

---

# Constant Regret, Generalized Mixability, and Mirror Descent

---

**Zakaria Mhammedi**

Research School of Computer Science  
Australian National University and DATA61  
zak.mhammedi@anu.edu.au

**Robert C. Williamson**

Research School of Computer Science  
Australian National University and DATA61  
bob.williamson@anu.edu.au

## Abstract

We consider the setting of prediction with expert advice; a learner makes predictions by aggregating those of a group of experts. Under this setting, and for the right choice of loss function and “mixing” algorithm, it is possible for the learner to achieve a constant regret regardless of the number of prediction rounds. For example, a constant regret can be achieved for *mixable* losses using the *aggregating algorithm*. The *Generalized Aggregating Algorithm* (GAA) is a name for a family of algorithms parameterized by convex functions on simplices (entropies), which reduce to the aggregating algorithm when using the *Shannon entropy*  $S$ . For a given entropy  $\Phi$ , losses for which a constant regret is possible using the GAA are called  $\Phi$ -mixable. Which losses are  $\Phi$ -mixable was previously left as an open question. We fully characterize  $\Phi$ -mixability and answer other open questions posed by [6]. We show that the Shannon entropy  $S$  is fundamental in nature when it comes to mixability; any  $\Phi$ -mixable loss is necessarily  $S$ -mixable, and the lowest worst-case regret of the GAA is achieved using the Shannon entropy. Finally, by leveraging the connection between the *mirror descent algorithm* and the update step of the GAA, we suggest a new *adaptive* generalized aggregating algorithm and analyze its performance in terms of the regret bound.

## 1 Introduction

Two fundamental problems in learning are how to aggregate information and under what circumstances can one learn fast. In this paper, we consider the problems jointly, extending the understanding and characterization of exponential mixing due to [10], who showed that not only does the “*aggregating algorithm*” learn quickly when the loss is suitably chosen, but that it is in fact a generalization of classical Bayesian updating, to which it reduces when the loss is log-loss [12]. We consider a general class of aggregating schemes, going beyond Vovk’s exponential mixing, and provide a complete characterization of the mixing behavior for general losses and general mixing schemes parameterized by an arbitrary entropy function.

In the *game of prediction with expert advice* a learner predicts the outcome of a random variable (outcome of the *environment*) by aggregating the predictions of a pool of experts. At the end of each prediction round, the outcome of the environment is announced and the learner and experts suffer losses based on their predictions. We are interested in algorithms that the learner can use to “aggregate” the experts’ predictions and minimize the *regret* at the end of the game. In this case, the regret is defined as the difference between the cumulative loss of the learner and that of the best expert in hindsight after  $T$  rounds.

The *Aggregating Algorithm* (AA) [10] achieves a constant regret — a precise notion of fast learning — for *mixable* losses; that is, the regret is bounded from above by a constant  $R_\ell$  which depends only on the loss function  $\ell$  and not on the number of rounds  $T$ . It is worth mentioning that mixability

is a weaker condition than exp-concavity, and contrary to the latter, mixability is an intrinsic, parametrization-independent notion [4].

Reid et al. [6] introduced the *Generalized Aggregating Algorithm* (GAA), going beyond the AA. The GAA is parameterized by the choice of a convex function  $\Phi$  on the simplex (entropy) and reduces to the AA when  $\Phi$  is the Shannon entropy. The GAA can achieve a constant regret for losses satisfying a certain condition called  $\Phi$ -mixability (characterizing when losses are  $\Phi$ -mixable was left as an open problem). This regret depends jointly on the *generalized mixability constant*  $\eta_\ell^\Phi$  — essentially the largest  $\eta$  such that  $\ell$  is  $(\frac{1}{\eta}\Phi)$ -mixable — and the divergence  $D_\Phi(e_\theta, \mathbf{q})$ , where  $\mathbf{q} \in \Delta_k$  is a prior distribution over  $k$  experts and  $e_\theta$  is the  $\theta$ th standard basis element of  $\mathbb{R}^k$  [6]. At each prediction round, the GAA can be divided into two steps; a *substitution step* where the learner picks a prediction from a set specified by the  $\Phi$ -mixability condition; and an *update step* where a new distribution  $\mathbf{q}$  over experts is computed depending on their performance. Interestingly, this update step is exactly the *mirror descent algorithm* [8, 5] which minimizes the weighted loss of experts.

**Contributions.** We introduce the notion of a *support loss*; given a loss  $\ell$  defined on any action space, there exists a proper loss  $\bar{\ell}$  which shares the same Bayes risk as  $\ell$ . When a loss is mixable, one can essentially work with a proper (support) loss instead — this will be the first stepping stone towards a characterization of (generalized) mixability.

The notion of  $\Phi$ -mixable and the GAA were previously restricted to finite losses. We extend these to allow for the use of losses which can take infinite values (such as the log-loss), and we show in this case that under the  $\Phi$ -mixability condition a constant regret is achievable using the GAA.

For an entropy  $\Phi$  and a loss  $\ell$ , we derive a necessary and sufficient condition (Theorems 13 and 14) for  $\ell$  to be  $\Phi$ -mixable. In particular, if  $\ell$  and  $\Phi$  satisfy some regularity conditions, then  $\ell$  is  $\Phi$ -mixable if and only if  $\eta_\ell\Phi - S$  is convex on the simplex, where  $S$  is the Shannon entropy and  $\eta_\ell$  is essentially the largest  $\eta$  such that  $\ell$  is  $\eta$ -mixable [10, 9]. This implies that a loss  $\ell$  is  $\Phi$ -mixable only if it is  $\eta$ -mixable for some  $\eta > 0$ . This, combined with the fact that  $\eta$ -mixability is equivalently  $(\frac{1}{\eta}S)$ -mixability (Theorem 12), reflects one fundamental aspect of the Shannon entropy.

Then, we derive an explicit expression for the generalized mixability constant  $\eta_\ell^\Phi$  (Corollary 17), and thus for the regret bound of the GAA. This allows us to compare the regret bound  $R_\ell^\Phi$  of any entropy  $\Phi$  with that of the Shannon entropy  $S$ . In this case, we show (Theorem 18) that  $R_\ell^S \leq R_\ell^\Phi$ ; that is, the GAA achieves the lowest worst-case regret when using the Shannon entropy — another result which reflects the fundamental nature of the Shannon entropy.

Finally, by leveraging the connection between the GAA and the mirror descent algorithm, we present a new algorithm — the *Adaptive Generalized Aggregating Algorithm* (AGAA). This algorithm consists of changing the entropy function at each prediction round similar to the *adaptive mirror descent algorithm* [8]. We analyze the performance of this algorithm in terms of its regret bound.

**Layout.** In §2, we give some background on loss functions and present new results (Theorem 4 and 5) based on the new notion of a *proper support loss*; we show that, as far as mixability is concerned, one can always work with a proper (support) loss instead of the original loss (which can be defined on an arbitrary action space). In §3, we introduce the notions of classical and generalized mixability and derive a characterization of  $\Phi$ -mixability (Theorems 13 and 14). We then introduce our new algorithm — the AGAA — and analyze its performance. We conclude the paper by a general discussion and direction for future work. All proofs, except for that of Theorem 16, are deferred to Appendix C.

**Notation.** Let  $m \in \mathbb{N}$ . We denote  $[m] := \{1, \dots, m\}$  and  $\tilde{m} := m - 1$ . We write  $\langle \cdot, \cdot \rangle$  for the standard inner product in Euclidean space. Let  $\Delta_m := \{\mathbf{p} \in [0, +\infty]^m : \langle \mathbf{p}, \mathbf{1}_m \rangle = 1\}$  be the *probability simplex* in  $\mathbb{R}^m$ , and let  $\tilde{\Delta}_m := \{\tilde{\mathbf{p}} \in [0, +\infty]^{\tilde{m}} : \langle \tilde{\mathbf{p}}, \mathbf{1}_{\tilde{m}} \rangle \leq 1\}$ . We will extensively make use of the affine map  $\Pi_m : \mathbb{R}^{\tilde{m}} \rightarrow \mathbb{R}^m$  defined by

$$\Pi_m(\mathbf{u}) := [u_1, \dots, u_{\tilde{m}}, 1 - \langle \mathbf{u}, \mathbf{1}_{\tilde{m}} \rangle]^\top. \quad (1)$$

We denote  $\text{int } C$ ,  $\text{ri } C$ , and  $\text{rbd } C$  the *interior*, *relative interior*, and *relative boundary* of a set  $C \in \mathbb{R}^m$ , respectively [2]. The *sub-differential* of a function  $f : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  at  $\mathbf{u} \in \mathbb{R}^m$  such that  $f(\mathbf{u}) < +\infty$  is defined by ([2])

$$\partial f(\mathbf{u}) := \{\mathbf{s}^* \in \mathbb{R}^m : f(\mathbf{v}) \geq f(\mathbf{u}) + \langle \mathbf{s}^*, \mathbf{v} - \mathbf{u} \rangle, \forall \mathbf{v} \in \mathbb{R}^m\}. \quad (2)$$

Table 1 on page 9 provides a list of the main symbols used in this paper.

## 2 Loss Functions

In general, a loss function is a map  $\ell: \mathcal{X} \times \mathcal{A} \rightarrow [0, +\infty]$  where  $\mathcal{X}$  is an outcome set and  $\mathcal{A}$  is an action set. In this paper, we only consider the case  $\mathcal{X} = [n]$ , i.e. finite outcome space. Overloading notation slightly, we define the mapping  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  by  $[\ell(\mathbf{a})]_x = \ell(x, \mathbf{a}), \forall x \in [n]$  and denote  $\ell_x(\cdot) := [\ell(\cdot)]_x$ . We further extend the new definition of  $\ell$  to the set  $\bigcup_{k \geq 1} \mathcal{A}^k$  such that for  $x \in [n]$  and  $A := [\mathbf{a}_\theta]_{1 \leq \theta \leq k}^\top \in \mathcal{A}^k$ ,  $\ell_x(A) := [\ell_x(\mathbf{a}_\theta)]_{1 \leq \theta \leq k}^\top \in [0, +\infty]^k$ . We define the *effective domain* of  $\ell$  by  $\text{dom } \ell := \{\mathbf{a} \in \mathcal{A}: \ell(\mathbf{a}) \in [0, +\infty]^n\}$ , and the *loss surface* by  $\mathcal{S}_\ell := \{\ell(\mathbf{a}): \mathbf{a} \in \text{dom } \ell\}$ . We say that  $\ell$  is *closed* if  $\mathcal{S}_\ell$  is closed in  $\mathbb{R}^n$ . The *superprediction* set of  $\ell$  is defined by  $\mathcal{S}_\ell^\infty := \{\ell(\mathbf{a}) + \mathbf{d}: (\mathbf{a}, \mathbf{d}) \in \mathcal{A} \times [0, +\infty]^n\}$ . Let  $\mathcal{S}_\ell := \mathcal{S}_\ell^\infty \cap [0, +\infty]^n$  be its *finite* part.

Let  $\mathbf{a}_0, \mathbf{a}_1 \in \mathcal{A}$ . The prediction  $\mathbf{a}_0$  is said to be *better* than  $\mathbf{a}_1$  if the component-wise inequality  $\ell(\mathbf{a}_0) \leq \ell(\mathbf{a}_1)$  holds and there exists some  $x \in [n]$  such that  $\ell_x(\mathbf{a}_0) < \ell_x(\mathbf{a}_1)$  [14]. A loss  $\ell$  is *admissible* if for any  $\mathbf{a} \in \mathcal{A}$  there are no better predictions.

For the rest of this paper (except for Theorem 4), we make the following assumption on losses;

**Assumption 1.**  $\ell$  is a closed, admissible loss such that  $\text{dom } \ell \neq \emptyset$ .

It is clear that there is no loss of generality in considering only admissible losses. The condition that  $\ell$  is closed is a weaker version of the more common assumption that  $\mathcal{A}$  is compact and that  $\mathbf{a} \mapsto \ell(x, \mathbf{a})$  is continuous with respect to the extended topology of  $[0, +\infty]$  for all  $x \in [n]$  [3, 1]. In fact, we do not make any explicit topological assumptions on the set  $\mathcal{A}$  ( $\mathcal{A}$  is allowed to be open in our case). Our condition simply says that if a sequence of points on the loss surface converges in  $[0, +\infty]^n$ , then there exists an action in  $\mathcal{A}$  whose image through the loss is equal to the limit. For example the 0-1 loss  $\ell_{0,1}$  is closed, yet the map  $\mathbf{p} \mapsto \ell_{0,1}(x, \mathbf{p})$  is not continuous on  $\Delta_2$ , for  $x \in \{0, 1\}$ .

In this paragraph let  $\mathcal{A}$  be the  $n$ -simplex, i.e.  $\mathcal{A} = \Delta_n$ . We define the *conditional risk*  $L_\ell: \Delta_n \times \Delta_n \rightarrow \mathbb{R}$  by  $L_\ell(\mathbf{p}, \mathbf{q}) = \mathbb{E}_{x \sim \mathbf{p}}[\ell_x(\mathbf{q})] = \langle \mathbf{p}, \ell(\mathbf{q}) \rangle$  and the *Bayes risk* by  $\underline{L}_\ell(\mathbf{p}) := \inf_{\mathbf{q} \in \Delta_n} L_\ell(\mathbf{p}, \mathbf{q})$ . In this case, the loss  $\ell$  is *proper* if  $\underline{L}_\ell(\mathbf{p}) = \langle \mathbf{p}, \ell(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \ell(\mathbf{q}) \rangle$  for all  $\mathbf{p} \neq \mathbf{q}$  in  $\Delta_n$  (and *strictly proper* if the inequality is strict). For example, the *log-loss*  $\ell_{\log}: \Delta_n \rightarrow [0, +\infty]^n$  is defined by  $\ell_{\log}(\mathbf{p}) = -\log \mathbf{p}$ , where the ‘log’ of a vector applies component-wise. One can easily check that  $\ell_{\log}$  is strictly proper. We denote  $\underline{L}_{\log}$  its Bayes risk.

The above definition of the Bayes risk is restricted to losses defined on the simplex. For a general loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ , we use the following definition;

**Definition 2** (Bayes Risk). Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss such that  $\text{dom } \ell \neq \emptyset$ . The Bayes risk  $\underline{L}_\ell: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty\}$  is defined by

$$\forall \mathbf{u} \in \mathbb{R}^n, \quad \underline{L}_\ell(\mathbf{u}) := \inf_{\mathbf{z} \in \mathcal{S}_\ell} \langle \mathbf{u}, \mathbf{z} \rangle. \quad (3)$$

The support function of a set  $\mathcal{C} \subseteq \mathbb{R}^n$  is defined by  $\sigma_{\mathcal{C}}(\mathbf{u}) := \sup_{\mathbf{z} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{z} \rangle$ ,  $\mathbf{u} \in \mathbb{R}^n$ , and thus it is easy to see that one can express the Bayes risk as  $\underline{L}_\ell(\mathbf{u}) = -\sigma_{\mathcal{S}_\ell}(-\mathbf{u})$ . Our definition of the Bayes risk is slightly different from previous ones ([3, 9, 1]) in two ways; 1) the Bayes risk is defined on all  $\mathbb{R}^n$  instead of  $[0, +\infty]^n$ ; and 2) the infimum is taken over the finite part of the superprediction set  $\mathcal{S}_\ell^\infty$ . The first point is a mere mathematical convenience and makes no practical difference since  $\underline{L}_\ell(\mathbf{p}) = -\infty$  for all  $\mathbf{p} \notin [0, +\infty]^n$ . For the second point, swapping  $\mathcal{S}_\ell$  for  $\mathcal{S}_\ell^\infty$  in (3) does not change the value of  $\underline{L}_\ell$  for mixable losses (see Appendix D). However, we chose to work with  $\mathcal{S}_\ell$  — a subset of  $\mathbb{R}^n$  — as it allows us to directly apply techniques from convex analysis.

**Definition 3** (Support Loss). We call a map  $\underline{\ell}: \Delta_n \rightarrow [0, +\infty]^n$  a support loss of  $\ell$  if

$$\begin{aligned} & \forall \mathbf{p} \in \text{ri } \Delta_n, \quad \underline{\ell}(\mathbf{p}) \in \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p}); \\ & \forall \mathbf{p} \in \text{rbd } \Delta_n, \exists (\mathbf{p}_m) \subset \text{ri } \Delta_n, \mathbf{p}_m \xrightarrow{m \rightarrow \infty} \mathbf{p} \text{ and } \underline{\ell}(\mathbf{p}_m) \xrightarrow{m \rightarrow \infty} \underline{\ell}(\mathbf{p}) \text{ component-wise,} \end{aligned}$$

where  $\partial \sigma_{\mathcal{S}_\ell}$  (see (2)) is the sub-differential of the support function —  $\sigma_{\mathcal{S}_\ell}$  — of the set  $\mathcal{S}_\ell$ .

**Theorem 4.** Any loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  such that  $\text{dom } \ell \neq \emptyset$ , has a proper support loss  $\underline{\ell}$  with the same Bayes risk,  $\underline{L}_\ell$ , as  $\ell$ .

Theorem 4 shows that regardless of the action space on which the loss is defined, there always exists a proper loss whose Bayes risk coincides with that of the original loss. This fact is useful in situations where the Bayes risk contains all the information one needs — such is the case for mixability. The next Theorem shows a stronger relationship between a loss and its corresponding support loss.

**Theorem 5.** *Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss and  $\underline{\ell}$  be a proper support loss of  $\ell$ . If the Bayes risk  $\underline{L}_\ell$  is differentiable on  $]0, +\infty[^n$ , then  $\underline{\ell}$  is uniquely defined on  $\text{ri } \Delta_n$  and*

$$\begin{aligned} \forall \mathbf{p} \in \text{dom } \underline{\ell}, \quad \exists \mathbf{a}_* \in \text{dom } \ell, \quad \ell(\mathbf{a}_*) &= \underline{\ell}(\mathbf{p}), \\ \forall \mathbf{a} \in \text{dom } \ell, \quad \exists (\mathbf{p}_m) \subset \text{ri } \Delta_n, \quad \underline{\ell}(\mathbf{p}_m) &\xrightarrow{m \rightarrow \infty} \ell(\mathbf{a}) \text{ component-wise.} \end{aligned}$$

Theorem 5 shows that when the Bayes risk is differentiable (a necessary condition for mixability — Theorem 12), the support loss is almost a reparametrization of the original loss, and in practice, it is enough to work with support losses instead. This will be crucial for characterizing  $\Phi$ -mixability.

### 3 Mixability in the Game of Prediction with Expert Advice

We consider the setting of prediction with expert advice [10]; there is a pool of  $k$  experts, parameterized by  $\theta \in [k]$ , which make predictions  $\mathbf{a}_\theta^t \in \mathcal{A}$  at each round  $t$ . In the same round, the learner predicts  $\mathbf{a}_\mathcal{M}^t := \mathcal{M}(\mathbf{a}_{1:k}^t, (x^s, \mathbf{a}_{1:k}^s)_{1 \leq s < t}) \in \mathcal{A}$ , where  $\mathbf{a}_{1:k}^t := [\mathbf{a}_\theta^t]_{1 \leq \theta \leq k}$ ,  $(x^s) \subset [n]$  are outcomes of the environment, and  $\mathcal{M}: \mathcal{A}^k \times ([n] \times \mathcal{A}^k)^* \rightarrow \mathcal{A}$  is a *merging strategy* [9]. At the end of round  $t$ ,  $x^t$  is announced and each expert  $\theta$  [resp. learner] suffers a loss  $\ell_{x^t}(\mathbf{a}_\theta)$  [resp.  $\ell_{x^t}(\mathbf{a}_\mathcal{M}^t)$ ], where  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ . After  $T > 0$  rounds, the cumulative loss of each expert  $\theta$  [resp. learner] is given by  $\text{Loss}_\theta^T(T) := \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_\theta^t)$  [resp.  $\text{Loss}_\mathcal{M}^T(T) := \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_\mathcal{M}^t)$ ]. We say that  $\mathcal{M}$  achieves a *constant regret* if  $\exists R > 0, \forall T > 0, \forall \theta \in [k], \text{Loss}_\mathcal{M}^T(T) \leq \text{Loss}_\theta^T(T) + R$ . In what follows, this game setting will be referred to by  $\mathfrak{G}_\ell^n(\mathcal{A}, k)$  and we only consider the case where  $k \geq 2$ .

#### 3.1 The Aggregating Algorithm and $\eta$ -mixability

**Definition 6** ( $\eta$ -mixability). *For  $\eta > 0$ , a loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  is said to be  $\eta$ -mixable, if  $\forall \mathbf{q} \in \Delta_k$ ,*

$$\forall \mathbf{a}_{1:k} \in \mathcal{A}^k, \exists \mathbf{a}_* \in \mathcal{A}, \forall x \in [n], \quad \ell_x(\mathbf{a}_*) \leq -\eta^{-1} \log \langle \mathbf{q}, \exp(-\eta \ell_x(\mathbf{a}_{1:k})) \rangle, \quad (4)$$

*where the  $\exp$  applies component-wise. Letting  $\mathfrak{H}_\ell := \{\eta > 0: \ell \text{ is } \eta\text{-mixable}\}$ , we define the mixability constant of  $\ell$  by  $\eta_\ell := \sup \mathfrak{H}_\ell$  if  $\mathfrak{H}_\ell \neq \emptyset$ ; and 0 otherwise.  $\ell$  is said to be mixable if  $\eta_\ell > 0$ .*

If a loss  $\ell$  is  $\eta$ -mixable for  $\eta > 0$ , the AA (Algorithm 1) achieves a constant regret in the  $\mathfrak{G}_\ell^n(\mathcal{A}, k)$  game [10]. In Algorithm 1, the map  $\mathfrak{S}_\ell: \mathcal{S}_\ell^\infty \rightarrow \mathcal{A}$  is a *substitution function* of the loss  $\ell$  [10, 4]; that is,  $\mathfrak{S}_\ell$  satisfies the component-wise inequality  $\ell(\mathfrak{S}_\ell(\mathbf{s})) \leq \mathbf{s}$ , for all  $\mathbf{s} \in \mathcal{S}_\ell^\infty$ .

It was shown by Chernov et al. [1] that the  $\eta$ -mixability condition (4) is equivalent to the convexity of the  $\eta$ -exponentiated superprediction set of  $\ell$  defined by  $\exp(-\eta \mathcal{S}_\ell^\infty) := \{\exp(-\eta \mathbf{s}): \mathbf{s} \in \mathcal{S}_\ell^\infty\}$ . Using this fact, van Erven et al. [9] showed that the mixability constant  $\eta_\ell$  of a strictly proper loss  $\ell: \Delta_n \rightarrow [0, +\infty]^n$ , whose Bayes risk is twice continuously differentiable on  $]0, +\infty[^n$ , is equal to

$$\underline{\eta}_\ell := \inf_{\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n} (\lambda_{\max}([\mathbf{H}\tilde{\underline{L}}_{\log}(\tilde{\mathbf{p}})]^{-1} \mathbf{H}\tilde{\underline{L}}_\ell(\tilde{\mathbf{p}})))^{-1}, \quad (5)$$

where  $\mathbf{H}$  is the Hessian operator and  $\tilde{\underline{L}} := \underline{L} \circ \Pi_n$  ( $\Pi_n$  was defined in (1)). The next theorem extends this result by showing that the mixability constant  $\eta_\ell$  of any loss  $\ell$  is lower bounded by  $\underline{\eta}_\ell$  in (5), as long as  $\ell$  satisfies Assumption 1 and its Bayes risk is twice differentiable.

**Theorem 7.** *Let  $\eta > 0$  and  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss. Suppose that  $\text{dom } \ell = \mathcal{A}$  and that  $\underline{L}_\ell$  is twice differentiable on  $]0, +\infty[^n$ . If  $\underline{\eta}_\ell > 0$  then  $\ell$  is  $\underline{\eta}_\ell$ -mixable. In particular,  $\eta_\ell \geq \underline{\eta}_\ell$ .*

We later show that, under the same conditions as Theorem 7, we actually have  $\eta_\ell = \underline{\eta}_\ell$  (Theorem 16) which indicates that the Bayes risk contains all the information necessary to characterize mixability.

**Remark 8.** *In practice, the requirement ‘ $\text{dom } \ell = \mathcal{A}$ ’ is not necessarily a strict restriction to finite losses; it is often the case that a loss  $\bar{\ell}: \bar{\mathcal{A}} \rightarrow [0, +\infty]^n$  only takes infinite values on the relative boundary of  $\bar{\mathcal{A}}$  (such is the case for the log-loss defined on the simplex), and thus the restriction  $\ell := \bar{\ell}|_{\mathcal{A}}$ , where  $\mathcal{A} = \text{ri } \bar{\mathcal{A}}$ , satisfies  $\text{dom } \ell = \mathcal{A}$ . It follows trivially from the definition of mixability (4) that if  $\ell$  is  $\eta$ -mixable and  $\bar{\ell}$  is continuous with respect to the extended topology of  $[0, +\infty]^n$  — a condition often satisfied — then  $\bar{\ell}$  is also  $\eta$ -mixable.*

### 3.2 The Generalized Aggregating Algorithm and $(\eta, \Phi)$ -mixability

A function  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  is an *entropy* if it is convex, its epigraph  $\text{epi } \Phi := \{(\mathbf{u}, h): \Phi(\mathbf{u}) \leq h\}$  is closed in  $\mathbb{R}^k \times \mathbb{R}$ , and  $\Delta_k \subseteq \text{dom } \Phi := \{\mathbf{u} \in \mathbb{R}^k: \Phi(\mathbf{u}) < +\infty\}$ . For example, the *Shannon entropy* is defined by  $S(\mathbf{q}) = +\infty$  if  $\mathbf{q} \notin [0, +\infty]^k$ , and

$$\forall \mathbf{q} \in [0, +\infty]^k, \quad S(\mathbf{q}) = \sum_{i \in [k]: q_i \neq 0} q_i \log q_i, \quad (6)$$

The *divergence* generated by an entropy  $\Phi$  is the map  $D_\Phi: \mathbb{R}^n \times \text{dom } \Phi \rightarrow [0, +\infty]$  defined by

$$D_\Phi(\mathbf{v}, \mathbf{u}) := \begin{cases} \Phi(\mathbf{v}) - \Phi(\mathbf{u}) - \Phi'(\mathbf{u}; \mathbf{v} - \mathbf{u}), & \text{if } \mathbf{v} \in \text{dom } \Phi; \\ +\infty, & \text{otherwise.} \end{cases} \quad (7)$$

where  $\Phi'(\mathbf{u}; \mathbf{v} - \mathbf{u}) := \lim_{\lambda \downarrow 0} [\Phi(\mathbf{u} + \lambda(\mathbf{v} - \mathbf{u})) - \Phi(\mathbf{u})]/\lambda$  (the limit exists since  $\Phi$  is convex [7]).

**Definition 9** ( $\Phi$ -mixability). *Let  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  be an entropy. A loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  is  $(\eta, \Phi)$ -mixable for  $\eta > 0$  if  $\forall \mathbf{q} \in \Delta_k, \forall \mathbf{a}_{1:k} \in \mathcal{A}^k, \exists \mathbf{a}_* \in \mathcal{A}$ , such that*

$$\forall x \in [n], \ell_x(\mathbf{a}_*) \leq \text{Mix}_\Phi^\eta(\ell_x(\mathbf{a}_{1:k}), \mathbf{q}) := \inf_{\hat{\mathbf{q}} \in \Delta_k} \langle \hat{\mathbf{q}}, \ell_x(\mathbf{a}_{1:k}) \rangle + \eta^{-1} D_\Phi(\hat{\mathbf{q}}, \mathbf{q}). \quad (8)$$

When  $\eta = 1$ , we simply say that  $\ell$  is  $\Phi$ -mixable and we denote  $\text{Mix}_\Phi := \text{Mix}_\Phi^1$ . Letting  $\mathfrak{H}_\ell^\Phi := \{\eta > 0: \ell \text{ is } (\eta, \Phi)\text{-mixable}\}$ , we define the generalized mixability constant of  $(\ell, \Phi)$  by  $\eta_\ell^\Phi := \sup \mathfrak{H}_\ell^\Phi$ , if  $\mathfrak{H}_\ell^\Phi \neq \emptyset$ ; and 0 otherwise.

Reid et al. [6] introduced the GAA (see Algorithm 2) which uses an entropy function  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  and a substitution function  $\mathfrak{S}_\ell$  (see previous section) to specify the learner's merging strategy  $\mathfrak{M}$ . It was shown that the GAA reduces to the AA when  $\Phi$  is the Shannon entropy  $S$ . It was also shown that under some regularity conditions on  $\Phi$ , the GAA achieves a constant regret in the  $\mathfrak{S}_\ell^n(\mathcal{A}, k)$  game for any finite,  $(\eta, \Phi)$ -mixable loss.

Our definition of  $\Phi$ -mixability differs slightly from that of Reid et al. [6] — we use directional derivatives to define the divergence  $D_\Phi$ . This distinction makes it possible to extend the GAA to losses which can take infinite values (such as the log-loss defined on the simplex). We show, in this case, that a constant regret is still achievable under the  $(\eta, \Phi)$ -mixability condition. Before presenting this result, we define the notion of  $\Delta$ -differentiability; for  $I \subseteq [k]$ , let  $\Delta_I := \{\mathbf{q} \in \Delta_k: q_\theta = 0, \forall \theta \notin I\}$ . We say that an entropy  $\Phi$  is  $\Delta$ -differentiable if  $\forall I \subseteq [k], \forall \mathbf{u}, \mathbf{u}_0 \in \text{ri } \Delta_I$ , the map  $\mathbf{z} \mapsto \Phi'(\mathbf{u}; \mathbf{z})$  is linear on  $\mathcal{L}_I^0 := \{\lambda(\mathbf{v} - \mathbf{u}_0): (\lambda, \mathbf{v}) \in \mathbb{R} \times \Delta_I\}$ .

**Theorem 10.** *Let  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  be a  $\Delta$ -differentiable entropy. Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss (not necessarily finite) such that  $\underline{L}_\ell$  is twice differentiable on  $]0, +\infty[^n$ . If  $\ell$  is  $(\eta, \Phi)$ -mixable then the GAA achieves a constant regret in the  $\mathfrak{S}_\ell^n(\mathcal{A}, k)$  game; for any sequence  $(x^t, \mathbf{a}_{1:k}^t)_{t=1}^T$ ,*

$$\text{Loss}_{\text{GAA}}^\ell(T) - \min_{\theta \in [k]} \text{Loss}_\theta^\ell(T) \leq R_\ell^\Phi := \inf_{\mathbf{q} \in \Delta_k} \max_{\theta \in [k]} D_\Phi(\mathbf{e}_\theta, \mathbf{q})/\eta_\ell^\Phi, \quad (9)$$

for initial distribution over experts  $\mathbf{q}^0 = \text{argmin}_{\mathbf{q} \in \Delta_k} \max_{\theta \in [k]} D_\Phi(\mathbf{e}_\theta, \mathbf{q})$ , where  $\mathbf{e}_\theta$  is the  $\theta$ th basis element of  $\mathbb{R}^k$ , and any substitution function  $\mathfrak{S}_\ell$ .

Looking at Algorithm 2, it is clear that the GAA is divided into two steps; 1) a *substitution step* which consists of finding a prediction  $\mathbf{a}_* \in \mathcal{A}$  satisfying the mixability condition (8) using a substitution function  $\mathfrak{S}_\ell$ ; and 2) an *update step* where a new distribution over experts is computed. Except for the case of the AA with the log-loss (which reduces to Bayesian updating [12]), there is not a unique choice of substitution function in general. An example of substitution function  $\mathfrak{S}_\ell$  is the *inverse loss* [13]. Kamalaruban et al. [4] discuss other alternatives depending on the curvature of the Bayes risk. Although the choice of  $\mathfrak{S}_\ell$  can affect the performance of the algorithm to some extent [4], the regret bound in (9) remains unchanged regardless of  $\mathfrak{S}_\ell$ . On the other hand, the update step is well defined and corresponds to a *mirror descent step* [6] (we later use this fact to suggest a new algorithm).

Algorithm 1: Aggregating Algorithm	Algorithm 2: Generalized Aggregating Algorithm
<b>input</b> : $q^0 \in \Delta_k; \eta > 0$ ; A $\eta$ -mixable loss $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ ; A substitution function $\mathfrak{S}_\ell$ . <b>output</b> : Learner's predictions $(a_*^t)$ <b>for</b> $t = 1$ <b>to</b> $T$ <b>do</b> Observe $A^t = a_{1:k}^t \in \mathcal{A}^k$ ; $a_*^t \leftarrow \mathfrak{S}_\ell \left( -\frac{1}{\eta} \log \sum_{\theta \in [k]} q_{\theta}^{t-1} e^{-\eta \ell(a_{\theta}^t)} \right)$ ; Observe outcome $x^t \in [n]$ ; $q_{\theta}^t \leftarrow \frac{q_{\theta}^{t-1} \exp(-\eta \ell_{x^t}(a_{\theta}^t))}{\langle q^{t-1}, \exp(-\eta \ell_{x^t}(A^t)) \rangle}, \forall \theta \in [k]$ ; <b>end</b>	<b>input</b> : $q^0 \in \Delta_k$ ; A $\Delta$ -differentiable entropy $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ ; $\eta > 0$ ; A $(\eta, \Phi)$ -mixable loss $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ ; A substitution function $\mathfrak{S}_\ell$ . <b>output</b> : Learner's predictions $(a_*^t)$ <b>for</b> $t = 1$ <b>to</b> $T$ <b>do</b> Observe $A^t = a_{1:k}^t \in \mathcal{A}^k$ ; $a_*^t \leftarrow \mathfrak{S}_\ell \left( [\text{Mix}_{\Phi}^{\eta}(\ell_x(A^t), q^{t-1})]_{1 \leq x \leq n}^T \right)$ ; Observe outcome $x^t \in [n]$ ; $q^t \leftarrow \underset{\mu \in \Delta_k}{\text{argmin}} \langle \mu, \ell_{x^t}(A^t) \rangle + \frac{1}{\eta} D_{\Phi}(\mu, q^{t-1})$ ; <b>end</b>

We conclude this subsection with two new and important results which will lead to a characterization of  $\Phi$ -mixability. The first result shows that  $(\eta, S)$ -mixability is equivalent to  $\eta$ -mixability, and the second rules out losses and entropies for which  $\Phi$ -mixability is not possible.

**Theorem 11.** *Let  $\eta > 0$ . A loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  is  $\eta$ -mixable if and only if  $\ell$  is  $(\eta, S)$ -mixable.*

**Proposition 12.** *Let  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  be an entropy and  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ . If  $\ell$  is  $\Phi$ -mixable, then the Bayes risk satisfies  $\underline{L}_{\ell} \in C^1([0, +\infty]^n)$ . If, additionally,  $\underline{L}_{\ell}$  is twice differentiable on  $]0, +\infty[^n$ , then  $\Phi$  must be strictly convex on  $\Delta_k$ .*

It should be noted that since the Bayes risk of a loss  $\ell$  must be differentiable for it to be  $\Phi$ -mixable for some entropy  $\Phi$ , Theorem 5 says that we can essentially work with a proper support loss  $\underline{\ell}$  of  $\ell$ . This will be crucial in the proof of the sufficient condition of  $\Phi$ -mixability (Theorem 14).

### 3.3 A Characterization of $\Phi$ -Mixability

In this subsection, we first show that given an entropy  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  and a loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  satisfying certain regularity conditions,  $\ell$  is  $\Phi$ -mixable if and only if

$$\boxed{\eta_{\ell} \Phi - S \text{ is convex on } \Delta_k.} \quad (10)$$

**Theorem 13.** *Let  $\eta > 0$ ,  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  a  $\eta$ -mixable loss, and  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  an entropy. If  $\eta_{\ell} \Phi - S$  is convex on  $\Delta_k$ , then  $\ell$  is  $\Phi$ -mixable.*

The converse of Theorem 13 also holds under additional smoothness conditions on  $\Phi$  and  $\ell$ ;

**Theorem 14.** *Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss such that  $\underline{L}_{\ell}$  is twice differentiable on  $]0, +\infty[^n$ , and  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  an entropy such that  $\tilde{\Phi} := \Phi \circ \Pi_k$  is twice differentiable on  $\text{int } \Delta_k$ . Then  $\ell$  is  $\Phi$ -mixable only if  $\eta_{\ell} \Phi - S$  is convex on  $\Delta_k$ .*

As consequence of Theorem 14, if a loss  $\ell$  is not classically mixable, i.e.  $\eta_{\ell} = 0$ , it cannot be  $\Phi$ -mixable for any entropy  $\Phi$ . This is because  $\eta_{\ell} \Phi - S \stackrel{*}{=} \underline{\eta}_{\ell} \Phi - S = -S$  is not convex (where equality “\*” is due to Theorem 7).

We need one more result before arriving at (10); Recall that the mixability constant  $\eta_{\ell}$  is defined as the supremum of the set  $\mathfrak{H}_{\ell} := \{\eta \geq 0: \ell \text{ is } \eta\text{-mixable}\}$ . The next lemma essentially gives a sufficient condition for this supremum to be attained when  $\mathfrak{H}_{\ell}$  is non-empty — in this case,  $\ell$  is  $\eta_{\ell}$ -mixable.

**Lemma 15.** *Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss. If  $\text{dom } \ell = \mathcal{A}$ , then either  $\mathfrak{H}_{\ell} = \emptyset$  or  $\eta_{\ell} \in \mathfrak{H}_{\ell}$ .*

**Theorem 16.** *Let  $\ell$  and  $\Phi$  be as in Theorem 14 with  $\text{dom } \ell = \mathcal{A}$ . Then  $\eta_{\ell} = \underline{\eta}_{\ell}$ . Furthermore,  $\ell$  is  $\Phi$ -mixable if and only if  $\eta_{\ell} \Phi - S$  is convex on  $\Delta_k$ .*

*Proof.* Suppose now that  $\ell$  is mixable. By Lemma 15, it follows that  $\ell$  is  $\eta_{\ell}$ -mixable, and from Theorem 11,  $\ell$  is  $(\eta_{\ell}^{-1} S)$ -mixable. Substituting  $\Phi$  for  $\eta_{\ell}^{-1} S$  in Theorem 14 implies that  $(\eta_{\ell}/\eta_{\ell} - 1) S$  is convex on  $\text{ri } \Delta_k$ . Thus,  $\eta_{\ell} \leq \underline{\eta}_{\ell}$ , and since from Theorem 7  $\underline{\eta}_{\ell} \leq \eta_{\ell}$ , we conclude that  $\eta_{\ell} = \underline{\eta}_{\ell}$ .

From Theorem 14, if  $\ell$  is  $\Phi$ -mixable then  $\eta_\ell \Phi - S$  is convex on  $\Delta_k$ . Now suppose that  $\eta_\ell \Phi - S$  is convex on  $\Delta_k$ . This implies that  $\eta_\ell > 0$ , and thus from Theorem 7,  $\ell$  is  $\eta_\ell$ -mixable. Now since  $\ell$  is  $\eta_\ell$ -mixable and  $\eta_\ell \Phi - S$  is convex on  $\Delta_k$ , Theorem 13 implies that  $\ell$  is  $\Phi$ -mixable.  $\square$

Note that the condition ‘ $\text{dom } \ell = \mathcal{A}$ ’ is in practice not a restriction to finite losses — see Remark 8. Theorem 16 implies that under the regularity conditions of Theorem 14, the Bayes risk  $\underline{L}_\ell$  [resp.  $(\underline{L}_\ell, \Phi)$ ] contains all necessary information to characterize classical [resp. generalized] mixability.

**Corollary 17** (The Generalized Mixability Constant). *Let  $\ell$  and  $\Phi$  be as in Theorem 16. Then the generalized mixability constant (see Definition 9) is given by*

$$\eta_\ell^\Phi = \eta_\ell \inf_{\tilde{q} \in \text{int } \tilde{\Delta}_k} \lambda_{\min}(\mathbf{H}\tilde{\Phi}(\tilde{q})(\mathbf{H}\tilde{S}(\tilde{q}))^{-1}), \quad (11)$$

where  $\tilde{\Phi} := \Phi \circ \Pi_k$ ,  $\tilde{S} = S \circ \Pi_k$ , and  $\Pi_k$  is defined in (1).

Observe that when  $\Phi = S$ , (11) reduces to  $\eta_\ell^S = \eta_\ell$  as expected from Theorem 11 and Theorem 16.

### 3.4 The (In)dependence Between $\ell$ and $\Phi$ and the Fundamental Nature of $S$

So far, we showed that the  $\Phi$ -mixability of losses satisfying Assumption 1 is characterized by the convexity of  $\eta\Phi - S$ , where  $\eta \in ]0, \eta_\ell]$  (see Theorems 13 and 14). As a result, and contrary to what was conjectured previously [6], the generalized mixability condition does not induce a correspondence between losses and entropies; for a given loss  $\ell$ , there is no particular entropy  $\Phi^\ell$  — specific to the choice of  $\ell$  — which minimizes the regret of the GAA. Rather, the Shannon entropy  $S$  minimizes the regret regardless of the choice of  $\ell$  (see Theorem 18 below). This reflects one fundamental aspect of the Shannon entropy.

Nevertheless, given a loss  $\ell$  and entropy  $\Phi$ , the curvature of the loss surface  $\mathcal{S}_\ell$  determines the maximum ‘learning rate’  $\eta_\ell^\Phi$  of the GAA; the curvature of  $\mathcal{S}_\ell$  is linked to  $\eta_\ell$  through the Hessian of the Bayes risk (see Theorem 30 in Appendix H.2), which is in turn linked to  $\eta_\ell^\Phi$  through (11).

Given a loss  $\ell$ , we now use the expression of  $\eta_\ell^\Phi$  in (11) to explicitly compare the regret bounds  $R_\ell^\Phi$  and  $R_\ell^S$  achieved with the GAA (see (9)) using entropy  $\Phi$  and the Shannon entropy  $S$ , respectively.

**Theorem 18.** *Let  $S, \Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ , where  $S$  is the Shannon entropy and  $\Phi$  is an entropy such that  $\tilde{\Phi} := \Phi \circ \Pi_k$  is twice differentiable on  $\text{int } \tilde{\Delta}_k$ . A loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  with  $\underline{L}_\ell$  twice differentiable on  $]0, +\infty[^n$ , is  $\Phi$ -mixable only if  $R_\ell^S \leq R_\ell^\Phi$ .*

Theorem 18 is consistent with Vovk’s result [10, §5] which essentially states that the regret bound  $R_\ell^S = \eta_\ell^{-1} \log k$  is in general tight for  $\eta$ -mixable losses.

## 4 Adaptive Generalized Aggregating Algorithm

In this section, we take advantage of the similarity between the GAA’s update step and the mirror descent algorithm (see Appendix E) to devise a modification to the GAA leading to improved regret bounds in certain cases. The GAA can be modified in (at least) two immediate ways; 1) changing the learning rate at each time step to speed-up convergence; and 2) changing the entropy, i.e. the regularizer  $\Phi$ , at each time step — similar to the *adaptive* mirror descent algorithm [8, 5]. In the former case, one can use Corollary 17 to calculate the maximum ‘learning rate’ under the  $\Phi$ -mixability constraint. Here, we focus on the second method; changing the entropy at each round. Algorithm 3 displays the modified GAA — which we call the *Adaptive Generalized Aggregating Algorithm* (AGAA) — in its most general form. In Algorithm 3,  $\Phi^*(z) := \sup_{q \in \Delta_k} \langle q, z \rangle - \Phi(q)$  is the *entropic dual* of  $\Phi$ .

Given a  $(\eta, \Phi)$ -mixable loss  $\ell$ , we verify that Algorithm 3 is well defined; for simplicity, assume that  $\text{dom } \ell = \mathcal{A}$  and  $\underline{L}_\ell$  is twice differentiable on  $]0, +\infty[^n$ . From the definition of an entropy,  $|\Phi| < +\infty$  on  $\Delta_k$ , and thus the entropic dual  $\Phi_t^*$  is defined and finite on all  $\mathbb{R}^k$  (in particular at  $\theta^t$ ). On the other hand, from Proposition 12,  $\Phi$  is strictly convex on  $\Delta_k$  which implies that  $\Phi^*$  (and thus  $\Phi_t^*$ ) is differentiable on  $\mathbb{R}^k$  (see e.g. [2, Thm. E.4.1.1]). It remains to check that  $\ell$  is  $(\eta, \Phi_t)$ -mixable. Since for  $\eta > 0$ ,  $(\eta, \Phi_t)$ -mixability is equivalent to  $(\frac{1}{\eta}\Phi_t)$ -mixability (by definition), Theorem 16 implies

---

**Algorithm 3:** Adaptive Generalized Aggregating Algorithm (AGAA)

---

**input** :  $\theta^1 = \mathbf{0} \in \mathbb{R}^k$ ; A  $\Delta$ -differentiable entropy  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ ;  $\eta > 0$ ; A  $(\eta, \Phi)$ -mixable loss  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ ; A substitution function  $\mathfrak{S}_\ell$ ; A protocol of choosing  $\beta^t$  at round  $t$ .  
**output** : Learner's predictions  $(\mathbf{a}_*^t)$

**for**  $t = 1$  **to**  $T$  **do**  
    Let  $\Phi_t(\mathbf{w}) := \Phi(\mathbf{w}) - \langle \mathbf{w}, \beta^t - \theta^t \rangle$ ; // New entropy  
    Observe  $A^t := \mathbf{a}_{1:k}^t \in \mathcal{A}^k$ ; // Experts' predictions  
     $\mathbf{a}_*^t \leftarrow \mathfrak{S}_\ell \left( [\text{Mix}_{\Phi_t}^\eta(\ell_x(A^t), \nabla \Phi_t^*(\theta^t))]_{1 \leq x \leq n}^T \right)$ ; // Learner's prediction  
    Observe  $x^t \in [n]$  and pick some  $\mathbf{v}^t \in \mathbb{R}^k$ ;  
     $\theta^{t+1} \leftarrow \theta^t - \eta \ell_{x^t}(A^t)$ ;  
**end**

---

that  $\ell$  is  $(\eta, \Phi_t)$ -mixable if and only if  $\eta \ell \eta^{-1} \Phi_t - \mathbb{S}$  is convex on  $\Delta_k$ . This is in fact the case since  $\Phi_t$  is an affine transformation of  $\Phi$ , and we have assumed that  $\ell$  is  $(\eta, \Phi)$ -mixable.

In what follows, we focus on a particular instantiation of Algorithm 3 where we choose  $\beta^t := -\eta \sum_{s=1}^{t-1} (\ell_{x^s}(A^s) + \mathbf{v}^s)$ , for some (arbitrary for now)  $(\mathbf{v}^s) \subset \mathbb{R}^k$ . The  $(\mathbf{v}^t)$  vectors act as correction terms in the update step of the AGAA. Using standard duality properties (see Appendix A), it is easy to show that the AGAA reduces to the GAA except for the update step where the new distribution over experts at round  $t \in [T]$  is now given by

$$\mathbf{q}^t = \nabla \Phi^*(\nabla \Phi(\mathbf{q}^{t-1}) - \eta \ell_{x^t}(A^t) - \eta \mathbf{v}^t).$$

**Theorem 19.** *Let  $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$  be a  $\Delta$ -differentiable entropy. Let  $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$  be a loss such that  $\underline{\ell}_\ell$  is twice differentiable on  $]0, +\infty[^n$ . Let  $\beta^t = -\eta \sum_{s=1}^{t-1} (\ell_{x^s}(A^s) + \mathbf{v}^s)$ , where  $\mathbf{v}^s \in \mathbb{R}^k$  and  $A^s := \mathbf{a}_{1:k}^s \in \mathcal{A}^k$ . If  $\ell$  is  $(\eta, \Phi)$ -mixable then for initial distribution  $\mathbf{q}^0 = \arg\min_{\mathbf{q} \in \Delta_k} \max_{\theta \in [k]} D_\Phi(e_\theta, \mathbf{q})$  and any sequence  $(x^t, \mathbf{a}_{1:k}^t)_{t=1}^T$ , the AGAA achieves the regret*

$$\forall \theta \in [k], \quad \text{Loss}_{\text{AGAA}}^\ell(T) - \text{Loss}_\theta^\ell(T) \leq R_\ell^\Phi + \Delta R_\theta(T), \quad (12)$$

where  $\Delta R_\theta(T) := \sum_{t=1}^{T-1} (v_\theta^t - \langle \mathbf{v}^t, \mathbf{q}^t \rangle)$ .

Theorem 19 implies that if the sequence  $(\mathbf{v}^t)$  is chosen such that  $\Delta R_{\theta^*}(T)$  is negative for the best expert  $\theta^*$  (in hindsight), then the regret bound ' $R_\ell^\Phi + \Delta R_{\theta^*}(T)$ ' of the AGAA is lower than that of the GAA (see (9)), and ultimately that of the AA (when  $\Phi = \mathbb{S}$ ). Unfortunately, due to Vovk's result [10, §5] there is no "universal" choice of  $(\mathbf{v}^t)$  which guarantees that  $\Delta R_{\theta^*}(T)$  is always negative. However, there are cases where this term is expected to be negative.

Consider a dataset where it is typical for the best experts (i.e., the  $\theta^*$ 's) to perform poorly at some point during the game, as measured by their average loss, for example. Under such an assumption, choosing the correction vectors  $\mathbf{v}^t$  to be negatively proportional to the average losses of experts, i.e.  $\mathbf{v}^t := -\frac{\alpha}{t} \sum_{s=1}^t \ell_{x^s}(A^s)$  (for small enough  $\alpha > 0$ ), would be consistent with the idea of making  $\Delta R_{\theta^*}(T)$  negative. To see this, suppose expert  $\theta^*$  is performing poorly during the game (say at  $t < T$ ), as measured by its instantaneous and average loss. At that point the distribution  $\mathbf{q}^t$  would put more weight on experts performing better than  $\theta^*$ , i.e. having a lower average loss. And since  $v_{\theta^*}^t$  is negatively proportional to the average loss of expert  $\theta^*$ , the quantity  $v_{\theta^*}^t - \langle \mathbf{v}^t, \mathbf{q}^t \rangle$  would be negative — consistent with making  $\Delta R_{\theta^*}(T) < 0$ . On the other hand, if expert  $\theta^*$  performs well during the game (say close to the best) then  $v_{\theta^*}^t - \langle \mathbf{v}^t, \mathbf{q}^t \rangle \simeq 0$ , since  $\mathbf{q}^t$  would put comparable weights between  $\theta^*$  and other experts (if any) with similar performance.

**Example 1.** (*A Negative Regret*). One can construct an example that illustrates the idea above. Consider the Brier game  $\mathfrak{G}_{\ell_{\text{Brier}}}^2(\Delta_2, 2)$ ; a probability game with 2 experts  $\{\theta_1, \theta_2\}$ , 2 outcomes  $\{0, 1\}$ , and where the loss  $\ell_{\text{Brier}}$  is the Brier loss [11] (which is 1-mixable). Assume that; expert  $\theta_1$  consistently predicts  $\Pr(x = 0) = 1/2$ ; expert  $\theta_2$  predicts  $\Pr(x = 0) = 1/4$  during the first 50 rounds, then switches to predicting  $\Pr(x = 0) = 3/4$  thereafter; the outcome is always  $x = 0$ . A straightforward simulation using the AGAA with the Shannon entropy, Vovk's substitution function for the Brier loss [11],  $\beta^t$  as in Theorem 19 with  $\mathbf{v}^t := -\frac{1}{8t} \sum_{s=1}^t \ell_{\text{Brier}}(x^s, A^s)$ , yields  $R_{\ell_{\text{Brier}}}^\Phi + \Delta R_{\theta^*}(T) \simeq -5$ ,



$\forall T \geq 150$ , where in this case  $\theta^* = \theta_2$  is the best expert for  $T \geq 150$ . The learner then does *better* than the best expert. If we use the AA instead, the learner does worse than  $\theta_2$  by  $\simeq R_{\text{Brier}}^S = \log 2$ .  $\square$

In real data, the situation described above — where the best expert does not necessarily perform optimally during the game — is typical, especially when the number of rounds  $T$  is large. We have tested the aggregating algorithms on real data as studied by Vovk [11]. We compared the performance of the AA with the AGAA, and found that the AGAA outperforms the AA, and in fact achieved a negative regret on two data sets. Details of the experiments are in Appendix J.

As pointed out earlier, there are situations where  $\Delta R_{\theta^*}(T) \geq 0$  even for the choice of  $(v^t)$  in Example 1, and this could potentially lead to a large positive regret for the AGAA. There is an easy way to remove this risk at a small price; the outputs of the AGAA and the AA can themselves be considered as expert predictions. These predictions can in turn be passed to a new instance of the AA to yield a *meta prediction*. The resulting worst case regret is guaranteed not to exceed that of the original AA instance by more than  $\eta^{-1} \log 2$  for an  $\eta$ -mixable loss. We test this idea in Appendix J.

## 5 Discussion and Future Work

In this work, we derived a characterization of  $\Phi$ -mixability, which enables a better understanding of when a constant regret is achievable in the game of prediction with expert advice. Then, borrowing techniques from mirror descent, we proposed a new “adaptive” version of the generalized aggregating algorithm. We derived a regret bound for a specific instantiation of this algorithm and discussed certain situations where the algorithm is expected to perform well. We empirically demonstrated the performance of this algorithm on football game predictions (see Appendix J).

Vovk [10, §5] essentially showed that given an  $\eta$ -mixable loss there is no algorithm that can achieve a lower regret bound than  $\eta^{-1} \log k$  on all sequences of outcomes. There is no contradiction in trying to design algorithms which perform well in expectation (maybe better than the AA) on “typical” data while keeping the worst case regret close to  $\eta^{-1} \log k$ . This was the motivation behind the AGAA. In future work, we will explore other choices for the correction vector  $v^t$  with the goal of lowering the (expected) bound in (12). In the present work, we did not study the possibility of varying the learning rate  $\eta$ . One might obtain better regret bounds using an adaptive learning rate as is the case with the mirror descent algorithm. Our Corollary 17 is useful in that it gives an upper bound on the maximal learning rate under the  $\Phi$ -mixability constraint. Finally, although our Theorem 18 states that worst-case regret of the GAA is minimized when using the Shannon entropy, it would be interesting to study the dynamics of the AGAA with other entropies.

Table 1: A short list of the main symbols used in the paper

Symbol	Description
$\ell$	A loss function defined on a set $\mathcal{A}$ and taking values in $[0, +\infty]^n$ (see Sec. 2)
$\mathcal{S}_\ell$	The finite part of the superprediction set of a loss $\ell$ (see Sec. 2)
$\underline{\ell}$	The support loss of a loss $\ell$ (see Def. 3)
$\underline{L}_\ell$	The Bayes risk corresponding to a loss $\ell$ (see Definition 2)
$\tilde{\underline{L}}_\ell$	The composition of the Bayes risk with an affine function; $\tilde{\underline{L}}_\ell := \underline{L}_\ell \circ \Pi_n$ (see (1))
$S$	The Shannon Entropy (see (6))
$\eta_\ell$	The mixability constant of $\ell$ (see Def. 6) ; essentially the largest $\eta$ s.t. $\ell$ is $\eta$ -mixable.
$\frac{\eta_\ell}{\eta_\ell^\Phi}$	Essentially the largest $\eta$ such that $\eta \underline{L}_\ell - \underline{L}_{\log}$ is convex (see (5) and [9])
$\eta_\ell^\Phi$	The generalized mixability constant (see Def. 9); the largest $\eta$ s.t. $\ell$ is $(\eta, \Phi)$ -mixable.
$\mathfrak{S}_\ell$	A substitution function of a loss $\ell$ (see Sec. 3.1)
$R_\ell^\Phi$	The regret achieved by the GAA using entropy $\Phi$ (see (9) and Algorithm 2)

## Acknowledgments

This work was supported by the Australian Research Council and DATA61.

## References

- [1] Alexey Chernov, Yuri Kalnishkan, Fedor Zhdanov, and Vladimir Vovk. Supermartingales in prediction with expert advice. *Theoretical Computer Science*, 411(29-30):2647–2669, 2010.
- [2] J-B. Hiriart-Urruty and C. Lemaréchal. *Fundamentals of Convex Analysis*. Springer, New York, 2001.
- [3] Yuri Kalnishkan, Volodya Vovk, and Michael V. Vyugin. Loss functions, complexities, and the Legendre transformation. *Theoretical Computer Science*, 313(2):195–207, 2004.
- [4] Parameswaran Kamalaruban, Robert Williamson, and Xinhua Zhang. Exp-concavity of proper composite losses. In *Conference on Learning Theory*, pages 1035–1065, 2015.
- [5] Francesco Orabona, Koby Crammer, and Nicolò Cesa-Bianchi. A generalized online mirror descent with applications to classification and regression. *Machine Learning*, 99(3):411–435, 2015.
- [6] Mark D. Reid, Rafael M. Frongillo, Robert C. Williamson, and Nishant Mehta. Generalized mixability via entropic duality. In *Conference on Learning Theory*, pages 1501–1522, 2015.
- [7] R. Tyrrell Rockafellar. *Convex analysis*. Princeton University Press, Princeton, NJ, 1997.
- [8] Jacob Steinhardt and Percy Liang. Adaptivity and optimism: An improved exponentiated gradient algorithm. In *International Conference on Machine Learning*, pages 1593–1601, 2014.
- [9] Tim van Erven, Mark D. Reid, and Robert C. Williamson. Mixability is Bayes risk curvature relative to log loss. *Journal of Machine Learning Research*, 13:1639–1663, 2012.
- [10] Vladimir Vovk. A game of prediction with expert advice. *Journal of Computer and System Sciences*, 56(2):153–173, 1998.
- [11] Vladimir Vovk and Fedor Zhdanov. Prediction with expert advice for the Brier game. *Journal of Machine Learning Research*, 10(Nov):2445–2471, 2009.
- [12] Volodya Vovk. Competitive on-line statistics. *International Statistical Review*, 69(2):213–248, 2001.
- [13] Robert C. Williamson. The geometry of losses. In *Conference on Learning Theory*, pages 1078–1108, 2014.
- [14] Robert C. Williamson, Elodie Vernet, and Mark D. Reid. Composite multiclass losses. *Journal of Machine Learning Research*, 17(223):1–52, 2016.