

Multicalibration yields better matchings

Riccardo Colini Baldeschi* Simone Di Gregorio[†] Simone Fioravanti[†]
Federico Fusco[†] Ido Guy* Daniel Haimovich* Stefano Leonardi[†]
Fridolin Linder* Lorenzo Perini* Matteo Russo[†] Niek Tax*

Abstract

Consider the problem of finding the best matching in a weighted graph where we only have access to predictions of the actual stochastic weights, based on an underlying context. If the predictor is the Bayes optimal one, then computing the best matching based on the predicted weights is optimal. However, in practice, this perfect information scenario is not realistic. Given an imperfect predictor, a suboptimal decision rule may compensate for the induced error and thus outperform the standard optimal rule.

In this paper, we propose *multicalibration* as a way to address this problem. This fairness notion requires a predictor to be unbiased on each element of a family of protected sets of contexts. Given a class of matching algorithms \mathcal{C} and any predictor γ of the edge-weights, we show how to construct a specific multicalibrated predictor $\hat{\gamma}$, with the following property. Picking the best matching based on the output of $\hat{\gamma}$ is competitive with the best decision rule in \mathcal{C} applied onto the original predictor γ . We complement this result by providing sample complexity bounds.

1 Introduction

The interplay of classical algorithms and machine learning routines in modern industry pipelines is a well-established phenomenon: optimization algorithms are used to *fit* learning models, which, in turn, are used to *generate* the input to algorithmic tasks, or to *guide* them. While one side of this synergy is well studied and understood, only recently has the theory community started investigating how machine learning may actually help to answer classical algorithmic questions.

In the *algorithms with predictions* framework (Mitzenmacher and Vassilvitskii, 2020; Balkanski et al., 2024), the goal is to investigate the extent to which additional information provided by some machine learning prediction can improve the worst-case theoretical guarantees. Clearly, if such additional information is correct, then the resulting performance should improve, while the whole pipeline should be robust with respect to bad quality predictions. Similarly, in *Data-Augmented Algorithm Design* (Balcan, 2020), the focus is on learning from data what the best algorithm is, on a specific input distribution. In this paper, we investigate a problem that is similar in spirit to these lines of work and is practically motivated. Imagine running an optimization task on an input that is not known, but whose relevant features are only predicted by some machine learning

*Meta Central Applied Science

[†]Sapienza University of Rome, Rome, Italy (Corresponding author: simone.digregorio@uniroma1.it)

black box. Our goal is to understand how such a prediction can be modified *ex-post*, to improve the overall quality of the algorithmic solution.

For instance, consider the problem of choosing the best out of k actions. The randomness of the environment is represented by $(V, X) \sim \mathcal{D}$ where $X \in \mathbb{R}^d$ represents context/features that can be used to guide the decision, and $V \in [0, 1]^k$ contains the actions' rewards. If we know \mathcal{D} , then the value-maximizing strategy entails choosing $\arg \max \mathbb{E}[V_i | X]$; however, in applications, we can only base our decision on a predictor $\gamma : \mathbb{R}^d \rightarrow [0, 1]^k$ of such quantity, generated by some black-box machine learning routine. Clearly, we are free to apply many decision rules to $\gamma(X) \in [0, 1]^k$, the most natural one being the arg max. However, it is fairly easy to construct examples where even unbiased estimators may perform arbitrarily badly.

Consider, for instance, two deterministic arms, one with value 1 and the other with value $1/\varepsilon$, while the context X is drawn independently in $[0, 1]$. The estimator γ that is always correct on the first action, but outputs $1/\varepsilon^2$ for the second arm if $X \in [0, \varepsilon]$ and zero otherwise, is unbiased (when taking the expectation with respect to X), but correctly guesses the relative ordering of the two actions with a probability ε that can be arbitrarily small. In particular, the naive decision rule “always choose the second action” strictly dominates the apparently optimal one of taking the action with the largest predicted value.

Given a family of candidate decision rules \mathcal{C} , and a black-box predictor γ , we want to find a way to *calibrate* γ in such a way that the arg max over the new predictor performs at least as well as the best decision rule in \mathcal{C} over the initial predictor γ . Note, we do not know the underlying distribution, nor how the predictor γ is actually computed, but we would like to modify the predictor in such a way that simply feeding its output in an optimization routine would perform as well as the best decision rule for γ within a given class \mathcal{C} .

1.1 Our Results

Our first contribution in this direction is conceptual: we argue that the right property for a predictor to have in this setting is a version of *multicalibration* (Hébert-Johnson et al., 2018), a notion of calibration stemming from the literature on algorithmic fairness. In words, multicalibration requires a predictor to be calibrated on each element of a family of protected sets of contexts. In this application, such a family depends on the structure of the problem at hand and the class of decision rules \mathcal{C} .

For generality, we consider the natural optimization task of finding a max-weight matching in an n -node graph where the edge-weights are stochastic, and we only have access to them via a black-box predictor γ . Given a finite class of decision rules \mathcal{C} , in Theorem 1 we show that a suitably multicalibrated predictor $\hat{\gamma}$ exists such that computing the arg max on the edges predicted by $\hat{\gamma}$ is optimal, up to an additive precision ε . This predictor can be constructed efficiently, with $\tilde{O}(n^8/\varepsilon^4 \log |\mathcal{C}|)^*$ many samples from \mathcal{D} .

An alternative to our approach would consist of estimating the expected performance of each decision rule sampling from \mathcal{D} , and then committing to the best one. While this would require fewer samples (i.e., $\tilde{O}(n^4/\varepsilon^2 \log |\mathcal{C}|)$ many), our modular result enjoys two desirable properties:

- (i) we improve on the original black-box predictor, both for the task at hand and in mean square error sense.

*The notation \tilde{O} hides terms that are polylogarithmic in $n, 1/\varepsilon$.

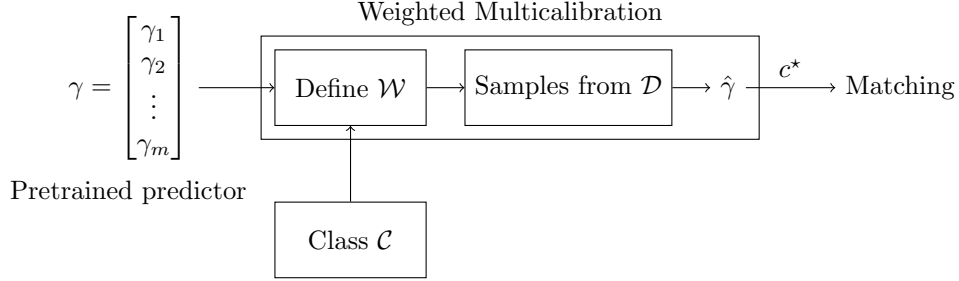


Figure 1: Proposed pipeline to obtain a matching: (1) the initial predictor γ is post-processed to a multicalibrated $\hat{\gamma}$, (2) the final matching is chosen according to c^* .

- (ii) we can decide ex-ante what the best algorithm will look like, as long as it is optimal with perfect information, meaning that we can reuse the same one for different learning instances.

Stated differently, we only need to consider the decision rules in \mathcal{C} during the preprocessing phase (in which we multicalibrate the predictor); afterwards, we only need to find the max-weight matching on the adjusted predictor. This means that in the “test phase” we only need one algorithm. The procedure is illustrated in Figure 1.

We stress that our sample complexity bound is worst-case with respect to the initial predictor γ . Since the analysis is based on a potential argument that uses the expected mean square error as potential, the closest the original predictor is to the Bayes optimal, the fewer samples are needed to multicalibrate. This is realistic for the quality of the predictors used in industry.

Beyond Matching. Our choice of max-weight matching as a running example is just for ease of presentation: our approach is flexible and applies easily to any linear maximization task with deterministic constraints (as in e.g., finding the max-weight independent set in a matroid or learning with rejection). Similarly, with proper adjustments, our analysis goes through even if we cannot solve optimally the underlying problem, but have an approximation routine.

Technical Challenges. The main technical idea revolves around adapting a notion coming from the fairness literature to an optimization task. Multicalibration is typically used to ensure that an estimator is unbiased over protected sets. Here, we use it to “protect” events of the type “a given edge e is chosen by a certain decision c ”.

1.2 Related Work

Multicalibration was first introduced in Hébert-Johnson et al. (2018), where the authors studied its sample complexity and its relationship to boosting. From its introduction, this concept has received much attention, especially because of its versatility: it has both inspired new notions like *omnipredictors* (Gopalan et al. (2022a)), *outcome indistinguishability* (Dwork et al. (2021); Gopalan et al. (2023)) and has been shown to be closely related to existing concepts in complexity theory (Casacuberta et al. (2024)). Several works have proposed variants and extensions to multicalibration as defined in Hébert-Johnson et al. (2018). Jung et al. (2021) introduce multicalibration requirements for higher moments. Gopalan et al. (2022b) propose the general framework of *weighted*

multicalibration, where two classes of functions come into play: a hypothesis class (representing the protected sets) and a weight class allowing to weaken or strengthen the probabilistic requirement as a function of the predictor. The authors frame it in general for multiclass classification: our definition of choice is an adaptation of weighted multicalibration for general vector-valued functions. Błasiok et al. (2024) introduce a more general notion of multicalibration on *auditing* functions, depending both on inputs and outputs of the predictor, which is an approach similar to our setting. Finally, Haghtalab et al. (2023) frame multicalibration in the broader context of multi-objective learning. Regarding the use of notions related to multicalibration in the context of optimization, it is worth mentioning Hu et al. (2023), where the authors adapt the definition of omnipredictor to tackle constrained loss minimization.

In an independent and concurrent work, Kiyani et al. (2025) investigate how to use predicted labels to make good decisions, when the underlying predictor is calibrated with respect to a generic family \mathcal{H} . Under some assumptions on \mathcal{H} , they find an explicit formula to ex-post modify the predicted labels so that taking the arg max (with respect to these modified predictors) is optimal, in some minimax sense. They further find a necessary and sufficient condition on \mathcal{H} so that directly taking the arg max (with respect to the predicted labels) is optimal.

2 Preliminaries

Multicalibration. Let \mathcal{X} be a feature space and $\mathcal{Y} = [0, 1]^d$ the label space. We consider point-label pairs (x, y) sampled from an unknown probability distribution \mathcal{D} supported on $\mathcal{X} \times \mathcal{Y}$. Given any predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$, *multicalibration* (Hébert-Johnson et al. (2018)) requires f to well approximate $\mathbb{E}[y|x]$ on average in each set of a given family of *protected* subsets of \mathcal{X} . To make it feasible to check this condition in practice, these sets must have some structure (i.e. finite Vapnik-Chervonenkis dimension) or be finite. Additionally, the guarantee should degrade with the probability of these sets, since we need to be able to account for low sample frequency.

The definition of multicalibration we use in this paper is arguably one of the most general ones that are present in the literature, from Gopalan et al. (2022b). We report it here for completeness, in a slightly different form, since we do not need hypothesis classes, we need binary weights and our task is a regression task: we thus do not map to the d -dimensional simplex[†]. In fact, notice that the notion of protected sets — normally induced by the hypothesis class — is in our case absorbed into the weights themselves, that are now allowed to depend only on x .

Definition 1 (Weighted Multicalibrated predictor). *Given a weight class $\mathcal{W} \subseteq \{[0, 1]^d \rightarrow \{0, 1\}^d\}$ and $\alpha \geq 0$, a predictor $f : \mathcal{X} \rightarrow [0, 1]^d$ is (\mathcal{W}, α) -multicalibrated if for every $w \in \mathcal{W}$ we have:*

$$|\mathbb{E}_{\mathcal{D}}[\langle w(f(x), x), y - f(x) \rangle]| \leq d \cdot \alpha. \quad (1)$$

Note that the expectation is with respect to the distribution \mathcal{D} , so that the above definition formalizes the intuition that f should be accurate on average, when summing over the contributions from the sets induced by each entry of a $w \in \mathcal{W}$; indeed, in our setting, by the tower property of conditional expectations, the expectation in Equation (1) does not change if y is replaced by the Bayes predictor $\mathbb{E}[y|x]$. We observe that the presence of indicator terms (the weights) inside the expectation in Equation (1) allows for this condition to be practically relevant, as it is robust against sets of low probability.

[†]Boosting-like approaches to attain multicalibration are not impacted by the switch in the setting. What changes is that we need to scale by d when using the same analysis.

3 Learning Matchings with Multicalibrated Predictions

In this section, we present our results in the maximum matching framework. We consider a complete undirected graph $G = (V, E)$ and denote with n and m the number of nodes and edges, respectively. We allow the graph to have null weights on the edges, so the assumption of completeness is without loss of generality.

3.1 Optimization Guarantees

Let \mathcal{C} be a finite and fixed family of algorithms for max-weight matching, and denote with c^* the optimal algorithm for such a combinatorial task (for instance, Edmonds algorithm—see Chapter 26 in Schrijver (2003)). We investigate a statistical scenario where $(x, y) \in \mathcal{X} \times [0, 1]^m$ are sampled i.i.d. from a joint distribution \mathcal{D} . The set \mathcal{X} is a generic space of context, while y_e for $e \in E$ represents the random weight of the edge. We are also given a generic predictor $\gamma : \mathcal{X} \rightarrow [0, 1]^m$, such that $\gamma_e(x)$ is a prediction for y_e . Note, we do not make any assumption on γ .

We want to build a new predictor $\hat{\gamma}$ that enjoys the following property:

$$\max_{c \in \mathcal{C}} \mathbb{E} \left[\sum_{e \in M_c(\gamma(x))} y_e \right] \lesssim \mathbb{E} \left[\sum_{e \in M_{c^*}(\hat{\gamma}(x))} y_e \right], \quad (2)$$

where $M_c(\gamma)$ and $M_c(\hat{\gamma})$ denote the matchings output by decision rule c on input weights γ and $\hat{\gamma}$ respectively.

We move to specifying the multicalibration setting needed to enforce Equation (2): we define the \mathcal{W} that appears in Definition 1 in two stages, one for each step of the proof of Theorem 1.

$$\begin{aligned} \mathcal{W}_1 &= \{(\mathbb{1}_{\{e \in M_c(\cdot)\}})_e\}_{c \in \mathcal{C}} \\ \mathcal{W}_2 &= \{(\mathbb{1}_{\{e \in M_{c^*}(\cdot)\}})_e\}_{e \in E} \\ \mathcal{W} &= \mathcal{W}_1 \cup \mathcal{W}_2 \end{aligned} \quad (3)$$

Notice that \mathcal{W}_2 is a singleton, with \mathcal{W}_1 and \mathcal{W}_2 being composed of functions taking values in $\{0, 1\}^m$ (with the clear equivalence between indexing with edges and indexing in $[m]$). Intuitively, the functions in \mathcal{W}_1 allow control of the error over all contexts $x \in \mathcal{X}$ that induce one $c \in \mathcal{C}$ — or c^* , for \mathcal{W}_2 — to pick a specific edge for the matching: the functions in \mathcal{W}_1 are meant to be applied onto γ , while the ones in \mathcal{W}_2 are meant to be applied onto $\hat{\gamma}$, as we show later.

We claim that if $\hat{\gamma}$ is a multicalibrated predictor w.r.t. \mathcal{W} , using c^* in the predictions given by $\hat{\gamma}(x)$ is on average (up to an additive constant depending on the α from Equation (1)) as good as taking the best $c \in \mathcal{C}$ according to γ .

Theorem 1. *Let $\varepsilon \in (0, 1)$ be any fixed precision and set the multicalibration parameter $\alpha = \varepsilon/2m$. If $\hat{\gamma}$ is (\mathcal{W}, α) -MC, it holds that:*

$$\max_{c \in \mathcal{C}} \mathbb{E} \left[\sum_{e \in M_c(\gamma(x))} y_e \right] \leq \varepsilon + \mathbb{E} \left[\sum_{e \in M_{c^*}(\hat{\gamma}(x))} y_e \right].$$

Proof. We suppress x as argument to the predictors, to simplify the notation. In the following, c is a generic matching function in \mathcal{C} . We start by arguing that the new predictor $\hat{\gamma}$ well estimates

the weight of the matching $M_c(\gamma)$, for any $c \in \mathcal{C}$. We have the following:

$$\mathbb{E} \left[\sum_e (y_e - \hat{\gamma}_e) \mathbb{1}_{\{e \in M_c(\gamma)\}} \right] = \mathbb{E} \left[\underbrace{\langle (\mathbb{1}_{\{e \in M_c(\gamma)\}})_e, y - \hat{\gamma} \rangle}_{w \in \mathcal{W}_1} \right] \leq \alpha m \quad (4)$$

The inner product appearing above is indeed upper-bounded by $\alpha \cdot m$ due to our notion of multicalibration and our choice of \mathcal{W}_1 . We now relate the expected weight of $M_c(\gamma)$ with respect to $\hat{\gamma}$ with the actual expected weight of $M_{c^*}(\hat{\gamma})$ (as measured by y_e).

$$\begin{aligned} \mathbb{E} \left[\sum_e \hat{\gamma}_e \mathbb{1}_{\{e \in M_c(\gamma)\}} \right] &\leq \mathbb{E} \left[\sum_e \hat{\gamma}_e \mathbb{1}_{\{e \in M_{c^*}(\hat{\gamma})\}} \right] && (c^*(\hat{\gamma}) \text{ optimal for } \hat{\gamma}) \\ &= \mathbb{E} \left[\sum_e (\hat{\gamma}_e - y_e) \mathbb{1}_{\{e \in M_{c^*}(\hat{\gamma})\}} \right] + \mathbb{E} \left[\sum_e y_e \mathbb{1}_{\{e \in M_{c^*}(\hat{\gamma})\}} \right] \\ &= \mathbb{E} \left[\underbrace{\langle (\mathbb{1}_{\{e \in M_{c^*}(\hat{\gamma})\}})_e, \hat{\gamma} - y \rangle}_{w \in \mathcal{W}_2} \right] + \mathbb{E} \left[\sum_e y_e \mathbb{1}_{\{e \in M_{c^*}(\hat{\gamma})\}} \right] \\ &\leq \alpha m + \mathbb{E} \left[\sum_{e \in M_{c^*}(\hat{\gamma})} y_e \right]. \end{aligned}$$

Note, the last inequality follows by our multicalibration setup. Plugging in the above inequality in Equation (4) yields the following:

$$\mathbb{E} \left[\sum_{e \in M_c(\gamma)} y_e \right] - \mathbb{E} \left[\sum_{e \in M_{c^*}(\hat{\gamma})} y_e \right] \leq 2\alpha m$$

Since the choice of $c \in \mathcal{C}$ was arbitrary, the result holds when taking the maximum over \mathcal{C} . \square

3.2 Sample Complexity Analysis

The algorithms and the complexity results that we report here are an adaptation of the ones in Gopalan et al. (2022b). The essential difference between their setting and ours, as already stressed in Section 2, is that we are dealing with functions supported on $[0, 1]^m$ (and not on the probability simplex as in the multiclass case). The algorithm proposed in Gopalan et al. (2022b) (like the ones presented in Hébert-Johnson et al. (2018); Jung et al. (2021)) is an iterative boosting procedure that initializes $\hat{\gamma}$ as a constant and then, at each iteration t , performs the following operations: (1) it checks for violations of the multicalibration conditions of the approximation $\hat{\gamma}^t$ and (2) it performs a projected gradient descent step on $\hat{\gamma}^t$. The pseudocode is reported in Algorithm 1. To make this effective in our case, some modifications are needed. First, we initialize $\hat{\gamma}$ as γ , that is, the algorithm effectively post-processes γ to achieve the multicalibration guarantees. Second, we need to project onto the entire $[0, 1]^m$: the function $\text{proj}_{[0,1]^m}(x)$ on line 13 denotes the vector such that each component $i = 1, \dots, m$ equals $\min(\max(x_i, 0), 1)$. Lastly, we need to modify the

Algorithm 1 WeightedMC($\alpha, \mathcal{W}, \{(x^j, y^j)\}_{j=1}^N, \gamma$)

```
1:  $\hat{\gamma}^0(\cdot) \leftarrow \gamma$ 
2:  $\eta \leftarrow \alpha/2$ 
3:  $mc \leftarrow \text{false}$ 
4:  $t \leftarrow 0$ 
5: while  $\neg mc$  do
6:    $mc \leftarrow \text{true}$ 
7:    $D \leftarrow \{x^j, y^j\}_{j=1}^N$ 
8:   if  $\text{CHECK}_{\mathcal{W}, \alpha, \hat{\gamma}^t}(D) = \perp$  then
9:     continue
10:  else
11:     $w_{t+1}, b_{t+1} \leftarrow \text{CHECK}_{\mathcal{W}, \alpha, \hat{\gamma}^t}(D)$ 
12:     $\delta_{t+1}(\cdot) \leftarrow b_{t+1} \cdot w_{t+1}(\hat{\gamma}^t(\cdot))$ 
13:     $\hat{\gamma}^{t+1}(\cdot) \leftarrow \text{proj}_{[0,1]^m}(\hat{\gamma}^t(\cdot) + \eta \cdot \delta_{t+1}(\cdot))$ 
14:     $mc \leftarrow \text{false}$ 
15:     $t \leftarrow t + 1$ 
16: return  $\hat{\gamma}^t$ 
```

function that checks for violations of the multicalibration requirement at iteration t , according to Definition 1.

By Definition 1, such a function should look for a $w \in \mathcal{W}$ such that $\langle w(\hat{\gamma}^t(\cdot)), y - \hat{\gamma}^t(\cdot) \rangle$ is higher than $\alpha \cdot m$ in expectation. Instead of treating this search as a learning problem, Algorithm 1 assumes oracle access to a function **CHECK**, which is a substitute for the weak agnostic learner appearing in both Hébert-Johnson et al. (2018) and Gopalan et al. (2022b)[‡]. The following is the formal definition:

Definition 2 (CHECK function). *For a data distribution \mathcal{D} supported on $\mathcal{X} \times [0, 1]^m$, a CHECK function for a weight class \mathcal{W} and a function $\hat{\gamma} : \mathcal{X} \rightarrow [0, 1]^m$ is a learning procedure that takes labeled data $D = \{(x^j, y^j)\}_{j=1}^N$, where each sample $(x, y) \sim \mathcal{D}$. For $\alpha > 0$, the learning procedure returns an element of $\mathcal{W} \times \{-1, +1\}$ or \perp :*

$$(w, b) \leftarrow \text{CHECK}_{\mathcal{W}, \alpha, \hat{\gamma}}(D)$$

satisfying the following properties:

1. *If there exists $w' \in \mathcal{W}$ such that $1/m \mathbb{E}_{\mathcal{D}}[\langle w'(\hat{\gamma}(x), x), y - \hat{\gamma}(x) \rangle] \notin [-\alpha, \alpha]$, then $w \neq \perp$ and $1/m |\mathbb{E}_{\mathcal{D}}[\langle w(\hat{\gamma}(x), x), y - \hat{\gamma}(x) \rangle]| \geq \alpha/2$ and $b = \text{sign}(\mathbb{E}_{\mathcal{D}}[\langle w(\hat{\gamma}(x), x), y - \hat{\gamma}(x) \rangle])$, where $\text{sign}(x) = x/|x|$;*
2. *If $w = \perp$, then we have that for all $w' \in \mathcal{W}$, $1/m \mathbb{E}_{\mathcal{D}}[\langle w'(\hat{\gamma}(x), x), y - \hat{\gamma}(x) \rangle] \in [-\alpha, \alpha]$.*

Definition 2 ensures that, if the algorithm converges, the final output satisfies the multicalibration conditions. If we do not make any assumption on γ , using a standard potential argument

[‡]We do not use the name *weak learner* because we are disregarding the hypotheses in this case, so the name would be technically improper, even though the idea is the same.

as in Hébert-Johnson et al. (2018), it is possible to show that the number of iterations needed by Algorithm 1 to converge is bounded by $1/m\alpha^2$.

We can actually provide a more fine-grained analysis, dependent on the mean square error of the initial predictor γ w.r.t. the Bayes predictor. The proof adapts the potential argument used in Gopalan et al. (2022b) (Lemma 5.3) to upper bound the number of iterations of the algorithm. Our results are summarized in the following Theorem.

Theorem 2. *Fix any failure probability $\delta \in (0, 1)$ and precision $\varepsilon \in (0, 1)$. Let \mathcal{C} be a collection of matching rules and γ a predictor with mean square error $r = \mathbb{E} [\|\gamma - \mathbb{E}[y|x]\|_2^2]$. Algorithm 1 enjoys the following properties:*

- *It converges in $T \leq 4r/m\alpha^2$ iterations;*
- *It requires $O\left(\frac{rn^6 \cdot \log\left(\frac{r|\mathcal{C}|}{\varepsilon\delta}\right)}{\varepsilon^4}\right)$ i.i.d. samples from distribution \mathcal{D} ;*
- *With probability $(1 - \delta)$ it returns a predictor $\hat{\gamma}$ such that:*

$$\max_{c \in \mathcal{C}} \mathbb{E} \left[\sum_{e \in M_c(\gamma(x))} y_e \right] \leq \varepsilon + \mathbb{E} \left[\sum_{e \in M_{c^*}(\hat{\gamma}(x))} y_e \right].$$

Proof. As a first step, we prove the bound on the number of iterations required for convergence. We define the value of the potential ϕ at the t^{th} iteration as

$$\phi(\hat{\gamma}^t) = \mathbb{E} [\|\hat{\gamma}^t - \mathbb{E}[y|x]\|_2^2] \quad (5)$$

By expanding this squared norm, we obtain that the decrease in potential between iteration t and $t + 1$ satisfies the following:

$$\phi(\hat{\gamma}^t) - \phi(\hat{\gamma}^{t+1}) \geq \alpha m \eta - \eta^2 m.$$

In particular, since $\eta = \alpha/2$, it simplifies to $\phi(\hat{\gamma}^t) - \phi(\hat{\gamma}^{t+1}) \geq m\alpha^2/4$. To bound T , we thus need to solve the inequality $\phi(\hat{\gamma}^0) - T \cdot m\alpha^2/4 \geq 0$. Since $\phi(\hat{\gamma}^0) = \phi(\gamma) = r$, the number of iterations linearly shrinks with r , and we get the first of our desired results.

Now we want to bound the sample complexity of statistically implementing the oracle call to CHECK. For any $w \in \mathcal{W}$, let $z_w^j = 1/m \langle w(\hat{\gamma}^t(x^j), x^j), y^j - \hat{\gamma}^t(x^j) \rangle$, denoting its empirical average and expectation with \hat{z}_w and \bar{z}_w , respectively. Clearly, $z_w^j \in [-1, 1] \forall j \in [N], w \in \mathcal{W}$, due to our normalization by m . We claim that, with enough samples, the procedure that returns the index of any $\hat{z}_w \notin [-\alpha/2, \alpha/2]$ is a correct implementation of a call to CHECK, with probability $1 - \delta_0$, where δ_0 is a failure parameter that we will fix later.

Due to Hoeffding's Inequality and union bounds over \mathcal{W} , we have:

$$\mathbb{P}(\exists w \in \mathcal{W} : |\hat{z}_w - \bar{z}_w| \geq \alpha/2) \leq 2|\mathcal{W}| \exp(-1/8N\alpha^2) = \delta_0.$$

Therefore, with probability $1 - \delta_0$, $N \in O(\log(|\mathcal{W}|/\delta_0)/\alpha^2)$ suffices to concentrate every \hat{z}_w in an interval of length α around its mean \bar{z}_w . This implies that, with the same probability, if a single $\bar{z}_w \notin [-\alpha, \alpha]$, $\hat{z}_w \notin [-\alpha/2, \alpha/2]$ and its weight w will be returned by the procedure. This proves our claim.

Now, as already argued, we call **CHECK** at most $4r/m\alpha^2$ times; therefore we need to union bound over the correctness of every call, meaning that we need to consider $\delta_0 = m\alpha^2\delta/4r$. Multiplying the number of calls by the sample complexity of each call, we get a sample complexity:

$$O\left(\frac{r \log(|\mathcal{W}|r/m\alpha\delta)}{m\alpha^4}\right).$$

Putting $\alpha = \varepsilon/m$ as required in Theorem 1, we get the claimed result. \square

Discussion of the sample complexity. Theorem 2 implies that the worst-case sample complexity (i.e. the one obtained with $r \in O(n^2)$) is $\tilde{O}(n^8/\varepsilon^4)$. Modern predictors (like γ in our setting) typically attain very small mean-squared error thanks to the enormous training datasets used in industry, which greatly ease generalization. Consequently, if the target precision in Theorems 1 and 2 is ε , it is natural to require the mean-squared error on each edge to be $O(\varepsilon^2)$. Under this realistic assumption, the sample complexity reduces to $\tilde{O}(n^8/\varepsilon^2)$.

4 Other Models

In this Section, we briefly explain how to generalize our approach to other linear optimization problems.

Finding the best action. This is the problem discussed in the introduction: there are k actions, each characterized by a value y_i , and the learner has access to a context vector $x \in \mathcal{X}$ that is drawn from a joint distribution \mathcal{D} over $\mathcal{X} \times [0, 1]^k$. Since choosing the best action is equivalent to finding the max-weight matching in a star graph of $k + 1$ nodes (a central node linked to the k action nodes), our results for matching do carry over immediately, with an improved (worst-case) sample complexity of $\tilde{O}(k^4/\varepsilon^4 \log |\mathcal{C}|)$, since there are only k edges in the graph.

Learning with Rejection. Consider now the supervised learning framework with *rejection* (Hendrickx et al., 2024), where the learner can either predict the label of the example seen, or “reject” it. Typically, the cost of rejecting a point is smaller than misclassifying it. If the underlying learning task is binary classification, this can be embedded in the best-action framework by assuming the existence of three actions: “predict YES”, “predict NO”, and “reject”. Since we only have a constant number of actions, we only need $\tilde{O}(1/\varepsilon^4 \log |\mathcal{C}|)$ samples.

Max-Weight Base. The basic setting with n elements, each characterized by a random value $y_i \in [0, 1]$ is similar to the ones above, with the difference that (i) a matroid[§] is defined on the elements and (ii) the algorithm designer can select any subset S of elements that is independent with respect to the matroid. Let r be the rank of the matroid[¶], then we can carry over the same analysis as in matching, with the difference that Equation (4) is upper bounded by $\alpha \cdot r$. This implies that setting $\alpha = \varepsilon/r$ is enough to get the same approximation result, for an overall sample complexity of $\tilde{O}(r^4/\varepsilon^4 \log |\mathcal{C}|)$.

[§]A family \mathcal{F} of subsets is called a matroid if (i) $\emptyset \in \mathcal{F}$, (ii) if $A \in \mathcal{F}$ and $B \subseteq A$, then $B \in \mathcal{F}$, and (iii) if $A, B \in \mathcal{F}$ with $|A| > |B|$, then $\exists a \in A$ such that $B \cup \{a\} \in \mathcal{F}$. In particular, maximizing a linear function with matroid constraints can be done efficiently using the greedy algorithm (Schrijver, 2003).

[¶]The cardinality of any set that is maximal with respect to the matroid property has a fixed cardinality that is called the rank of a matroid

Acknowledgments

S.D.G. is supported by the PNRR MUR project IR0000013-SoBigData.it. S.D.G., S.F., F.F., S.L., M.R. are supported by the Meta/Sapienza project on “Online Constrained Optimization and Multi-Calibration in Algorithm and Mechanism Design”.

References

- M. Balcan. Data-driven algorithm design. In *Beyond the Worst-Case Analysis of Algorithms*, pages 626–645. Cambridge University Press, 2020.
- E. Balkanski, V. Gkatzelis, and X. Tan. Mechanism design with predictions: An annotated reading list. *SIGecom Exch.*, 21(1):54–57, Oct. 2024.
- J. Błasiok, P. Gopalan, L. Hu, A. T. Kalai, and P. Nakkiran. Loss minimization yields multicalibration for large neural networks. In *15th Innovations in Theoretical Computer Science Conference (ITCS 2024)*, volume 287, pages 17–1. Schloss Dagstuhl-Leibniz Zentrum für Informatik, 2024.
- S. Casacuberta, C. Dwork, and S. P. Vadhan. Complexity-theoretic implications of multicalibration. In *Proc. of the 56th Annual ACM Symposium on Theory of Computing, STOC*, pages 1071–1082, 2024.
- C. Dwork, M. P. Kim, O. Reingold, G. N. Rothblum, and G. Yona. Outcome indistinguishability. In *Proc. of the 53rd Annual Symposium on Theory of Computing (STOC)*, pages 1095–1108, 2021.
- P. Gopalan, A. T. Kalai, O. Reingold, V. Sharan, and U. Wieder. Ominipredictors. In *13th Innovations in Theoretical Computer Science Conference (ITCS 2022)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2022a.
- P. Gopalan, M. P. Kim, M. Singhal, and S. Zhao. Low-degree multicalibration. In *Proc. of the 35th Conference on Learning Theory (COLT)*, 2022b.
- P. Gopalan, L. Hu, M. P. Kim, O. Reingold, and U. Wieder. Loss minimization through the lens of outcome indistinguishability. In *Proc. of the 14th Innovations in Theoretical Computer Science Conference, ITCS*, pages 60:1–60:20, 2023.
- N. Haghtalab, M. I. Jordan, and E. Zhao. A unifying perspective on multi-calibration: Game dynamics for multi-objective learning. In *Proc. of the 36th Annual Conference on Neural Information Processing Systems, NeurIPS*, 2023.
- U. Hébert-Johnson, M. Kim, O. Reingold, and G. Rothblum. Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning*, pages 1939–1948. PMLR, 2018.
- K. Hendrickx, L. Perini, D. V. der Plas, W. Meert, and J. Davis. Machine learning with a reject option: a survey. *Mach. Learn.*, 113(5):3073–3110, 2024.
- L. Hu, I. R. L. Navon, O. Reingold, and C. Yang. Ominipredictors for constrained optimization. In *Proc. of the 40th ICML*, 2023.

- C. Jung, C. Lee, M. M. Pai, A. Roth, and R. Vohra. Moment multicalibration for uncertainty estimation. In M. Belkin and S. Kpotufe, editors, *Proc. of the 34th Conf. on Learning Theory, COLT*, volume 134, pages 2634–2678, 2021.
- S. Kiyani, H. Hassani, G. Pappas, and A. Roth. Robust decision making with partially calibrated forecasts, 2025. URL <https://arxiv.org/abs/2510.23471>.
- M. Mitzenmacher and S. Vassilvitskii. Algorithms with predictions. In *Beyond the Worst-Case Analysis of Algorithms*, pages 646–662. Cambridge University Press, 2020.
- A. Schrijver. *Combinatorial Optimization: Polyhedra and Efficiency*, volume A. Springer, 2003.