \mu \otimes \eta).$$

458 To finish the lower bound, note that  $\mathbb{E}[|Y_n - m_{\infty}^*|] \leq \mathcal{O}(\sqrt{n}^{-1})$  follows from (21) provided we show  
 459 that  $m_{\infty}^* = \hat{m}_{\infty}^* - \lambda \mathcal{L}_{\mathbf{F}}(\pi_{\infty}^*)$ ; this is the content of Lemma D.4 that can be found in Appendix D.2.2.

460 **Upper bound.** Define the random positive measure

$$\hat{\pi}_n^* = \frac{1}{n^2} \sum_{ij} p_{\infty}^*((x_i, s_i), (y_j, w_j)) \delta_{((x_i, s_i), (y_j, w_j))} \quad (22)$$

461 and define  $\hat{Z}_n$  by

$$\hat{Z}_n := \langle c, \hat{\pi}_n^* \rangle + \varepsilon \mathbf{KL}(\hat{\pi}_n^* \| \mu_n \otimes \eta_n) + \lambda \mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*).$$

462 Observe that  $\hat{\pi}_n^*$  does not necessarily satisfies the marginal constraints that define  $\Pi(\mu, \eta)$ , and,  
 463 consequently,  $\hat{Z}_n$  is not necessarily an upper bound of  $m_n^*$ . A natural idea is to project  $\hat{\pi}_n^*$  onto the  
 464 constraint set. In what follows we interchangeably use the notation  $\hat{\pi}_n^*$  to denote (22) and the matrix  
 465 with  $ij$  entry given by  $p_{\infty}^*((x_i, s_i), (y_j, w_j))/n^2$ . An interesting projection to consider is the one  
 466 proposed in [1, Algorithm 2] named *Round*, reviewed in Appendix D.2.2. Let

$$\bar{\pi}_n^* = \text{Round}(\hat{\pi}_n^*),$$

467 where, similar to  $\hat{\pi}_n^*$ , the notation  $\bar{\pi}_n^*$  denotes both the matrix and the measure  
 468  $\sum_{ij} (\bar{\pi}_n^*)_{ij} \delta_{(x_i, s_i), (y_j, w_j)}$ . Define

$$Z_n := \int c d\bar{\pi}_n^* + \varepsilon \mathbf{KL}(\bar{\pi}_n^* \| \mu_n \otimes \eta_n) + \lambda \mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*). \quad (23)$$

469 By construction,  $\bar{\pi}_n^*$  satisfies the marginal constraints, hence it holds that  $Z_n \geq m_n^*$ . The remaining  
 470 goal is thus to prove  $\mathbb{E}[|Z_n - m_n^*|] \leq \mathcal{O}(\log(n)/\sqrt{n})$ . To do so the idea is to compare  $Z_n$  with  $\hat{Z}_n$   
 471 and  $\hat{Z}_n$  with  $m_{\infty}^*$ . In fact, from

$$|Z_n - m_{\infty}^*| \leq |Z_n - \hat{Z}_n| + |\hat{Z}_n - m_{\infty}^*|,$$

472 it is sufficient to prove (i)  $\mathbb{E}[|Z_n - \hat{Z}_n|] \in \mathcal{O}(\log(n)/\sqrt{n}^{-1})$  and (ii)  $\mathbb{E}[|\hat{Z}_n - m_{\infty}^*|] \in \mathcal{O}(\log(n)/\sqrt{n}^{-1})$ .

473 (i).  $\mathbb{E}|\widehat{Z}_n - m_\infty^*|$ . Let  $V_{i,j} := ((x_i, s_i), (y_j, w_j))$  and since the samples are i.i.d.,

$$\mathbb{E}\langle c, \widehat{\pi}_n^* \rangle = \frac{1}{n^2} \sum_{ij} \mathbb{E}_{\mu \otimes \eta} [p_\infty^*(V_{i,j})c(V_{i,j})] = \langle cp_\infty^*, \mu \otimes \eta \rangle = \langle c, \pi_\infty^* \rangle.$$

474 Moreover, let  $d\widehat{\pi}_n^*/d(\mu_n \otimes \eta_n)$  be the Radon-Nikodym derivative of  $\widehat{\pi}_n^*$  with respect to  $\mu_n \otimes \eta_n$  and,  
475 by definition,

$$\begin{aligned} \mathbb{E}[\mathbf{KL}(\widehat{\pi}_n^* || \mu_n \otimes \eta_n)] &= \mathbb{E}\langle \log \frac{d\widehat{\pi}_n^*}{d(\mu_n \otimes \eta_n)}, \widehat{\pi}_n^* \rangle = \frac{1}{n^2} \sum_{ij} \mathbb{E}_{\mu \otimes \eta} \log p_\infty^*(V_{i,j}) p_\infty^*(V_{i,j}) \\ &= \langle \log(p_\infty^*) p_\infty^*, \mu \otimes \eta \rangle = \langle \log \frac{d\pi_\infty^*}{d(\mu \otimes \eta)}, \pi_\infty^* \rangle = \mathbf{KL}(\pi_\infty^* || \mu \otimes \eta). \end{aligned}$$

476 Since  $\mathcal{X}$  and  $\mathcal{Y}$  are compact and  $c$  is continuous it follows that  $c(V_{i,j})$  is uniformly bounded. Further-  
477 more, by Lemma D.6, there exists  $r > 0$  such that,  $\mu \otimes \eta$ -a.s.,  $\frac{1}{r} \leq p_\infty^*(V_{i,j}) \leq r$ . Consequently, it  
478 follows by Lemma D.5,

$$\max\{\mathbb{E}|\langle c, \widehat{\pi}_n^* - \pi_\infty^* \rangle|, \mathbb{E}|\mathbf{KL}(\widehat{\pi}_n^* || \mu_n \otimes \eta_n) - \mathbf{KL}(\pi_\infty^* || \mu \otimes \eta)|\} \in \mathcal{O}(\sqrt{n}^{-1}). \quad (24)$$

479 Finally, for the last term, note that  $x^2 - y^2 = (x - y)^2 + 2y(x - y)$ , hence for  $s, w \in \{0, 1\}$

$$\mathbb{E}|\langle \xi_{sw}, \widehat{\pi}_n^* \rangle^2 - \langle \xi_{sw}, \pi_\infty^* \rangle^2| \leq \mathbb{E}\langle \xi_{sw}, \widehat{\pi}_n^* - \pi_\infty^* \rangle^2 + 2\mathbb{E}[|\langle \xi_{sw}, \pi_\infty^* \rangle \langle \xi_{sw}, \widehat{\pi}_n^* - \pi_\infty^* \rangle|] \in \mathcal{O}(\sqrt{n}^{-1}), \quad (25)$$

480 where we conclude similarly as above using Lemma D.5.

481 Inequalities (24) and (25) are enough to establish  $\mathbb{E}|\widehat{Z}_n - m_\infty^*| \in \mathcal{O}(\sqrt{n}^{-1})$  and conclude the proof  
482 of (i).

483 (ii).  $\mathbb{E}|Z_n - \widehat{Z}_n|$ . To obtain a bound on  $\mathbb{E}|Z_n - \widehat{Z}_n|$  we first translate it into a bound on  $\mathbb{E}\|\widehat{\pi}_n^* - \pi_n^*\|_1$ .  
484 In fact, Lemma D.7 shows that

$$\mathbb{E}|Z_n - \widehat{Z}_n| \leq \mathcal{O}(\log(n))\mathbb{E}\|\widehat{\pi}_n^* - \pi_n^*\|_1$$

485 Moreover [1, Lemma 7] implies that

$$\|\widehat{\pi}_n^* - \pi_n^*\|_1 \leq 2\|\widehat{\pi}_n^* \mathbf{1} - 1/n\mathbf{1}\|_1 + 2\|(\widehat{\pi}_n^*)^T \mathbf{1} - 1/n\mathbf{1}\|_1, \quad (26)$$

486 and, consequently, to obtain the bound on  $\mathbb{E}|Z_n - \widehat{Z}_n|$  it is enough to show that, in expectation, the  
487 right-hand-side of (26) is  $\mathcal{O}(\sqrt{n}^{-1})$ . This last bound essentially reduces to the proof of Lemma 14 of  
488 [65] together with the Cauchy Schwartz inequality – this is made precise in Lemma D.8.  $\square$

## 489 D.2.1 Review of Round

490 This section reviews the Round map of [1] defined as: Given  $M \in \mathbb{R}_{\geq 0}^{n^2}$ , the matrix  $G := \text{Round}(M)$   
491 is obtained via

- 492 1. Let  $M' := \text{diag}(x)M$ , with  $x_i := \min(1, \frac{1/n}{r_i(M)})$  and  $r_i(M) := (M\mathbf{1})_i$
- 493 2. Let  $M'' := M' \text{diag}(y)$ , with  $y_j := \min(1, \frac{1/n}{c_j(M')})$  and  $c_j(M') := ((M')^T \mathbf{1})_j$
- 494 3. Output  $G := M'' + \frac{1}{\|\text{err}_r\|_1} \text{err}_r \text{err}_c^T$ , where

$$\text{err}_r := 1/n\mathbf{1} - r(M'') \quad \text{and} \quad \text{err}_c := 1/n\mathbf{1} - c(M'').$$

495 **Lemma D.3.** Suppose that  $\exp(-K)/n^2 \leq M_{ij} \leq \exp(K)/n^2$  and let  $G := \text{Round}(M)$ . Then  
496  $G_{ij} \geq \exp(-4K)/n^2$ .

497 *Proof.* First observe that

498 **Claim 1.**  $\exp(-2K)/n^2 \leq M'_{ij} \leq M_{ij} \leq \exp(K)/n^2$ .

499 To see this let  $M_i$  denote the  $i$ th line of  $M$  and similarly for  $M'$ . We have that

$$M'_i = x_i M_i \leq M_i, \quad \text{since } x_i \leq 1.$$

500 which shows that  $M'_{ij} \leq M_{ij} \leq \exp(K)/n^2$ .

501 To obtain the lower bound suppose that  $x_i < 1$  (otherwise there is nothing to prove). In this case we  
502 have  $M'_i = \frac{1/n}{r_i(M)} M_i$  and, from the upper bounds on  $M$ ,

$$r_i(M) = \sum_j M_{ij} \leq n(\exp(K)/n^2).$$

503 As a consequence,

$$M'_{ij} \geq \frac{1/n}{n(\exp(K)/n^2)} M_{ij} = \exp(-K) M_{ij} \geq \exp(-2K)/n^2,$$

504 thus establishing claim 1.

505 **Claim 2.**  $\exp(-4K)/n^2 \leq M''_{ij} \leq M'_{ij} \leq \exp(K)/n^2$ .

506 To prove Claim 2 we argue as in Claim 1.

507 **Claim 3.** The vectors  $\text{err}_r$  and  $\text{err}_c$  are both non-negative.

508 To see this just observe that after step 1, row  $i$  of  $M'$  satisfies  $r_i(M') \leq 1/n$  and, since  $\text{diag}(y) \leq$   
509  $\text{diag}(1)$ , step 2 can only decrease  $r_i$ . A similar reasoning shows that  $c_j(M'') \leq 1/n$ .  $\square$

## 510 D.2.2 Technical Lemmas

511 **Lemma D.4.** Let  $\pi_\infty^*$  be a minimizer of

$$\min_{\substack{\pi_1 = \mu \\ \pi_2 = \eta}} F(\pi) := \langle c, \pi \rangle + \varepsilon \mathbf{KL}(\pi \| \mu \otimes \eta) + \lambda \mathcal{L}_{\mathbf{F}}(\pi) := m_\infty^*. \quad (27)$$

512 Then  $\pi_\infty^*$  also minimizes

$$\begin{aligned} \min_{\substack{\pi_1 = \mu \\ \pi_2 = \eta}} F_L(\pi) &:= \langle c, \pi \rangle + \varepsilon \mathbf{KL}(\pi \| \mu \otimes \eta) + \lambda \mathcal{L}_{\mathbf{F}}(\pi_\infty^*) + \langle 2\lambda \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_\infty^* \rangle \xi_{sw}, \pi - \pi_\infty^* \rangle \\ &= \min_{\substack{\pi_1 = \mu \\ \pi_2 = \eta}} \langle c + 2\lambda \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_\infty^* \rangle \xi_{sw}, \pi \rangle + \varepsilon \mathbf{KL}(\pi \| \mu \otimes \eta) - \lambda \mathcal{L}_{\mathbf{F}}(\pi_\infty^*) \\ &= \widehat{m}_\infty^* - \lambda \mathcal{L}_{\mathbf{F}}(\pi_\infty^*). \end{aligned} \quad (28)$$

513 and, consequently,  $\widehat{m}_\infty^* - \lambda \mathcal{L}_{\mathbf{F}}(\pi_\infty^*) = m_\infty^*$ .

514 *Proof.* The proof is inspired by the one from [64, Proposition 2.1]. Assume  $\pi_\infty^*$  minimizes Eq. (27).

515 The presence of entropic regularization implies that  $\pi_\infty^* \ll \mu \otimes \eta$ . Let

$$g_L = 2 \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \pi_\infty^* \rangle \xi_{sw}.$$

516 We already know that for any  $\pi \in \Pi(\mu, \eta)$

$$\mathcal{L}_{\mathbf{F}}(\pi) \geq \mathcal{L}_{\mathbf{F}}(\pi_\infty^*) + \langle g_L, \pi - \pi_\infty^* \rangle.$$

517 Define

$$g_{\mathbf{KL}} = \log \left( \frac{d\pi_\infty^*}{d(\mu \otimes \eta)} \right),$$

518 we prove, for any  $\pi \in \Pi(\mu, \eta)$ ,

$$\mathbf{KL}(\pi \| \mu \otimes \eta) \geq \mathbf{KL}(\pi_\infty^* \| \mu \otimes \eta) + \langle g_{\mathbf{KL}}, \pi - \pi_\infty^* \rangle,$$

519 with convex analysis<sup>1</sup>, and thus

$$F(\pi) \geq F(\pi_\infty^*) + \langle c + \lambda g_L + \varepsilon g_{\text{KL}}, \pi - \pi_\infty^* \rangle.$$

520 Consequently, for any  $\pi \in \Pi(\mu, \eta)$  such that  $\pi \ll \mu \otimes \eta^2$ ,

$$\langle c + \lambda g_L + \varepsilon g_{\text{KL}}, \pi - \pi_\infty^* \rangle \geq 0,$$

521 otherwise, for  $r \in [0, 1]$ ,  $\pi_r = (1-r)\pi_\infty^* + r\pi \in \Pi(\mu, \eta) \cap \{\nu; \nu \ll \mu \otimes \eta\}$ , we could find  $r^* \in (0, 1]$ ,  
 522 such that  $F(\pi_{r^*}) < F(\pi_\infty^*)$ . In fact, by contradiction, assume  $\langle c + \lambda g_L + \varepsilon g_{\text{KL}}, \pi - \pi_\infty^* \rangle < 0$ , as  
 523  $\lim_{r \rightarrow 0^+} \frac{F(\pi_r) - F(\pi_\infty^*)}{r} = \langle c + \lambda g_L + \varepsilon g_{\text{KL}}, \pi - \pi_\infty^* \rangle$ , we have, by continuity, the existence of a  
 524 sufficiently small  $r^* \in (0, 1]$  such that  $F(\pi_{r^*}) < F(\pi_\infty^*)$ . Moreover, as  $F_L(\pi_\infty^*) = F(\pi_\infty^*)$ , for any  
 525  $\pi \in \Pi(\mu, \eta)$

$$F_L(\pi) \geq F_L(\pi_\infty^*) + \langle c + \lambda g_L + \varepsilon g_{\text{KL}}, \pi - \pi_\infty^* \rangle \geq F_L(\pi_\infty^*).$$

526 Hence,  $\pi_\infty^*$  is also a minimizer of Eq. (28).  $\square$

527 **Lemma D.5.** Suppose  $a \in L^\infty(\mu \otimes \eta)$  is such that  $\langle a, \mu \otimes \eta \rangle = 0$ . Then,

$$\mathbb{E} |\langle a, \mu_n \otimes \eta_n \rangle| \in \mathcal{O}(\sqrt{n}^{-1}) \quad \text{and} \quad \mathbb{E} \langle a, \mu_n \otimes \eta_n \rangle^2 \in \mathcal{O}(n^{-1}),$$

528 *Proof.* Let  $\mathfrak{S}_n$  denote the set of all permutation on  $n$  elements. By counting terms, we observe that

$$\langle a, \mu_n \otimes \eta_n \rangle = \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} \frac{1}{n} \sum_{k=1}^n a((X_k, S_k), (Y_{\sigma(k)}, W_{\sigma(k)})).$$

529 In fact, there are  $(n-1)!$  permutations  $\sigma \in \mathfrak{S}_n$  satisfying  $\sigma(i) = j$  because there are  $(n-1)!$   
 530 ways to permute the rest. Hence, the contribution of the specific  $a((X_k, S_k), (Y_{\sigma(k)}, W_{\sigma(k)}))$  is  
 531 weighted by  $\frac{1}{n!} \frac{(n-1)!}{n} = \frac{1}{n^2}$ . Now, we observe that for any fixed permutation  $\sigma$ , the joint law  
 532  $((X_1, S_1), (Y_{\sigma(1)}, W_{\sigma(1)}), \dots, (X_n, S_n), (Y_{\sigma(n)}, W_{\sigma(n)}))$  is identical to that of  $V_1, \dots, V_n$  where  
 533  $V_k \sim \mu \otimes \eta$  are independent and identically distributed. Thus,

$$\begin{aligned} \mathbb{E} [\langle a, \mu_n \otimes \eta_n \rangle] &\leq \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} \frac{1}{n} \mathbb{E} \left| \sum_{k=1}^n a(V_k) \right| = \frac{1}{n} \mathbb{E} \left| \sum_{k=1}^n a(V_k) \right| \leq \frac{1}{n} \sqrt{\mathbb{E} \left[ \left( \sum_{k=1}^n a(V_k) \right)^2 \right]} \\ &= \frac{1}{n} \sqrt{\mathbb{E} \left[ \sum_{k=1}^n a(V_k)^2 \right]} \leq \frac{\|a\|_{L^\infty(\mu \otimes \eta)}}{\sqrt{n}}. \end{aligned}$$

534 Moreover, by applying Jensen inequality,

$$\begin{aligned} \mathbb{E} \langle a, \mu_n \otimes \eta_n \rangle^2 &\leq \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} \mathbb{E} \left( \frac{1}{n} \sum_{k=1}^n a(V_k) \right)^2 = \mathbb{E} \left( \frac{1}{n} \sum_{k=1}^n a(V_k) \right)^2 \leq \frac{1}{n^2} \mathbb{E} \left( \sum_{k=1}^n a(V_k) \right)^2 \\ &\leq \frac{1}{n^2} \mathbb{E} \sum_{k=1}^n a(V_k)^2 \leq \frac{\|a\|_{L^\infty(\mu \otimes \eta)}^2}{n}. \quad \square \end{aligned}$$

535 **Lemma D.6.** If  $\mathcal{X}$  and  $\mathcal{Y}$  are compact, then there exists  $r, s > 0^3$  such that,  $\mu \otimes \eta$ -a.s.,

$$\frac{1}{r} \leq p_\infty^*(V_{i,j}) \leq r, \quad \text{and} \quad \frac{1}{s} \leq \frac{d\bar{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \leq n^2.$$

<sup>1</sup>If  $\pi$  is singular with respect to  $\mu \otimes \eta$ , the inequality is trivial. Assume  $\pi \ll \mu \otimes \eta$ , we define for  $r \in [0, 1]$ ,  $\pi_r = (1-r)\pi_\infty^* + r\pi \in \Pi(\mu, \eta)$  and  $\phi(r) = \mathbf{KL}(\pi_r \| \mu \otimes \eta)$ . By convexity of  $\phi$ ,  $\phi(r) \geq \phi(0) + r \lim_{h \rightarrow 0^+} \frac{\phi(h) - \phi(0)}{h}$ . We compute the previous directional limit using the monotone convergence theorem as in Csiszár [19, Lemma 2.1].

<sup>2</sup>If instead  $\pi$  is singular with respect to  $\mu \otimes \eta$ , then  $F(\pi) = F_L(\pi) = \infty$ , so such measures are trivially excluded from being minimizers of either functional.

<sup>3</sup>only depending on the size of the ambient space and the entropic regularization magnitude.

536 *Proof.* u Lemma D.4 shows that  $\pi_\infty^*$  is a solution of an optimal transport problem with a continuous  
537 cost  $\hat{c}$ . Moreover,  $p_\infty^* = \exp(\varphi_\infty^* \oplus \psi_\infty^* - \hat{c})$  (see Theorem 4.2. [57]) and, from Lemma 4.11 in  
538 [57], the functions  $\varphi_\infty^*$  and  $\psi_\infty^*$  can be assumed to be continuous (since  $\hat{c}$  is). The result for  $p_\infty^*$  now  
539 follows from the compactness of  $\mathcal{X}$  and  $\mathcal{Y}$ .

540 The result for  $\bar{\pi}_n^*$  follows from Lemma D.3 and the fact that  $\bar{\pi}_n^*$  is a probability, i.e.,  $(\bar{\pi}_n^*)_{ij} \leq 1$ .  $\square$

541 **Lemma D.7.** If  $\mathcal{X}$  and  $\mathcal{Y}$  are compact, then

$$|Z_n - \hat{Z}_n| \leq \mathcal{O}(\log(n)) \|\hat{\pi}_n^* - \bar{\pi}_n^*\|_1.$$

542 *Proof.* Denote, for any  $i, j \in [n]^2$ ,  $V_{i,j} := ((x_i, s_i), (y_j, w_j))$

$$\begin{aligned} |Z_n - \hat{Z}_n| &= |c \bar{\pi}_n^* - \hat{\pi}_n^* + \varepsilon (\mathbf{KL}(\bar{\pi}_n^* \|\mu_n \otimes \eta_n) - \mathbf{KL}(\hat{\pi}_n^* \|\mu_n \otimes \eta_n)) + \lambda (\mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) - \mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*))| \\ &\leq \|c\|_\infty \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1 + \varepsilon |\mathbf{KL}(\bar{\pi}_n^* \|\mu_n \otimes \eta_n) - \mathbf{KL}(\hat{\pi}_n^* \|\mu_n \otimes \eta_n)| + \lambda |\mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) - \mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*)|. \end{aligned}$$

543 We bound the two remaining terms as follows. First, by Lemma D.6, there exists  $r, s > 0$  such that,

544  $\mu \otimes \eta$ -a.s.,  $\frac{1}{r} \leq p_\infty^*(V_{i,j}) \leq r$  and  $\frac{1}{s} \leq \frac{d\bar{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \leq n^2$ .

545 (i)

$$\mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*) - \mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) \geq \langle 2 \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \bar{\pi}_n^* \rangle \xi_{sw}, \hat{\pi}_n^* - \bar{\pi}_n^* \rangle,$$

546 and

$$\mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) - \mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*) \geq \langle 2 \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \hat{\pi}_n^* \rangle \xi_{sw}, \bar{\pi}_n^* - \hat{\pi}_n^* \rangle,$$

547 hence, for  $n$  sufficiently large,

$$\begin{aligned} |\mathcal{L}_{\mathbf{F}}(\bar{\pi}_n^*) - \mathcal{L}_{\mathbf{F}}(\hat{\pi}_n^*)| &\leq 2 \max\left\{ \left\| \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \hat{\pi}_n^* \rangle \xi_{sw} \right\|_\infty, \left\| \sum_{s,w \in \{0,1\}^2} \langle \xi_{sw}, \bar{\pi}_n^* \rangle \xi_{sw} \right\|_\infty \right\} \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1 \\ &\leq 8 \max\{\|\hat{\pi}_n^*\|_\infty, \|\bar{\pi}_n^*\|_\infty\} \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1 \\ &\leq \mathcal{O}(1) \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1, \end{aligned}$$

548 where we used the fact that  $\xi_{sw} \in [-1, 1]$ ,  $\forall s, w \in \{0, 1\}$ .

549 (ii) Similarly as in Lemma D.4, we can prove

$$\mathbf{KL}(\hat{\pi}_n^* \|\mu_n \otimes \eta_n) - \mathbf{KL}(\bar{\pi}_n^* \|\mu_n \otimes \eta_n) \geq \langle \log \left( \frac{d\bar{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \right), \hat{\pi}_n^* - \bar{\pi}_n^* \rangle,$$

550 and

$$\mathbf{KL}(\bar{\pi}_n^* \|\mu_n \otimes \eta_n) - \mathbf{KL}(\hat{\pi}_n^* \|\mu_n \otimes \eta_n) \geq \langle \log \left( \frac{d\hat{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \right), \bar{\pi}_n^* - \hat{\pi}_n^* \rangle,$$

551 hence, for  $n$  sufficiently large,

$$\begin{aligned} &|\mathbf{KL}(\bar{\pi}_n^* \|\mu_n \otimes \eta_n) - \mathbf{KL}(\hat{\pi}_n^* \|\mu_n \otimes \eta_n)| \\ &\leq \max\left\{ \left\| \log \left( \frac{d\hat{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \right) \right\|_\infty, \left\| \log \left( \frac{d\bar{\pi}_n^*}{d(\mu_n \otimes \eta_n)} \right) \right\|_\infty \right\} \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1 \leq \mathcal{O}(\log(n)) \|\bar{\pi}_n^* - \hat{\pi}_n^*\|_1. \end{aligned}$$

552  $\square$

553 **Lemma D.8.** Let  $\hat{\pi}_n^*$  be an  $n \times n$  matrix with  $ij$  entry defined by  $(\hat{\pi}_n^*)_{ij} = p_\infty^*(x_i, y_j)/n^2$  where  
554  $p_\infty^*$  is a continuous function satisfying  $\langle p_\infty^*(x, \cdot), \eta \rangle = \langle p_\infty^*(\cdot, y), \mu \rangle = 1^4$ . Then

$$\mathbb{E} \|\hat{\pi}_n^* \mathbf{1} - 1/n \mathbf{1}\|_1 \leq \mathcal{O}(\sqrt{n}^{-1}) \quad \text{and} \quad \mathbb{E} \|(\hat{\pi}_n^*)^T \mathbf{1} - 1/n \mathbf{1}\|_1 \leq \mathcal{O}(\sqrt{n}^{-1}).$$

<sup>4</sup>This condition is a consequence of  $\hat{\pi}_\infty^* \in \Pi(\mu, \eta)$ .

555 *Proof.* We consider only the bound of  $\|\hat{\pi}_n^* \mathbf{1} - 1/n \mathbf{1}\|_1$  since the other is analogous. Note that

$$\begin{aligned} \|\hat{\pi}_n^* \mathbf{1} - 1/n \mathbf{1}\|_1 &= \sum_i |1/n - 1/n^2 \sum_j (\hat{\pi}_n^*)_{ij}| = 1/n \sum_i |1 - 1/n \sum_j p_\infty^*(x_i, y_j)| \\ &\leq \frac{\sqrt{n}}{n} \sqrt{\sum_i (1 - 1/n \sum_j p_\infty^*(x_i, y_j))^2} = \sqrt{\frac{1}{n} \sum_i (1 - 1/n \sum_j p_\infty^*(x_i, y_j))^2} \end{aligned}$$

556 where the inequality follows from Cauchy-Schwartz. To finish note that the proof of Lemma 14 in  
557 [65] only relies on the marginal constraints  $\langle p_\infty^*(x, \cdot), \eta \rangle = \langle p_\infty^*(\cdot, y), \mu \rangle = 1$  and the boundedness  
558 of  $p_\infty^*$ . Consequently,

$$\mathbb{E} \frac{1}{n} \sum_i (1 - 1/n \sum_j p_\infty^*(x_i, y_j))^2 \leq \mathcal{O}(n^{-1})$$

559 and the result then follows by Jensen's inequality.  $\square$

## 560 E Experimental details and results for cost learning

### 561 E.1 Experimental details

562 Data  $X$  and  $Y$  are generated from mixtures of two bivariate Gaussian components, the sensitive  
563 attribute serving as an indicator of the Gaussian component from which each sample originates. In  
564 all the experiments, the strength of the entropic regularization is fixed at  $\epsilon = 1$  and only fairness  
565 regularization varies. The base cost used is the squared euclidean distance.

566 For the penalized approach, the problem is solved by using the implementation of the Generalized Con-  
567 ditional Gradient algorithm provided by the Python package Python Optimal Transport (POT) [28, 29].

568 For the cost learning approach, the problem is solved using Adam algorithm with a learning rate fixed  
569 at  $5 \times 10^{-3}$ . The inner problem is solved using the POT implementation of the Sinkhorn algorithm.  
570 The gradient of the bilevel objective is computed by backpropagation in the Sinkhorn iterations.

### 571 E.2 Results for the cost learning approach

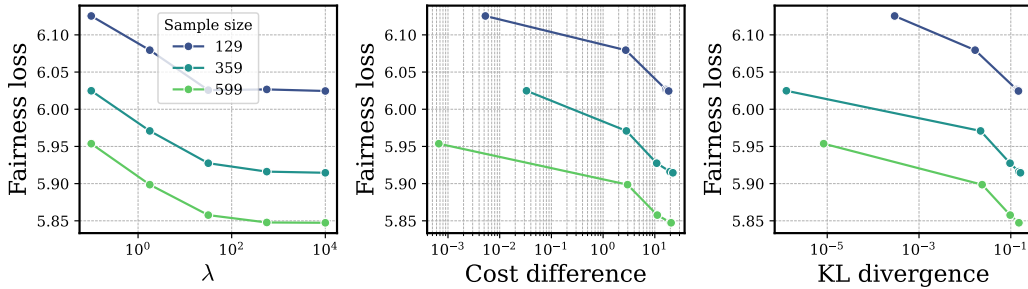


Figure 3: Fairness vs. penalty (**left**), and fairness-KL divergence (**center**) and fairness-cost difference (**right**) trade-offs (w.r.t. the non penalized problem) for varying sample sizes for the **cost-learning algorithm** with cost-sensitive fairness loss.