

PAC-Bayes and Fairness: Risk and Fairness Bounds on Distribution Dependent Fair Priors

Luca Oneto¹, Michele Donini², Massimiliano Pontil²

¹DIBRIS - University of Genova - Italy

²Istituto Italiano di Teconologia - Italy

Abstract. We address the problem of algorithmic fairness: ensuring that sensitive information does not unfairly influence the outcome of a classifier. We face this issue in the PAC-Bayes framework and we present an approach which trades off and bounds the risk and the fairness of the Gibbs Classifier measured with respect to different state-of-the-art fairness measures. For this purpose, we further develop the idea that the PAC-Bayes prior can be defined based on the data-generating distribution without actually needing to know it. In particular, we define a prior and a posterior which gives more weight to functions which exhibit good generalization and fairness properties.

1 Introduction

In recent years there has been a lot of interest in the problem of enhancing learning methods with “fairness” requirements, see [1] and references therein. The general aim is to ensure that sensitive information (e.g. knowledge about gender of an individual) does not “unfairly” influence the outcome of a learning algorithm. Several notions of fairness and associated learning methods have been introduced in machine learning in the past few years, including Demographic Parity [2], Equal Odds and Equal Opportunities [3], Disparate Treatment, Impact, and Mistreatment [4]. The underlying idea behind such notions is to balance decisions of a classifier among the different sensitive groups and label sets.

Contemporary, it is well known that combining the output of several classifiers results in much better generalization performance than using any one of them alone [5]. The major open problem in this scenario is how to weight the different classifiers in order to obtain good performance and properly assess them. The PAC-Bayes approach [6] is one of the sharpest analysis frameworks in this context, since it can provide a tight bound on the risk of the Gibbs Classifier (GC). The GC chooses a classifier in a set according to a posterior distribution each time a new sample has to be classified. In particular, in the PAC-Bayes framework, a prior distribution over the classifiers must be defined before seeing the data, then, based on the available data, a posterior distribution is chosen, and the risk of the associate GC is estimated, based on the empirical risk and the divergence between the prior and posterior distributions.

The major weakness in the conventional PAC-Bayes approach is that a posterior distribution that minimizes both the empirical risk of the GC and the

divergence between prior and posterior distributions must be chosen, since this divergence is part of the bound [7]. In order to address this issue, Catoni [8] proposed a localized PAC-Bayes analysis, which exploits a Boltzmann prior distribution defined in terms of the unknown data distribution. Note that, since the prior depends on the data generating distribution, the PAC-Bayes analysis is still valid because the prior is defined before observing the data. By tuning the prior to the distribution, Catoni was able to remove the divergence term from the bound, hence significantly reducing the complexity penalty.

For this reason, in this work, we propose a distribution dependent prior and a data dependent posterior distributions which balance the trade off between accuracy and fairness of the resulting GC measured with different notions of fairness. Then we will propose a bound on both the risk and the fairness of the solution chosen according to this prior and posterior and, similar to the result of Catoni [8], we will show that it is possible to remove the divergence term from the bound.

2 Preliminaries

Let $\mathcal{D} = \{(x_1, s_1, y_1), \dots, (x_n, s_n, y_n)\}$ be a sequence of n samples drawn independently from an unknown probability distribution μ over $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$, where $\mathcal{Y} = \{\pm 1\}$ is the set of binary output labels, $\mathcal{S} = \{1, \dots, k\}$ represents group membership, and \mathcal{X} is the input space. For every $g \in \mathcal{S}$ and operator $\diamond \in \{\pm\}$, we define the subset of training points negatively or positively labeled which belongs to the group g as $\mathcal{D}_{g, \diamond} = \{(x, s, y) : (x, s, y) \in \mathcal{D}, s = g, y = \diamond 1\}$ where $n_{g, \diamond} = |\mathcal{D}_{g, \diamond}|$.

Let us consider a function (or model) $f: \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}$ chosen from a set \mathcal{F} of possible models. The error (risk) of f is measured by a prescribed $[0, 1]$ -bounded loss function $\ell: \mathbb{R} \times \mathcal{Y} \rightarrow [0, 1]$. Then the risk of f with respect to ℓ , and its empirical estimator, can be defined as $L^\ell(f) = \mathbb{E}_{x, s, y} \{\ell(f(x, s), y)\}$ and $\hat{L}^\ell(f) = \hat{\mathbb{E}}_{\mathcal{D}} \ell(f(x, s), y)$. Moreover, the risk of f , and its empirical estimator, over the negatively and positively labeled point $\diamond \in \{\pm\}$ of group membership equal to $g \in \mathcal{S}$ can be defined as $L_{g, \diamond}^\ell(f) = \mathbb{E}_{x, s, y} \{\ell(f(x, s), y) | s = g, y = \diamond 1\}$ and $\hat{L}_{g, \diamond}^\ell(f) = \hat{\mathbb{E}}_{\mathcal{D}_{g, \diamond}} \ell(f(x, s), y)$.

The fairness of the model can be measured with respect to many notions of fairness as mentioned in the introduction. In this work we choose to opt for the Equal Opportunity (EOp) and the Equal Odds (EOd). For $\diamond \in \{\pm\}$, the EOp $^\diamond$ constraint is defined as [3]

$$\mathbb{P}\{\diamond f(x, s) > 0 | s = 1, y = \diamond 1\} = \dots = \mathbb{P}\{\diamond f(x, s) > 0 | s = k, y = \diamond 1\},$$

where $\diamond \in \{\pm\}$, since we can define the EOp of the positively (EOp $^+$) or negatively (EOp $^-$) labeled samples. The EOd, instead, is just the concurrent verification of the EOp $^+$ and EOp $^-$. Since a model f , in general, will not be able to exactly fulfill the EOp $^\diamond$ with $\diamond \in \{\pm\}$ nor the EOd constraints we define the Difference of EOp $^\diamond$, namely DEOp $^\diamond(f)$, with $\diamond \in \{\pm\}$ as

$$1/k \sum_{g_1 \in \mathcal{S}} \left| \mathbb{P}\{\diamond f(x, s) > 0 | s = g_1, y = \diamond 1\} - 1/k \sum_{g_2 \in \mathcal{S}} \mathbb{P}\{\diamond f(x, s) > 0 | s = g_2, y = \diamond 1\} \right|,$$

and the Difference of EOd, namely $\text{DEOd}(f)$, which is defined as the average value between the $\text{DEOp}^+(f)$ and $\text{DEOp}^-(f)$.

Exploiting the recent work of [1] it is possible to reformulate the EOp^\diamond constraint and the $\text{DEOp}^\diamond(f)$ with $\diamond \in \{\pm\}$, and consequently the EOd and the $\text{DEOp}^-(f)$, in terms of risks when the hard loss function, namely the function which detects a classification, $\ell_h(f(x, s), y) = \mathbb{1}\{yf(x, s) \leq 0\}$ is exploited. Specifically, the EOp^\diamond constraint can be reformulated as $L_{1,\diamond}^{\ell_h}(f) = \dots = L_{k,\diamond}^{\ell_h}(f)$, and consequently the $\text{DEOp}^\diamond(f) = 1/k \sum_{g_1 \in \mathcal{S}} |L_{g_1,\diamond}^{\ell_h}(f) - \sum_{g_2 \in \mathcal{S}} L_{g_2,\diamond}^{\ell_h}(f)|$ where $\diamond \in \{\pm\}$.

The GC draws a function $f \in \mathcal{F}$, according to a probability distribution Q over \mathcal{F} , each time a label for an input $x \in \mathcal{X}$ is required. For the GC, referred as G_Q , we can define its risk $L^\ell(G_Q) = \mathbb{E}_{f \sim Q}\{L^\ell(f)\}$ together with the empirical counterpart $\widehat{L}^\ell(G_Q) = \mathbb{E}_{f \sim Q}\{\widehat{L}^\ell(f)\}$ [7] and, analogously, its risk over the negatively and positively labeled point $\diamond \in \{\pm\}$ and group membership equal to $g \in \mathcal{S}$ $L_{g,\diamond}^\ell(G_Q) = \mathbb{E}_{f \sim Q}\{L_{g,\diamond}^\ell(f)\}$ together with the empirical counterpart $\widehat{L}_{g,\diamond}^\ell(G_Q) = \mathbb{E}_{f \sim Q}\{\widehat{L}_{g,\diamond}^\ell(f)\}$. With a simple analogy we can also reformulate the EOp^\diamond constraint and the DEOp^\diamond with $\diamond \in \{\pm\}$, and consequently the EOd and the $\text{DEOp}^-(G_Q)$, for the GC. Specifically, the EOp^\diamond constraint can be reformulated as $L_{1,\diamond}^{\ell_h}(G_Q) = \dots = L_{k,\diamond}^{\ell_h}(G_Q)$, and consequently the $\text{DEOp}^\diamond(G_Q) = 1/k \sum_{g_1 \in \mathcal{S}} |L_{g_1,\diamond}^{\ell_h}(G_Q) - \sum_{g_2 \in \mathcal{S}} L_{g_2,\diamond}^{\ell_h}(G_Q)|$ where $\diamond \in \{\pm\}$. The empirical counterpart of these quantities is indicated with an hat and can be computed by replacing the deterministic quantities with their empirical estimators.

Finally, given two probability distributions Q and P over \mathcal{F} , we will denote with $\text{KL}[Q||P]$ the Kullback-Leibler Divergence (KLD) between P and Q .

2.1 PAC-Bayes Risk Bounds

Based on the preliminaries we can recall the state of the art bound on the risk of the GC¹.

Theorem 1 ([7]). *For any probability distribution P over \mathcal{F} , chosen before seeing \mathcal{D} , $\forall Q$ we have $\mathbb{P}\{|\widehat{L}^\ell(G_Q) - L^\ell(G_Q)| \geq \sqrt{1/2n}(\text{KL}[Q||P] + \ln(2\sqrt{n}/\delta))\} \leq \delta$.*

The main problem of the PAC-Bayes Theory regards the choice of P and Q . Q should fit our observations, but, at the same time, Q should be close to P , in order to minimize the KLD. The milestone result of [8], later extended by [7], proposes to use a Boltzmann prior distribution P which depends on the data generating distribution μ . In particular, let us suppose that the density function associated to P is $p(f) = Z_P e^{-\gamma L^\ell(f)}$, where $\gamma \in [0, \infty)$ and Z_P is a normalization term. Basically, this distribution gives more importance to functions that possess small risk. If we choose as posterior Q a distribution which gives more importance to functions with small empirical risk with the following density function $q(f) = Z_Q e^{-\gamma \widehat{L}^\ell(f)}$, where Z_Q is a normalization term, it can be proved that this theorem, built on the result of Theorem 1, holds.

¹Better bounds in terms of rates of convergence and constants can be derived but this is out of the scope of this work.

Theorem 2 ([7]). *Given the prior P and the posterior Q defined above, we can state that $\mathbb{P}\{\text{KL}[Q||P] \geq \text{KL}_1(\gamma, \delta, n) \doteq \gamma^2/n + \gamma\sqrt{2/n \ln(2\sqrt{n}/\delta)}\} \leq 2\delta$. Consequently, we have that $\mathbb{P}\{|\widehat{L}^\ell(G_Q) - L^\ell(G_Q)| \geq \sqrt{1/2n(\text{KL}_1(\gamma, \delta, n) + \ln(2\sqrt{n}/\delta))}\} \leq 3\delta$.*

3 Contribution: Risk and Fairness Bounds on Fair Priors

In our case the scope is not to simply fit the data minimizing the risk of the GC but we require also the fairness of the solution w.r.t the EOp^\diamond with $\diamond \in \{\pm\}$ or the EOd constraints. In other words we want to contemporary minimize the risk of the GC and the $\text{DEOp}^+(G_Q)^{2,3}$. In order to achieve this goal, first we have to bound the $\text{DEOp}^+(G_Q)$ analogously to what has been done with $\widehat{L}^\ell(G_Q)$ in Theorem 1; then we will have to define a P and an Q able to both reduce the risk, the fairness, and the KLD. Let us start with the first objective with the following theorem.

Theorem 3. *For any probability distribution P over \mathcal{F} , chosen before seeing \mathcal{D} , $\forall Q$ we have $\mathbb{P}\{|\widehat{\text{DEOp}}^+(G_Q) - \text{DEOp}^+(G_Q)| \geq \sqrt{1/2n_{1,+}(\text{KL}[Q||P] + \ln(2\sqrt{n_{1,+}}/\delta))} + \sqrt{1/2n_{2,+}(\text{KL}[Q||P] + \ln(2\sqrt{n_{2,+}}/\delta))}\} \leq 2\delta$.*

Proof. In order to prove our statement we have to note that, thanks to the reverse triangle inequality, we have that

$$\begin{aligned} |\widehat{\text{DEOp}}^+(G_Q) - \text{DEOp}^+(G_Q)| &= |\widehat{L}_{1,+}^{\ell_h}(G_Q) - \widehat{L}_{2,+}^{\ell_h}(G_Q)| - |L_{1,+}^{\ell_h}(G_Q) - L_{2,+}^{\ell_h}(G_Q)| \\ &\leq |\widehat{L}_{1,+}^{\ell_h}(G_Q) - L_{1,+}^{\ell_h}(G_Q)| + |\widehat{L}_{2,+}^{\ell_h}(G_Q) - L_{2,+}^{\ell_h}(G_Q)|, \end{aligned}$$

and by exploiting the Theorem 1 the statement of the theorem is proved. \square

At this point, we have to define our P and Q . Exploiting the idea developed in [1], a good function should minimize the empirical risk subject to fairness constraints such that

$$f : \arg \min_{f \in \mathcal{F}} \widehat{L}^\ell(f) \quad \text{s.t.} \quad \widehat{\text{DEOp}}^+(f) \leq \epsilon,$$

or equivalently, for a particular value of $\lambda \in [0, \infty]$

$$f : \arg \min_{f \in \mathcal{F}} \widehat{L}^\ell(f) + \lambda \widehat{\text{DEOp}}^+(f), \quad (1)$$

where $\epsilon \in [0, 1]$ is necessary since for $\epsilon=0$ there may be no solution [9]. Consequently ϵ and λ regulate the trade off between accuracy and fairness of the solution. Then, following the ideas in [8, 7] we propose to use the following probability density function for Q

$$q(f) = Z_Q e^{-\gamma(\widehat{L}^\ell(f) + \lambda \widehat{\text{DEOp}}^+(f))}, \quad (2)$$

²From now on we will exploit the DEOp^+ and its empirical estimator, the extension to DEOp^- and DEOd is simple.

³In order to simplify the presentation and for readability we will deal with the case when $k=2$, namely the sensitive feature can assume just two values (e.g. male and female).

and consequently the following one for P

$$p(f) = Z_P e^{-\gamma(L^\ell(f) + \lambda \text{DEOp}^+(f))}, \quad (3)$$

where $Z_Q = 1/\int_{\mathcal{F}} q(f)df$ and $Z_P = 1/\int_{\mathcal{F}} p(f)df$. Basically our posterior distribution weights more the optimal solution of the Problem (1) and exponentially less the other ones based on their distance, in terms of cost function, from the optimal one.

If the Q and P defined respectively in Eqns. (3) and (2) are exploited we can prove the following theorem.

Theorem 4. *Given the prior P and the posterior Q defined in Eqns. (3) and (2), we can state that $\mathbb{P}\{\text{KL}[Q||P] \geq U(\delta, n, n_{1,+}, n_{2,+})\} \leq 6\delta$ where*

$$\begin{aligned} U(\delta, n, n_{1,+}, n_{2,+}) &= a^2 + 2a\sqrt{b} + b, \\ a &= \gamma \left(\sqrt{1/2n} + \lambda \left(\sqrt{1/2n_{1,+}} + \sqrt{1/2n_{2,+}} \right) \right), \\ b &= 2\gamma \left(\sqrt{1/2n} \ln(2\sqrt{n}/\delta) + \lambda \left(\sqrt{1/2n_{1,+}} \ln(2\sqrt{n_{1,+}}/\delta) + \sqrt{1/2n_{2,+}} \ln(2\sqrt{n_{2,+}}/\delta) \right) \right). \end{aligned}$$

Proof. The proof consists in noting that:

$$\begin{aligned} \text{KL}[Q||P] &= \mathbb{E}_{f \sim Q} \gamma \left(L^\ell(f) - \widehat{L}^\ell(f) + \lambda \left(\text{DEOp}^+(f) - \widehat{\text{DEOp}}^+(f) \right) \right) - \ln(Z_P/Z_Q) \\ &= \gamma \left(L^\ell(G_Q) - \widehat{L}^\ell(G_Q) + \lambda \left(\text{DEOp}^+(G_Q) - \widehat{\text{DEOp}}^+(G_Q) \right) \right) \\ &\quad - \ln \left(\int_{\mathcal{F}} p(f) e^{-\gamma(L^\ell(f) - \widehat{L}^\ell(f) + \lambda(\text{DEOp}^+(f) - \widehat{\text{DEOp}}^+(f)))} df \right) \\ &\leq \gamma \left(L^\ell(G_Q) - \widehat{L}^\ell(G_Q) + \lambda \left(\text{DEOp}^+(G_Q) - \widehat{\text{DEOp}}^+(G_Q) \right) \right) \\ &\quad + \gamma \left(L^\ell(G_P) - \widehat{L}^\ell(G_P) + \lambda \left(\text{DEOp}^+(G_P) - \widehat{\text{DEOp}}^+(G_P) \right) \right), \end{aligned}$$

where the last step follows from the Jensen's inequality. By exploiting this last result and Theorems 1 and 3 we have that the following inequality holds with probability at least $(1 - 6\delta)$

$$\begin{aligned} \text{KL}[Q||P] &\leq \gamma \left(\sqrt{\frac{\text{KL}[Q||P] + \ln\left(\frac{2\sqrt{n}}{\delta}\right)}{2n}} + \lambda \left(\sqrt{\frac{\text{KL}[Q||P] + \ln\left(\frac{2\sqrt{n_{1,+}}}{\delta}\right)}{2n_{1,+}}} + \right. \right. \\ &\quad \left. \left. \sqrt{\frac{\text{KL}[Q||P] + \ln\left(\frac{2\sqrt{n_{2,+}}}{\delta}\right)}{2n_{2,+}}} \right) \right) + \gamma \left(\sqrt{\frac{\ln\left(\frac{2\sqrt{n}}{\delta}\right)}{2n}} + \lambda \left(\sqrt{\frac{\ln\left(\frac{2\sqrt{n_{1,+}}}{\delta}\right)}{2n_{1,+}}} + \sqrt{\frac{\ln\left(\frac{2\sqrt{n_{2,+}}}{\delta}\right)}{2n_{2,+}}} \right) \right). \end{aligned}$$

The statement of the theorem is obtained by solving with respect to $\text{KL}[Q||P]$. \square

By plugging the results of Theorem 4 into Theorems 1 and 3 it is possible to obtain a fully empirical bound on the risk and the DEO⁺ of the GC where the prior P and the posterior Q are defined respectively in Eqns. (3) and (2).

Note that the final rate of the bound is $O(\sqrt{\ln(\min(n_{1,+}, n_{2,+}))/\min(n_{1,+}, n_{2,+})})$, which is optimal in the general case [6, 7] (see the state-of-the-art bound of Theorem 2) since we are contemporary controlling the risk and the DEO⁺ based on three empirical estimator $L^\ell(G_Q)$, $L_{1,+}^\ell(G_Q)$, and $L_{1,+}^\ell(G_Q)$ which exploit respectively n , $n_{1,+}$, and $n_{2,+}$ samples (note that $n \geq \max(n_{1,+}, n_{2,+})$).

4 Conclusion

In this paper we dealt with the problem of ensuring that sensitive information does not unfairly influence the outcome of a classifier. In particular we dealt with this problem in the PAC-Bayes framework by proposing a prior defined in terms of the data generating distribution and a posterior defined in terms of the observed one which gives more weight to functions which exhibit good generalization and fairness properties measured with respect to different state-of-the-art notions of fairness. Then we derived bounds on both the risk and the fairness of the resulting Gibbs Classifier. Results show optimal rate on convergence, at least in the general case, and we were able to remove the divergence term from the bound.

As future work we will deal with the problem of improving the constants and the rate of convergence of the bounds in the lucky case of small empirical error [6, 7] and to test the proposed approach in real world applications.

References

- [1] M. Donini, L. Oneto, S. Ben-David, J. Shawe-Taylor, and M. Pontil. Empirical risk minimization under fairness constraints. In *Advances in Neural Information Processing Systems*, 2018.
- [2] T. Calders, F. Kamiran, and M. Pechenizkiy. Building classifiers with independency constraints. In *IEEE international conference on Data mining*, 2009.
- [3] M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, 2016.
- [4] M. B. Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *International Conference on World Wide Web*, 2017.
- [5] Z. H. Zhou. *Ensemble methods: foundations and algorithms*. Chapman and Hall/CRC, 2012.
- [6] D. McAllester. A pac-bayesian tutorial with a dropout bound. *arXiv preprint arXiv:1307.2118*, 2013.
- [7] G. Lever, F. Laviolette, and J. Shawe-Taylor. Tighter pac-bayes bounds through distribution-dependent priors. *Theoretical Computer Science*, 473:4–28, 2013.
- [8] O. Catoni. *PAC-Bayesian Supervised Classification*. Institute of Mathematical Statistics, 2007.
- [9] G. Pleiss, M. Raghavan, F. Wu, J. Kleinberg, and K. Q. Weinberger. On fairness and calibration. In *Advances in Neural Information Processing Systems*, 2017.