

---

# Statistical Guarantees for Fairness Aware Plug-In Algorithms

---

Anonymous Authors<sup>1</sup>

## Abstract

A plug-in algorithm to estimate Bayes Optimal Classifiers for fairness-aware binary classification has been proposed in (Menon & Williamson, 2018). However, the statistical efficacy of their approach has not been established. We prove that the plug-in algorithm is statistically consistent. We also derive finite sample guarantees associated with learning the Bayes Optimal Classifiers via the plug-in algorithm. Finally, we propose a protocol that modifies the plug-in approach, so as to simultaneously guarantee fairness and differential privacy with respect to a binary feature deemed sensitive.

## 1. Introduction and Related Work

Bayes Optimal Classifiers (BOCs) (Devroye et al., 1996) are of significant importance, since they achieve the least average error possible for any classification task. However, BOCs are generally specified in terms of unknown distributional quantities. Constructing sound estimators for BOCs, provided access to only a finite training sample, is thus of utmost practical relevance. One approach to estimating the BOC is through constructing 'plug-in' estimators. The plug-in principle applied to a broad class of problems, including that of binary classification, is well studied in the statistics literature (Audibert et al., 2007; Denis & Hebiri, 2017; Yang, 1999). Indeed, the existence of a plug-in classifier that is optimal in the minimax sense is established in (Audibert et al., 2007; Yang, 1999). In their work, (Menon & Williamson, 2018) propose a plug-in algorithm to estimate the BOCs corresponding to fairness-aware learning (FAL) tasks. However, (Menon & Williamson, 2018) do not provide guarantees on the statistical efficacy of their algorithm. In this paper, we plug this gap by *proving that the plug-in algorithm of (Menon & Williamson, 2018), is indeed statistically consistent*. We also *characterise the*

*sample complexity associated with the task of learning a low regret classifier via the plug-in algorithm*. Closest to our work is that of (Chzhen et al., 2019), wherein an asymptotic study (for a different fairness aware plug-in classifier) is carried out. The work of (Chzhen et al., 2019) however, focuses on settings wherein perfect fairness constraints are imposed. It is well established that due to inherent fairness-accuracy trade-offs, ensuring perfect fairness without considerable loss in accuracy is generally not possible (Menon & Williamson, 2018; Zhao & Gordon, 2019; Chen et al., 2018). We thus focus on approximate notions of two fairness metrics, Demographic Parity (DPar) and Equality of Opportunity (EO). Further, the approach of (Chzhen et al., 2019) requires access to the sensitive variable (denoted  $\bar{Y}$  hereon) at test time which is often not permitted. The plug-in approach of (Menon & Williamson, 2018) however does not necessitate test-time access to  $\bar{Y}$ . Indeed, real-world settings may impose even more stringent requirements on  $\bar{Y}$ . For example, we may be required to ensure that our model does not leak information about the sensitive attribute,  $\bar{Y}$ , corresponding to any individual. In such cases, a possible solution is to protect individuals via Differential Privacy (DP) (Dwork et al., 2006). The literature combining fairness and privacy (Jagielski et al., 2019; Cummings et al., 2019; Mozannar et al., 2020), is emerging and limited. Such settings motivate us to *propose an easy to deploy, modified version of the plug-in algorithm, referred to as DP Plug-in*. The framework ensures that  $\bar{Y}$  is protected via DP. Using publicly available data sets, we demonstrate empirically, that the *DP Plug-in algorithm achieves strong privacy-fairness-accuracy guarantees*, as it outperforms the private, fair approach of (Jagielski et al., 2019) across 3 out of 4 experimental setups considered.

## 2. Background and Notation

For brevity, we only introduce the main features from (Menon & Williamson, 2018) pertinent to our study in this section. We present other useful definitions and results from (Menon & Williamson, 2018) in section A of the supplement. Additionally, our focus in the main thesis of this paper will be on the approximate EO criterion for the case when  $\bar{Y}$  is unavailable during test-time. Analogous (and relatively simpler) analyses for 1) the case when  $\bar{Y}$  is available at test time, and for 2) the approximate DPar criterion are presented in sections B and C of the

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

supplement for completeness.

Access to a finite training sample,  $S = \{x_i, y_i, \bar{y}_i\}_{i=1}^n$  drawn i.i.d from some unknown distribution  $\mathbb{P}$  is assumed in (Menon & Williamson, 2018).  $\forall i \in [n]$  the triplet  $(x_i, y_i, \bar{y}_i)$  is a realisation of the random variable triplet  $(X, Y, \bar{Y})$  comprising of the feature, label and sensitive attribute respectively. Let  $\pi = \mathbb{P}(Y = 1)$ ,  $\bar{\pi} = \mathbb{P}(\bar{Y} = 1)$  and  $\beta = \mathbb{P}(\bar{Y} = 1|Y = 1)$ . Assume that  $\pi, \bar{\pi}, \beta > 0$

Let  $(X, Y) \sim \mathcal{D}$ ,  $(X, \bar{Y}|Y = 1) \sim \bar{\mathcal{D}}_{EO}$ .  $f: \mathcal{X} \rightarrow [0, 1]$  denotes a randomised classifier on measurable domain  $\mathcal{X}$ .  $f$  yields predictions via  $(\hat{Y}|X = x) \sim \text{Bernoulli}(f(x))$ . Regression functions w.r.t.  $\mathcal{D}, \bar{\mathcal{D}}_{EO}$  are given by  $\eta(x) = \mathbb{P}(Y = 1|X = x)$ ,  $\bar{\eta}_{EO}(x, y) = \mathbb{P}(\bar{Y} = 1|X = x, Y = y)$  respectively.

A central object of interest in (Menon & Williamson, 2018), is the notion of cost-sensitive risks (CSR). Denoting false positive and negative rates of  $f$  w.r.t  $\mathcal{D}$ , by  $FPR_{\mathcal{D}}(f)$  and  $FNR_{\mathcal{D}}(f)$  respectively, the CSR of a classifier  $f$  w.r.t a distribution  $\mathcal{D}$ , parameterised by  $c \in [0, 1]$  is given by:

$$CS(f; \mathcal{D}, c) := c(1 - \pi)FPR_{\mathcal{D}}(f) + \pi(1 - c)FNR_{\mathcal{D}}(f)$$

**Definition 2.2** A binary classifier  $f$ , with corresponding predictor  $\hat{Y}$  admits Equality of Opportunity if:  $\mathbb{P}(\hat{Y} = 1|Y = 1, \bar{Y} = 0) = \mathbb{P}(\hat{Y} = 1|Y = 1, \bar{Y} = 1)$

Thus, EO requires parity in the TPRs between groups as explicated in Definition 2.2. Obtaining perfect fairness while retaining non-trivial accuracy is generally not possible, and so (Menon & Williamson, 2018) introduce approximate measures of fairness which require the additive or multiplicative disparity between prediction rates to be small. A key lemma in (Menon & Williamson, 2018) draws an equivalence between the super-level sets of approximate fairness measures and CSRs. This in turn leads to a reduction of the FAL problem to a problem with constraints on cost-sensitive risks (The reader may refer to section A of the supplement for a more thorough presentation of these results). The FAL problem in (Menon & Williamson, 2018) is thus posed in terms of CSRs as follows:

**Problem 2.1 (Cost-sensitive FAL)** For trade-off parameter  $\lambda \in \mathbb{R}$ , and cost parameters  $c, \bar{c} \in (0, 1)^2$ , minimise the fairness-aware cost-sensitive risk:

$$R_{FA}(f; \mathcal{D}, \bar{\mathcal{D}}_{EO}, c, \bar{c}, \lambda) = CS(f; \mathcal{D}, c) - \lambda CS(f; \bar{\mathcal{D}}_{EO}, \bar{c})$$

Equipped with this soft constrained FAL problem formulated in terms of CSRs, (Menon & Williamson, 2018) derive the BOCs corresponding to such problems. We present the BOC for approx. EO and the plug-in algorithm of (Menon & Williamson, 2018) to estimate this BOC, in Theorem 2.2 and Algorithm 1 respectively. It is this (optimal) classifier and algorithm that will make for the key objects of our theoretical analysis in Section 3.

**Theorem 2.2 (BOC for FAL)** Pick any costs  $c, \bar{c} \in (0, 1)^2$  and trade-off parameter  $\lambda \in \mathbb{R}$ . Then:

$$\text{Argmin}_{R_{FA}}(f; \mathcal{D}, \bar{\mathcal{D}}_{EO}) = \{H_{\alpha} \circ s^*(x) | \alpha \in [0, 1]\}$$

$$\text{where, } s^*(x) = \left\{1 - \frac{\lambda}{\pi} (\bar{\eta}_{EO}(x, 1) - \bar{c})\right\} \eta(x) - c$$

$$\text{and, } H_{\alpha}(z) = \mathbb{I}(z > 0) + \alpha \mathbb{I}(z = 0)$$

---

### Algorithm 1 Plugin approach to FAL, EO setting

---

**Input:** Sample  $S = \{x_i, y_i, \bar{y}_i\}_{i=1}^n$  from distribution  $\mathbb{P}$ ; cost parameters  $c, \bar{c}$ ; trade-off parameter  $\lambda$

**Estimate:**  $\pi$  via  $\hat{\pi} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y_i = 1\}$

**Estimate:**  $\eta: \mathcal{X} \rightarrow [0, 1]$  using appropriate CPE on  $\{x_i, y_i\}_{i=1}^n$

**Estimate:**  $\bar{\eta}_{EO}: (\mathcal{X}, Y) \rightarrow [0, 1]$  using appropriate CPE on  $S$

**Compute:**  $\hat{s}(x) = \left\{1 - \frac{\lambda}{\hat{\pi}} (\hat{\eta}_{EO}(x, 1) - \bar{c})\right\} \hat{\eta}(x) - c$

**Return:**  $\hat{f}(x) = H_{\alpha}(\hat{s}(x))$  for any  $\alpha \in [0, 1]$

---

## 3. Theory

In this section, we analyse the asymptotic and non-asymptotic properties of the plug-in algorithm (Algorithm 1). Recall from Problem 2.1, that our goal is to minimise the fairness aware cost-sensitive risk, which for a given choice of cost parameters  $c, \bar{c} \in [0, 1]^2$  and trade-off parameter  $\lambda \in \mathbb{R}$  is given by  $CS(f; \mathcal{D}, c) - \lambda CS(f; \bar{\mathcal{D}}_{EO}, \bar{c})$ .

In the language of (Narasimhan et al., 2014), we introduce the notion of performance measure. A performance measure, defined w.r.t a distribution  $\mathbb{P}$ , and performance metric  $\Psi$ , is a mapping from the space of measurable functions  $\mathcal{F}$  to the reals, i.e.,  $\mathfrak{P}_{\mathbb{P}}^{\Psi}: \mathcal{F} \rightarrow \mathbb{R}$ . In our setting, the performance metric is simply given by the negative of the objective function of Problem 2.1. Unless stated otherwise, we will denote  $\bar{\mathcal{D}}_{EO}$  by  $\bar{\mathcal{D}}$  and  $\bar{\eta}_{EO}$  by  $\bar{\eta}$  from hereon. Our performance measure is given by:

$$\begin{aligned} \mathfrak{P}_{\mathbb{P}}^{\Psi}(f) &= \Psi[TPR_{\mathcal{D}}(f), TNR_{\mathcal{D}}(f), \pi, TPR_{\bar{\mathcal{D}}}(f), \\ &TNR_{\bar{\mathcal{D}}}(f), \beta] = -\{CS(f; \mathcal{D}, c) - \lambda CS(f; \bar{\mathcal{D}}, \bar{c})\} \end{aligned}$$

Thus, a classifier's performance measure explains its merit with regards to the combined, fairness-utility objective, appropriately balanced by cost and trade-off parameters. Performance metric  $\Psi$ , makes explicit that our performance measure is a function of the classifier's TPRs and TNRs with respect to  $\mathcal{D}$  and  $\bar{\mathcal{D}}_{EO}$ , as well as distributional quantities  $\pi$  and  $\beta$ . The CSRs, and thus the performance measure, are linear in TPRs, TNRs and class probabilities implying the performance measure is continuous in the arguments of  $\Psi$ . The regret of a classifier  $f$ , w.r.t performance measure  $\mathfrak{P}_{\mathbb{P}}^{\Psi}$  is defined as:  $\text{regret}_{\mathbb{P}}^{\Psi}(f) = \mathfrak{P}_{\mathbb{P}}^{\Psi,*} - \mathfrak{P}_{\mathbb{P}}^{\Psi}(f)$  where,  $\mathfrak{P}_{\mathbb{P}}^{\Psi,*} = \mathfrak{P}_{\mathbb{P}}^{\Psi}(f^*)$ . In our case,  $f^*$  is the BOC introduced in Theorem 2.2

### 3.1. Asymptotic Analysis

In this sub-section, we prove that the plug-in procedure yields an estimator  $\hat{f}$  which is  $\Psi$ -consistent, implying that  $\text{regret}_{\mathbb{P}}^{\Psi}(\hat{f}) \xrightarrow{P} 0$ , where  $\xrightarrow{P}$  denotes convergence in probability. We denote the estimators of  $\eta, \bar{\eta}$  by  $\hat{\eta}$  and  $\hat{\bar{\eta}}$  respectively. In order to proceed we make the following assumptions:

**Assumption 1**  $\mathbb{P}_{X|Y=1}(\gamma(x) \leq c), \mathbb{P}_{X|Y=-1}(\gamma(x) \leq c)$ ,

$\mathbb{P}_{X|Y=1, \bar{Y}=1}(\gamma(x) \leq c)$  and  $\mathbb{P}_{X|Y=1, \bar{Y}=-1}(\gamma(x) \leq c)$  are continuous at  $c$ , where  $\gamma(x) = (1 + \frac{\lambda \bar{c}}{\pi})\eta(x) - \frac{\lambda}{\pi}\bar{\eta}(x, 1)\eta(x)$ , i.e.,  $\gamma(x)$  is  $s^*(x)$  in Theorem 2.2 without the constant term  $c$

**Assumption 2** Class probability estimators (CPEs)  $\hat{\eta}, \hat{\bar{\eta}}$  are  $L$ -1 consistent, i.e.,  $\mathbb{E}_X[|\eta(x) - \hat{\eta}(x)|] \xrightarrow{p} 0$ ;  $\mathbb{E}_{X,Y}[|\bar{\eta}(x, y) - \hat{\bar{\eta}}(x, y)|] \xrightarrow{p} 0$

*Remark:* As noted in (Narasimhan et al., 2014; Chzhen et al., 2019), Assumption 2 is not a very strong one, as an appropriately regularized ERM yields an  $L$ -1 consistent class probability estimator for proper losses (Menon et al., 2013; Agarwal, 2013).

**Assumption 3** Domain  $\chi$  is compact and there exist constants  $a, B \in \mathbb{R}_+$ , such that the PDFs,  $f_{X|Y=-1}, f_X, f_{X|Y=1}$  satisfy  $\forall x \in \chi, 0 < a \leq f_{X|Y=-1}(x), f_X(x), f_{X|Y=1}(x) \leq B$

*Remark:* We make this assumption for technical convenience. This is akin to the 'strong density assumption' defined in (Audibert et al., 2007). This assumption is not necessary for the case when  $\bar{Y}$  is available at test time, or for either case relating to the approximate Demographic Parity criterion

We now state our key lemma that facilitates the consistency result. Denoting the estimator derived via the plug-in procedure for  $\gamma(x)$  by  $\hat{\gamma}(x) = (1 + \frac{\lambda \bar{c}}{\pi})\hat{\eta}(x) - \frac{\lambda}{\pi}\hat{\bar{\eta}}(x, 1)\hat{\eta}(x)$ , we have:

**Lemma 3.1** Provided Assumptions 2 and 3 hold,  $\hat{\gamma}$  is  $L$ -1 consistent, i.e.,  $\mathbb{E}_X[|\gamma(x) - \hat{\gamma}(x)|] \xrightarrow{p} 0$

The validity of Lemma 3.1 allows us to leverage the proof template of (Narasimhan et al., 2014) which in turn proves the plug-in algorithm's consistency.

**Theorem 3.2** Provided Assumptions 1, 2 and 3 hold, the plug-in algorithm is  $\Psi$ -consistent, i.e., the algorithm yields  $\hat{f} = \text{sign} \circ \{\hat{\gamma} - c\}$ , s.t.,  $\mathfrak{P}_{\mathbb{P}}^{\Psi}(\hat{f}) \xrightarrow{p} \mathfrak{P}_{\mathbb{P}}^{\Psi,*}$ , i.e.,  $\text{regret}_{\mathbb{P}}^{\Psi}(\hat{f}) \xrightarrow{p} 0$

*Proof sketch:* Lemma 3.1 and Assumption 1 allow us to show that the plug-in yields an estimator  $\hat{f}$  which is s.t.:

$$\begin{aligned}
 TPR_{\mathcal{D}}(\hat{f}) &\xrightarrow{p} TPR_{\mathcal{D}}(f^*); \quad TNR_{\mathcal{D}}(\hat{f}) \xrightarrow{p} TNR_{\mathcal{D}}(f^*) \\
 TPR_{\bar{\mathcal{D}}}(\hat{f}) &\xrightarrow{p} TPR_{\bar{\mathcal{D}}}(f^*); \quad TNR_{\bar{\mathcal{D}}}(\hat{f}) \xrightarrow{p} TNR_{\bar{\mathcal{D}}}(f^*)
 \end{aligned}$$

The result then follows by the Continuous Mapping Theorem (Mann & Wald, 1943), since  $\Psi$  is continuous in its arguments. Complete proofs and detailed discussion for the results presented in this section can be found in section B of the supplement.

### 3.2. Non-Asymptotic Analysis

In this section, our objective is to characterise the sample complexity requirements associated with learning a clas-

sifier that yields small regret, via the plug-in algorithm of (Menon & Williamson, 2018). In our problem formulation, the performance measure of a classifier, is a linear function of its true positive and true negative rates. This implies that the performance measure is non-decomposable, since it cannot be expressed as a summation/ expectation over individual instances. This is contrary to the case associated with most standard loss functions that feature in the ML literature, and thus the finite-sample analysis for our performance measure is non-standard. We provide a strategy that allows us to precisely relate the sample complexity of this task to the sample complexity associated with learning the regression functions,  $\eta$  and  $\bar{\eta}$ , as well as other distributional quantities. We defer the detailed derivation of this strategy to section C of the supplement. We assume in this section that,  $\pi = \mathbb{P}(Y = 1)$  is known. While we can remove this assumption and modify our analysis to obtain equivalent results, we found doing so makes the underlying algebra/ geometry much more convoluted, without adding significant insight. Thus, for simplicity, we proceed by assuming  $\pi$  is known.

Recall, by Assumption 2, that we are working in a setting wherein the class probability estimators (CPEs),  $\hat{\eta}, \hat{\bar{\eta}}$  are  $L$ -1 consistent. Convergence in the  $L$ -1 norm implies convergence in probability, so we can meaningfully define the sample complexity associated with learning the regression function  $\eta$  via the CPE  $\hat{\eta}$ :

**Definition 3.4** The sample complexity of learning  $\eta$ , is a mapping  $m_{\eta} : (0, 1)^3 \rightarrow \mathbb{N}$ , where  $m_{\eta}((\epsilon, \delta'), \delta)$  is the minimal (integer) number of training samples required to ensure that, with probability  $\geq (1 - \delta)$ :  $\mathbb{P}_X(|\eta(x) - \hat{\eta}(x)| \geq \epsilon) \leq \delta'$

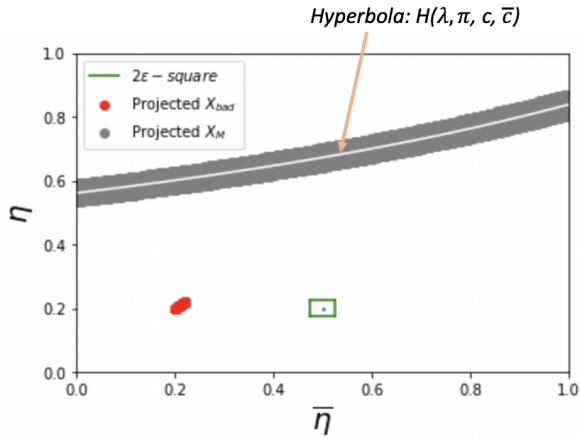
Note, we show in Lemma B.1 of the supplement that  $\hat{\bar{\eta}}(\cdot, 1)$  is also  $L$ -1 consistent. The sample complexity of learning  $\bar{\eta}(\cdot, 1)$  is thus analogously defined, and we denote this by  $m_{\bar{\eta}}$ . Our non-asymptotic result is derived via a geometric argument based in the plane of regression functions, i.e., the  $(\bar{\eta}(\cdot, 1), \eta)$ -plane. We define some key objects pertaining to our derivation. Consider in the  $(\bar{\eta}(\cdot, 1), \eta)$ -plane, the hyperbola  $H(\lambda, \pi, c, \bar{c}) := \{(1 + \frac{\lambda \bar{c}}{\pi})\eta - \frac{\lambda}{\pi}\bar{\eta}(\cdot, 1)\eta - c = 0\}$ . Also, for  $\epsilon \in (0, \frac{1}{2})$ , let  $X_M := \{x \in \chi : \text{the square of length } 2\epsilon \text{ centred at } (\bar{\eta}(x, 1), \eta(x)) \text{ intersects the hyperbola } H(\lambda, \pi, c, \bar{c}) \text{ in the } (\bar{\eta}(\cdot, 1), \eta)\text{-plane}\}$ . Having defined  $H(\lambda, \pi, c, \bar{c})$  and  $X_M$ , we now state our non-asymptotic result:

**Theorem 3.5** Let  $\delta, \delta', \epsilon \in (0, \frac{1}{2})$ . Pick any  $t > Q = 4G\{\max\{c(1 - \pi), (1 - c)\pi, |\lambda|\bar{c}(1 - \beta), |\lambda|(1 - \bar{c})\beta\}\}$ , where  $G = \max\{\frac{B}{\pi}, \frac{B}{1 - \pi}, \frac{B}{\pi\beta}, \frac{B}{\pi(1 - \beta)}\}$ , and  $B = \delta' + \mathbb{P}_X(X_M)$ . Provided access to  $n \geq \max\{m_{\eta}((\epsilon, \frac{\delta'}{2}), \frac{\delta}{8}), m_{\bar{\eta}}((\epsilon, \frac{\delta'}{2}), \frac{\delta}{8})\}$  training samples drawn i.i.d. from  $\mathbb{P}$ , the plug-in algorithm yields an estimator  $\hat{f}$ , such that, with probability at least  $(1 - \delta)$  :  $\text{regret}_{\mathbb{P}}^{\Psi}(\hat{f}) \leq t$



165 *Proof sketch:* Our proof entails showing that, for appropriate  
 166  $t$ ,  $\mathbb{P}_{S \sim \mathbb{P}^n} [\text{regret}_{\mathbb{P}}^{\mathbb{V}} > t] \leq \delta$  holds, so long as we can mean-  
 167 ingfully upper bound  $\mathbb{P}_X [\text{sign} \circ \{f^*(x)\} \neq \text{sign} \circ \{\hat{f}(x)\}]$   
 168 with probability  $\geq (1 - \frac{\delta}{4})$ . Denoting the upper bound by  $B$ ,  
 169 we characterise its form via a geometric argument. Roughly,  
 170  $B \subset \mathbb{P}_X(X_M)$ , where  $X_M$  is a specific region enclosing  
 171 the hyperbola,  $H(\lambda, \pi, c, \bar{c})$  in the  $(\bar{\eta}(\cdot, 1), \eta)$ -plane. Then  
 172 setting  $t$  as described, the result follows. A visual simu-  
 173 lation of the underlying geometry can be found in Figure 1.

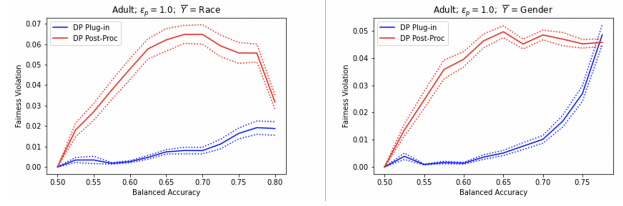
175 Theorem 3.5 tells us that the regret can be made to  
 176 decrease arbitrarily, provided a sufficient increase in the  
 177 number of training samples. The precise rate of decay,  
 178 depends on 1) the sample complexities associated with  
 179 learning the regression functions and 2) the rate at which the  
 180 probability measure, i.e.,  $\mathbb{P}_X$ , decays around the hyperbola  
 181  $H(\lambda, \pi, c, \bar{c})$  (in the  $(\bar{\eta}(\cdot, 1), \eta)$ -plane) upon shrinking the  
 182 region of consideration around it (i.e., the region akin to the  
 183 'Projected  $X_M$ ' region in Figure 1). Refer to section C of  
 184 supplement for a detailed proof and discussion.



187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201 *Figure 1.* For any point  $x \in \mathcal{X}$ : corresponding projected coordi-  
 202 nates in the  $(\bar{\eta}(\cdot, 1), \eta)$ -plane, i.e.,  $(\bar{\eta}(x, 1), \eta(x))$  lie outside  
 203 of  $\{\text{projected } X_M \cup \text{projected } X_{\text{bad}}\}$ ; we can be certain that  
 204  $\text{sign} \circ \{f^*(x)\} = \text{sign} \circ \{\hat{f}(x)\}$  - the point centred within the  
 205 green square of length  $2\epsilon$  is one such point.

## 206 4. Fairness under Differential Privacy

207  
208  
209 In this section, we work in a setting wherein there is  
 210 an additional requirement for our modelling pipeline to  
 211 mitigate information leakage about the sensitive attribute,  
 212  $\bar{Y}$ . To meet such a requirement we make use of a notion of  
 213 privacy known as differential privacy (Dwork et al., 2006),  
 214 which roughly speaking, ensures that an algorithm's output  
 215 does not differ significantly on data sets that differ in only a  
 216 single instance.



217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000

### Algorithm 2 DP Plugin approach to FAL, EO setting

**Input:** Sample  $S = \{x_i, y_i, \bar{y}_i\}_{i=1}^n$  from distribution  $\mathbb{P}$ ; cost parameters  $c, \bar{c}$ ; trade-off parameter  $\lambda$ ; privacy parameter  $\epsilon_p$

**Estimate:**  $\pi$  via  $\hat{\pi} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y_i = 1\}$

**Estimate:**  $\eta: \mathcal{X} \rightarrow [0, 1]$  using appropriate CPE on  $\{x_i, y_i\}_{i=1}^n$

**Estimate:**  $\bar{\eta}_{EO}: (\mathcal{X}, Y) \rightarrow [0, 1]$  using appropriate CPE on  $S$

**Private:**  $\hat{\eta}_{EO}^{priv}$  via appropriate privacy preserving protocol yielding,  $\epsilon_p$ -DP protected  $\hat{\eta}_{EO}^{priv}$

**Compute:**  $\hat{s}^{priv}(x) = \left\{1 - \frac{\lambda}{\bar{\pi}} (\hat{\eta}_{EO}^{priv}(x, 1) - \bar{c})\right\} \hat{\eta}(x) - c$

**Return:**  $\hat{f}^{priv}(x) = H_\alpha(\hat{s}^{priv}(x))$  for any  $\alpha \in [0, 1]$

We detail the DP Plugin protocol in Algorithm 2. We compare DP Plug-in's performance against that of the private, fair, post-processing approach (DP Post-Proc) of (Jagielski et al., 2019). We found that our method outperforms the DP Post-Proc approach in three, out of four experimental set ups considered. We present our evaluations on the Adult data set (Dheeru & Taniskidou, 2017) in Figure 2. Complete details for the DP Plug-in algorithm, and for our experimental set up and methodology, can be found in section D of the supplement.

## 5. Conclusion and Future Work

Our main contributions in this paper included (1) proving the plug-in algorithm of (Menon & Williamson, 2018) is consistent, (2) characterising the sample complexity of learning fairness-aware BOCs via the plug-in algorithm, and (3) proposing an easy to deploy, privacy-preserving protocol for the plug-in algorithm. As future directions, we believe it would be valuable to extend our analysis to the case where the sensitive attribute is non-binary; the case where multiple attributes are deemed sensitive. It would also be useful to study the statistical properties of learning algorithms across other settings, such as those demanding individual fairness, model explainability, or intersections between such areas of ethical and practical importance.

## References

- Agarwal, S. Surrogate regret bounds for the area under the roc curve via strongly proper losses. In *Conference on Learning Theory*, pp. 338–353. PMLR, 2013.
- Audibert, J.-Y., Tsybakov, A. B., et al. Fast learning rates for plug-in classifiers. *The Annals of statistics*, 35(2): 608–633, 2007.
- Chen, I., Johansson, F. D., and Sontag, D. Why is my classifier discriminatory? *arXiv preprint arXiv:1805.12002*, 2018.
- Chzhen, E., Denis, C., Hebiri, M., Oneto, L., and Pontil, M. Leveraging labeled and unlabeled data for consistent fair binary classification. *arXiv preprint arXiv:1906.05082*, 2019.
- Cummings, R., Gupta, V., Kimpara, D., and Morgenstern, J. On the compatibility of privacy and fairness. In *Adjunct Publication of the 27th Conference on User Modeling, Adaptation and Personalization*, pp. 309–315, 2019.
- Denis, C. and Hebiri, M. Confidence sets with expected sizes for multiclass classification. *The Journal of Machine Learning Research*, 18(1):3571–3598, 2017.
- Devroye, L., Györfi, L., and Lugosi, G. A probabilistic theory of pattern recognition. Technical report, 1996.
- Dheeru, D. and Taniskidou, E. K. Uci machine learning repository. 2017.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pp. 265–284. Springer, 2006.
- Jagielski, M., Kearns, M., Mao, J., Oprea, A., Roth, A., Sharifi-Malvajerdi, S., and Ullman, J. Differentially private fair learning. In *International Conference on Machine Learning*, pp. 3000–3008. PMLR, 2019.
- Mann, H. B. and Wald, A. On stochastic limit and order relationships. *The Annals of Mathematical Statistics*, 14(3):217–226, 1943.
- Menon, A., Narasimhan, H., Agarwal, S., and Chawla, S. On the statistical consistency of algorithms for binary classification under class imbalance. In *International Conference on Machine Learning*, pp. 603–611. PMLR, 2013.
- Menon, A. K. and Williamson, R. C. The cost of fairness in binary classification. In *Conference on Fairness, Accountability and Transparency*, pp. 107–118. PMLR, 2018.
- Mozannar, H., Ohannessian, M., and Srebro, N. Fair learning with private demographic data. In *International Conference on Machine Learning*, pp. 7066–7075. PMLR, 2020.
- Narasimhan, H., Vaish, R., and Agarwal, S. On the statistical consistency of plug-in classifiers for non-decomposable performance measures. In *NIPS*, volume 27, pp. 1493–1501, 2014.
- Yang, Y. Minimax nonparametric classification. i. rates of convergence. *IEEE Transactions on Information Theory*, 45(7):2271–2284, 1999.
- Zhao, H. and Gordon, G. J. Inherent tradeoffs in learning fair representations. *arXiv preprint arXiv:1906.08386*, 2019.