
Supplementary Material for ‘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

Contents

A	Notation and Preliminaries	3
B	Technical Lemmas	4
C	Proofs of Results in the Main Body	12
C.1	Proof of Theorem 4	12
C.2	Proofs of Theorem 5 and Proposition 12	13
C.3	Proof of Theorem 7	17
C.4	Proof of Theorem 10	18
C.5	Proof of Theorem 11	21
C.6	Proof of Theorem 13	22
C.7	Proof of Theorem 14	24
C.8	Proof of Lemma 15	27
C.9	Proof of Theorem 17	27
C.10	Proof of Theorem 18	28
C.11	Proof of Theorem 19	29
D	Defining the Bayes Risk Using the Superprediction Set	30
E	The Update Step of the GAA and the Mirror Descent Algorithm	31
F	The Generalized Aggregating Algorithm Using the Shannon Entropy S	31
G	Legendre Φ, but no Φ-mixable ℓ	32
H	Loss Surface and Superprediction Set	33
H.1	Convexity of the Superprediction Set	34
H.2	Curvature of the Loss Surface	34

I	Classical Mixability Revisited	35
J	An Experiment on Football Prediction Dataset	36
J.1	Testing the AGAA	36
J.2	Testing a AA-AGAA Meta-Learner	37

A Notation and Preliminaries

For $n \in \mathbb{N}$, we define $\tilde{n} = n - 1$. We denote $[n] := \{1, \dots, n\}$ the set of integers between 1 and n . Let $\langle \cdot, \cdot \rangle$ denote the standard inner product in \mathbb{R}^n and $\|\cdot\|$ the corresponding norm. Let I_n and $\mathbf{1}_n$ denote the $n \times n$ identity matrix and the vector of all ones in \mathbb{R}^n . Let e_1, \dots, e_n denote the *standard basis* for \mathbb{R}^n . For a set $\mathcal{I} \subseteq \mathbb{N}$ and $\mathbf{r}_1, \dots, \mathbf{r}_n \in \mathbb{R}^k$, we denote $[\mathbf{r}_i]_{i \in \mathcal{I}} := [\mathbf{r}_{i_1}, \dots, \mathbf{r}_{i_k}] \in \mathbb{R}^{n \times k}$, where $\mathcal{I} = \{i_1, \dots, i_k\}$ and $i_1 < \dots < i_k$. We denote its transpose by $[\mathbf{r}_i]_{i \in \mathcal{I}}^\top \in \mathbb{R}^{k \times n}$. For two vectors $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$, we write $\mathbf{p} \leq \mathbf{q}$ [resp. $\mathbf{p} < \mathbf{q}$], if $\forall i \in [n], p_i \leq q_i$ [resp. $p_i < q_i$]. We also denote $\mathbf{p} \odot \mathbf{q} = [p_i q_i]_{1 \leq i \leq n} \in \mathbb{R}^n$ the *Hadamard product* of \mathbf{p} and \mathbf{q} . If (\mathbf{c}_k) is a sequence of vectors in $\mathcal{C} \subseteq \mathbb{R}^n$, we simply write $(\mathbf{c}_k) \subset \mathcal{C}$. For a sequence $(\mathbf{v}_m) \subset \mathbb{R}^n$, we write $\mathbf{v}_m \xrightarrow{m \rightarrow \infty} \mathbf{v}$ or $\lim_{m \rightarrow \infty} \mathbf{v}_m = \mathbf{v}$, if $\forall i \in [n], \lim_{m \rightarrow \infty} [\mathbf{v}_m]_i = v_i$. For a square matrix $A \in \mathbb{R}^{n \times n}$, $\lambda_{\min}(A)$ [resp. $\lambda_{\max}(A)$] denotes its minimum [resp. maximum] eigenvalue. For $k \geq 1$, $\mathbf{u} \in [0, +\infty]^k$ and $\mathbf{w} \in \mathbb{R}^k$, we define $\log \mathbf{u} := [\log u_i]_{1 \leq i \leq k}^\top \in \mathbb{R}^k$ and $\exp \mathbf{w} := [\exp w_i]_{1 \leq i \leq k}^\top \in \mathbb{R}^k$.

Let $\Delta_n := \{\mathbf{p} \in [0, 1]^n : \langle \mathbf{p}, \mathbf{1}_n \rangle = 1\}$ be the *probability simplex* in \mathbb{R}^n . We also define $\tilde{\Delta}_n := \{\tilde{\mathbf{p}} \in [0, +\infty]^{\tilde{n}} : \langle \tilde{\mathbf{p}}, \mathbf{1}_{\tilde{n}} \rangle \leq 1\}$. We will use the notations $\Delta_n^k := (\Delta_n)^k$ and $\tilde{\Delta}_n^k := (\tilde{\Delta}_n)^k$. For $\mathcal{I} \subseteq [n]$, the set $\Delta_{\mathcal{I}} = \{\mathbf{q} \in \Delta_n : q_i = 0, \forall i \in [n] \setminus \mathcal{I}\}$ is a $|\mathcal{I}|$ -*face* of Δ_n . We denote $\Pi_{\mathcal{I}}^n : \mathbb{R}^n \rightarrow \mathbb{R}^{|\mathcal{I}|}$ the linear projection operator satisfying $\Pi_{\mathcal{I}}^n \mathbf{u} = [u_i]_{i \in \mathcal{I}}^\top$. If there is no ambiguity from the context, we may simply write $\Pi_{\mathcal{I}}$ instead of $\Pi_{\mathcal{I}}^n$. It is easy to verify that $\Pi_{\mathcal{I}} \Pi_{\mathcal{I}}^\top = I_{|\mathcal{I}|}$ and that $\mathbf{q} \mapsto \Pi_{\mathcal{I}} \mathbf{q}$ is a bijection from $\Delta_{\mathcal{I}} \subseteq \Delta_n$ to $\Delta_{|\mathcal{I}|}$. In the special case where $\mathcal{I} = [\tilde{n}]$, we write $\Pi_n := \Pi_{[\tilde{n}]}$ and we define the affine operator $\Pi_n : \mathbb{R}^{\tilde{n}} \rightarrow \mathbb{R}^n$ by $\Pi_n(\mathbf{u}) := [u_1, \dots, u_{\tilde{n}}, 1 - \langle \mathbf{u}, \mathbf{1}_{\tilde{n}} \rangle]^\top = J_n \mathbf{u} + \mathbf{e}_n$, where $J_n := \begin{bmatrix} I_{\tilde{n}} \\ -\mathbf{1}_{\tilde{n}}^\top \end{bmatrix} \in \mathbb{R}^{n \times \tilde{n}}$.

For $\mathbf{u} \in \mathbb{R}^n$ and $c \in \mathbb{R}$, we denote $\mathcal{H}_{\mathbf{u}, c} := \{\mathbf{y} \in \mathbb{R}^n : \langle \mathbf{y}, \mathbf{u} \rangle \leq c\}$ and $\mathcal{B}(\mathbf{u}, c) := \{\mathbf{v} \in \mathbb{R}^n : \|\mathbf{u} - \mathbf{v}\| \leq c\}$. $\mathcal{H}_{\mathbf{u}, c}$ is a closed half space and $\mathcal{B}(\mathbf{u}, c)$ is the c -ball in \mathbb{R}^n centered at \mathbf{u} . Let $\mathcal{C} \subseteq \mathbb{R}^n$ be a non-empty set. We denote $\text{int } \mathcal{C}$, $\text{ri } \mathcal{C}$, $\text{bd } \mathcal{C}$, and $\text{rbd } \mathcal{C}$ the *interior*, *relative interior*, *boundary*, and *relative boundary* of a set $\mathcal{C} \in \mathbb{R}^n$, respectively [7]. We denote the *indicator function* of \mathcal{C} by $\iota_{\mathcal{C}}$, where for $\mathbf{u} \in \mathcal{C}$, $\iota_{\mathcal{C}}(\mathbf{u}) = 0$, otherwise $\iota_{\mathcal{C}}(\mathbf{u}) = +\infty$. The *support function* of \mathcal{C} is defined by

$$\sigma_{\mathcal{C}}(\mathbf{u}) := \sup_{\mathbf{s} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{s} \rangle, \quad \mathbf{u} \in \mathbb{R}^n.$$

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$. We denote $\text{dom } f := \{\mathbf{u} \in \mathbb{R}^n : f(\mathbf{u}) < +\infty\}$ the *effective domain* of f . The function f is *proper* if $\text{dom } f \neq \emptyset$. The function f is *convex* if $\forall (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^n$ and $\lambda \in]0, 1[$, $f(\lambda \mathbf{u} + (1 - \lambda)\mathbf{v}) \leq \lambda f(\mathbf{u}) + (1 - \lambda)f(\mathbf{v})$. When the latter inequality is strict for all $\mathbf{u} \neq \mathbf{v}$, f is *strictly convex*. When f is convex, it is *closed* if it is *lower semi-continuous*; that is, for all $\mathbf{u} \in \mathbb{R}^n$, $\liminf_{\mathbf{v} \rightarrow \mathbf{u}} f(\mathbf{v}) \geq f(\mathbf{u})$. The function f is said to be *1-homogeneous* if $\forall (\mathbf{u}, \alpha) \in \mathbb{R}^n \times]0, +\infty[$, $f(\alpha \mathbf{u}) = \alpha f(\mathbf{u})$, and it is said to be *1-coercive* if $\frac{f(\mathbf{u})}{\|\mathbf{u}\|} \rightarrow +\infty$ as $\|\mathbf{u}\| \rightarrow \infty$. Let f be proper. The *sub-differential* of f is defined by

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

Any element $\mathbf{s} \in \partial f(\mathbf{u})$ is called a *sub-gradient* of f at \mathbf{u} . We say that f is *directionally differentiable* if for all $(\mathbf{u}, \mathbf{v}) \in \text{dom } f \times \mathbb{R}^n$ the limit $\lim_{t \downarrow 0} \frac{f(\mathbf{u} + t\mathbf{v}) - f(\mathbf{u})}{t}$ exists in $[-\infty, +\infty]$. In this case, we denote the limit by $f'(\mathbf{u}; \mathbf{v})$. When f is convex, it is directionally differentiable [11]. Let f be proper and directionally differentiable. The *divergence* generated by f is the map $D_f : \mathbb{R}^n \times \text{dom } f \rightarrow [0, +\infty]$ defined by

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

For $\mathcal{I} \subset [n]$ and $f_{\mathcal{I}} := f \circ \Pi_{\mathcal{I}}^\top$, it is easy to verify that $f'_{\mathcal{I}}(\Pi_{\mathcal{I}} \mathbf{p}; \Pi_{\mathcal{I}} \mathbf{q} - \Pi_{\mathcal{I}} \mathbf{p}) = f'(\mathbf{p}; \mathbf{q} - \mathbf{p})$, $\forall (\mathbf{p}, \mathbf{q}) \in \Delta_{\mathcal{I}}$. In this case, it holds that $D_f(\mathbf{q}, \mathbf{p}) = D_{f_{\mathcal{I}}}(\Pi_{\mathcal{I}} \mathbf{q}, \Pi_{\mathcal{I}} \mathbf{p})$. If f is differentiable [resp. twice differentiable] at $\mathbf{u} \in \text{int dom } f$, we denote $\nabla f(\mathbf{u}) \in \mathbb{R}^n$ [resp. $\text{Hf}(\mathbf{u}) \in \mathbb{R}^{n \times n}$] its *gradient vector* [resp. *Hessian matrix*] at \mathbf{u} . A vector-valued function $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is differentiable at \mathbf{u} if for all $i \in [m]$, g_i is differentiable at \mathbf{u} . In this case, the *differential* of g at \mathbf{u} is the linear operator $\text{D}g(\mathbf{u}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined by $\text{D}g(\mathbf{u}) := [\nabla g_i(\mathbf{u})]_{1 \leq i \leq m}^\top$. If f has k continuous derivatives on a set $\Omega \subset \mathbb{R}^k$, we write $f \in C^k(\Omega)$.

We define $\tilde{f} : \mathbb{R}^{\tilde{n}} \rightarrow \mathbb{R} \cup \{+\infty\}$ by $\tilde{f} := f \circ \Pi_n + \iota_{\tilde{\Delta}_n}$. That is,

$$\tilde{f}(\tilde{\mathbf{u}}) := \begin{cases} f(J_n \tilde{\mathbf{u}} + \mathbf{e}_n), & \text{for } \tilde{\mathbf{u}} \in \tilde{\Delta}_n; \\ +\infty, & \text{for } \tilde{\mathbf{u}} \in \mathbb{R}^{n-1} \setminus \tilde{\Delta}_n. \end{cases} \quad (1)$$

If \tilde{f} is directionally differentiable, then $f'(\mathbf{p}, \mathbf{q} - \mathbf{p}) = \tilde{f}'(\tilde{\mathbf{p}}, \tilde{\mathbf{q}} - \tilde{\mathbf{p}})$, for $\mathbf{p}, \mathbf{q} \in \Delta_n$. If \tilde{f} is differentiable at $\tilde{\mathbf{p}} = \Pi_n(\mathbf{p})$, then $\tilde{f}'(\tilde{\mathbf{p}}, \tilde{\mathbf{q}} - \tilde{\mathbf{p}}) = \langle \nabla \tilde{f}(\tilde{\mathbf{p}}), \tilde{\mathbf{q}} - \tilde{\mathbf{p}} \rangle$. If, additionally, f is differentiable at $\mathbf{p} \in \text{ri } \Delta_k$, the chain rule yields $\nabla \tilde{f}(\tilde{\mathbf{p}}) = J_n^\top \nabla f(\mathbf{p})$. Since $J_n(\tilde{\mathbf{p}} - \tilde{\mathbf{q}}) = \Pi_n(\tilde{\mathbf{p}} - \tilde{\mathbf{q}}) = \mathbf{p} - \mathbf{q}$, it also follows that $\langle \tilde{\mathbf{p}} - \tilde{\mathbf{q}}, \nabla \tilde{f}(\tilde{\mathbf{p}}) \rangle = \langle \mathbf{p} - \mathbf{q}, \nabla f(\mathbf{p}) \rangle$.

The *Fenchel dual* of a (proper) function f is defined by $f^*(\mathbf{v}) := \sup_{\mathbf{u} \in \text{dom } f} \langle \mathbf{u}, \mathbf{v} \rangle - f(\mathbf{u})$, and it is a closed, convex function on \mathbb{R}^n [7]. The following proposition gives some useful properties of the Fenchel dual which will be used in several proofs.

Proposition 1 ([7]). *Let $f, h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$. If f and h are proper and there are affine functions minorizing them on \mathbb{R}^n , then for all $\mathbf{v}_0 \in \mathbb{R}^n$*

$$\begin{aligned} (i) \quad & g(\mathbf{u}) = f(\mathbf{u}) + r, \forall \mathbf{u} & \implies & g^*(\mathbf{v}) = f^*(\mathbf{v}) - r, \forall \mathbf{v} \\ (ii) \quad & g(\mathbf{u}) = f(\mathbf{u}) + \langle \mathbf{v}_0, \mathbf{u} \rangle, \forall \mathbf{u} & \implies & g^*(\mathbf{v}) = f^*(\mathbf{v} - \mathbf{v}_0), \forall \mathbf{v} \\ (iii) \quad & f \leq h & \implies & f^* \geq h^*, \\ (iv) \quad & \mathbf{s} \in \partial f^*(\mathbf{v}), \mathbf{v} \in \mathbb{R}^n & \implies & f^*(\mathbf{v}) = \langle \mathbf{v}, \mathbf{s} \rangle - f(\mathbf{s}), \\ (v) \quad & g(\mathbf{u}) = f(t\mathbf{u}), t > 0, \forall \mathbf{u} & \implies & g^*(\mathbf{v}) = f^*(\mathbf{v}/t), \end{aligned}$$

A function $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ is an *entropy* if it is closed, convex, and $\Delta_k \subseteq \text{dom } \Phi$. Its *entropic dual* $\Phi^* : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ is defined by $\Phi^*(\mathbf{z}) := \sup_{\mathbf{q} \in \Delta_k} \langle \mathbf{q}, \mathbf{z} \rangle - \Phi(\mathbf{q})$, $\mathbf{z} \in \mathbb{R}^k$. For the remainder of this paper, we consider entropies defined on \mathbb{R}^k , where $k \geq 2$.

Let $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy and $\Phi_\Delta := \Phi + \iota_{\Delta_k}$. In this case, $\Phi^* = \Phi_\Delta^*$. It is clear that Φ_Δ is 1-coercive, and therefore, $\text{dom } \Phi^* = \text{dom } \Phi_\Delta^* = \mathbb{R}^k$ [7, Prop. E.1.3.8]. The entropic dual of Φ can also be expressed using the Fenchel dual of $\tilde{\Phi} : \mathbb{R}^{k-1} \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by (1) after substituting f by Φ and n by k . In fact,

$$\begin{aligned} \Phi^*(\mathbf{z}) &= \sup_{\tilde{\mathbf{q}} \in \tilde{\Delta}_k} \langle J_k \tilde{\mathbf{q}} + \mathbf{e}_k, \mathbf{z} \rangle - \Phi(J_k \tilde{\mathbf{q}} + \mathbf{e}_k), \\ &= \langle \mathbf{e}_k, \mathbf{z} \rangle + \sup_{\tilde{\mathbf{q}} \in \tilde{\Delta}_k} \langle \tilde{\mathbf{q}}, J_k^\top \mathbf{z} \rangle - \tilde{\Phi}(\tilde{\mathbf{q}}), \\ &= \langle \mathbf{e}_k, \mathbf{z} \rangle + \tilde{\Phi}^*(J_k^\top \mathbf{z}), \end{aligned} \quad (2)$$

where (2) follows from the fact that $\text{dom } \tilde{\Phi} = \tilde{\Delta}_k$. Note that when Φ is an entropy, $\tilde{\Phi}$ is a closed convex function on \mathbb{R}^{k-1} . Hence, it holds that $\tilde{\Phi}^{**} = \tilde{\Phi}$ [11].

The *Shannon entropy* by $S(\mathbf{q}) := \sum_{i \in [k]: q_i \neq 0} q_i \log q_i$,¹ if $\mathbf{q} \in [0, +\infty[^k$; and $+\infty$ otherwise.

We will also make use of the following lemma.

Lemma 2 ([3]). $\forall m \geq 1, \forall A, B \in \mathbb{R}^{m \times m}$, $\lambda_{\max}(AB) = \lambda_{\max}(BA)$ and $\lambda_{\min}(AB) = \lambda_{\min}(BA)$.

B Technical Lemmas

This appendix presents technical lemmas which will be needed in various proofs of results from the main body of the paper.

For an open convex set Ω in \mathbb{R}^n and $\alpha > 0$, a function $\phi : \Omega \rightarrow \mathbb{R}$ is said to be α -strongly convex if $\mathbf{u} \mapsto \phi(\mathbf{u}) - \alpha \|\mathbf{u}\|^2$ is convex on Ω [8]. The next lemma is a characterization of a generalization of α -strong convexity, where $\mathbf{u} \mapsto \|\mathbf{u}\|^2$ is replaced by any strictly convex function.

Lemma 3. *Let $\Omega \subseteq \mathbb{R}^n$ be an open convex set. Let $\phi, \psi : \Omega \rightarrow \mathbb{R}$ be twice differentiable.*

If ψ is strictly convex, then $\forall \mathbf{u} \in \Omega$, $H\psi(\mathbf{u})$ is invertible, and for any $\alpha > 0$

$$\forall \mathbf{u} \in \Omega, \lambda_{\min}(H\phi(\mathbf{u})(H\psi(\mathbf{u}))^{-1}) \geq \alpha \iff \phi - \alpha\psi \text{ is convex}, \quad (3)$$

Furthermore, if $\alpha > 1$, then the left hand side of (3) implies that $\phi - \psi$ is strictly convex.

¹The Shannon entropy is usually defined with a minus sign. However, it will be more convenient for us to work without it.

Proof. Suppose that $\inf_{\mathbf{u} \in \Omega} \lambda_{\min}(\mathbf{H}\phi(\mathbf{u})(\mathbf{H}\psi(\mathbf{u}))^{-1}) \geq \alpha$. Since g is strictly convex and twice differentiable on Ω , $\mathbf{H}\psi(\mathbf{u})$ is symmetric positive definite, and thus invertible. Therefore, there exists a symmetric positive definite matrix $G \in \mathbb{R}^{n \times n}$ such that $GG = \mathbf{H}\psi(\mathbf{u})$. Lemma 2 implies

$$\begin{aligned}
& \inf_{\mathbf{u} \in \Omega} \lambda_{\min}(\mathbf{H}\phi(\mathbf{u})(\mathbf{H}\psi(\mathbf{u}))^{-1}) && \geq \alpha, \\
\iff & \inf_{\mathbf{u} \in \Omega} \lambda_{\min}(G^{-1}\mathbf{H}\phi(\mathbf{u})G^{-1}) && \geq \alpha, \\
\iff & \forall \mathbf{u} \in \Omega, \forall \mathbf{v} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, \frac{\mathbf{v}^\top G^{-1}(\mathbf{H}\phi(\mathbf{u}))G^{-1}\mathbf{v}}{\mathbf{v}^\top \mathbf{v}} && \geq \alpha, \\
\iff & \forall \mathbf{u} \in \Omega, \forall \mathbf{w} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, \mathbf{w}^\top (\mathbf{H}\phi(\mathbf{u}))\mathbf{w} && \geq \alpha \mathbf{w}^\top GG\mathbf{w} = \mathbf{w}^\top (\alpha \mathbf{H}\psi(\mathbf{u}))\mathbf{w}, \\
\iff & \forall \mathbf{u} \in \Omega, \mathbf{H}\phi(\mathbf{u}) && \succeq \alpha \mathbf{H}\psi(\mathbf{u}), \\
\iff & \forall \mathbf{u} \in \Omega, \mathbf{H}(\phi - \alpha\psi)(\mathbf{u}) && \succeq 0,
\end{aligned}$$

where in the third and fifth lines we used the definition of minimum eigenvalue and performed the change of variable $\mathbf{w} = G^{-1}\mathbf{v}$, respectively. To conclude the proof of (3), note that the positive semi-definiteness of $\mathbf{H}(\phi - \alpha\psi)$ is equivalent to the convexity of $\phi - \alpha\psi$ [7, Thm B.4.3.1].

Finally, note that the equivalences established above still hold if we replace α , “ \geq ”, and “ \succeq ” by 1, “ $>$ ”, and “ \succ ”, respectively. The strict convexity of $\phi - \psi$ then follows from the positive definiteness of $\mathbf{H}(\phi - \psi)$ (ibid.). \square

The following result due to [5] will be crucial to prove the convexity of the superprediction set (Theorem 29).

Lemma 4 ([5]). *Let $\Delta(\Omega)$ be the set of distributions over some set $\Omega \subseteq \mathbb{R}$. Let a function $Q : \Delta(\Omega) \times \Omega \rightarrow \mathbb{R}$ be such that $Q(\cdot, \omega)$ is continuous for all $\omega \in \Omega$. If for all $\boldsymbol{\pi} \in \Delta(\Omega)$ it holds that $\mathbb{E}_{\omega \sim \boldsymbol{\pi}} Q(\boldsymbol{\pi}, \omega) \leq r$, where $r \in \mathbb{R}$ is some constant, then*

$$\exists \boldsymbol{\pi} \in \Delta(\Omega), \forall \omega \in \Omega, Q(\boldsymbol{\pi}, \omega) \leq r.$$

Note that when Ω in the lemma above is $[n]$, $\Delta([n]) \equiv \Delta_n$.

The next crucial lemma is a slight modification of a result due to [5].

Lemma 5. *Let $f : \text{ri } \Delta_n \times [n] \rightarrow \mathbb{R}$ be a continuous function in the first argument and such that $\forall (\mathbf{q}, x) \in \text{ri } \Delta_n \times [n], -\infty < f(\mathbf{q}, x)$. Suppose that $\forall \mathbf{p} \in \text{ri } \Delta_n, \mathbb{E}_{x \sim \mathbf{p}}[f(\mathbf{p}, x)] \leq 0$, then*

$$\forall \epsilon > 0, \exists \mathbf{p}_\epsilon \in \text{ri } \Delta_n, \forall x \in [n], f(\mathbf{p}_\epsilon, x) \leq \epsilon.$$

Proof. Pick any $\delta > 0$ such that $\delta(n-1) < 1$, and $c_0 < 0$ such that $\forall (\mathbf{q}, x) \in \text{ri } \Delta_n \times [n], c_0 \leq f(\mathbf{q}, x)$. We define $\Delta_n^\delta := \{\mathbf{p} \in \Delta_n : \forall x \in [n], p_x \geq \delta\}$ and $g(\mathbf{q}, \mathbf{p}) := \mathbb{E}_{x \sim \mathbf{q}}[f(\mathbf{p}, x)]$. For a fixed $\mathbf{q}, \mathbf{p} \mapsto g(\mathbf{q}, \mathbf{p})$ is continuous, since f is continuous in the first argument. For a fixed $\mathbf{p}, \mathbf{q} \mapsto g(\mathbf{q}, \mathbf{p})$ is linear, and thus concave. Since Δ_n^δ is convex and compact, g satisfies Ky Fan’s minimax Theorem [1, Thm. 11.4], and therefore, there exists $\mathbf{p}^\delta \in \Delta_n^\delta$ such that

$$\forall \mathbf{q} \in \Delta_n^\delta, \mathbb{E}_{x \sim \mathbf{q}}[f(\mathbf{p}^\delta, x)] = g(\mathbf{q}, \mathbf{p}^\delta) \leq \sup_{\boldsymbol{\mu} \in \Delta_n^\delta} g(\boldsymbol{\mu}, \boldsymbol{\mu}) = \sup_{\boldsymbol{\mu} \in \Delta_n^\delta} \mathbb{E}_{x \sim \boldsymbol{\mu}}[f(\boldsymbol{\mu}, x)] \leq 0. \quad (4)$$

For $x_0 \in [n]$, let $\hat{\mathbf{q}} \in \Delta_n^\delta$ be such that $\hat{q}_{x_0} = 1 - \delta(n-1)$ and $\hat{q}_x = \delta$ for $x \neq x_0$ (this is a legitimate distribution since $\delta(n-1) < 1$ by construction). Substituting $\hat{\mathbf{q}}$ for \mathbf{q} in (4) gives

$$\begin{aligned}
& (1 - \delta(n-1))f(\mathbf{p}^\delta, x_0) + \delta \sum_{x \neq x_0} f(\mathbf{p}^\delta, x) && \leq 0, \\
\implies & (1 - \delta(n-1))f(\mathbf{p}^\delta, x_0) && \leq -c_0\delta(n-1), \\
\implies & f(\mathbf{p}^\delta, x_0) && \leq [-c_0\delta(n-1)]/[1 - \delta(n-1)].
\end{aligned}$$

Choosing $\delta^* := \epsilon/[-c_0 + \epsilon](n-1)$, and $\mathbf{p}_\epsilon := \mathbf{p}^{\delta^*}$ gives the desired result. \square

Lemma 6. *Let $f, g : I \rightarrow \mathbb{R}^n$, where $I \subseteq \mathbb{R}$ is an open interval containing 0. Suppose g [resp. f] is continuous [resp. differentiable] at 0. Then $t \mapsto \langle f(t), g(t) \rangle$ is differentiable at 0 if and only if $t \mapsto \langle f(0), g(t) \rangle$ is differentiable at 0, and we have*

$$\left. \frac{d}{dt} \langle f(t), g(t) \rangle \right|_{t=0} = \left\langle \left. \frac{d}{dt} f(t) \right|_{t=0}, g(0) \right\rangle + \left. \frac{d}{dt} \langle f(0), g(t) \rangle \right|_{t=0}.$$

Proof. We have

$$\begin{aligned} \frac{\langle f(t), g(t) \rangle - \langle f(0), g(0) \rangle}{t} &= \frac{\langle f(t), g(t) \rangle - \langle f(0), g(t) \rangle}{t} + \frac{\langle f(0), g(t) \rangle - \langle f(0), g(0) \rangle}{t}, \\ &= \left\langle \frac{f(t) - f(0)}{t}, g(t) \right\rangle + \frac{\langle f(0), g(t) \rangle - \langle f(0), g(0) \rangle}{t}. \end{aligned}$$

But since g [resp. f] is continuous [resp. differentiable] at 0, the first term on the right hand side of the above equation converges to $\langle \frac{d}{dt}f(t)|_{t=0}, g(0) \rangle$ as $t \rightarrow 0$. Therefore, $\frac{1}{t}(\langle f(0), g(t) \rangle - \langle f(0), g(0) \rangle)$ admits a limit when $t \rightarrow 0$ if and only if $\frac{1}{t}(\langle f(t), g(t) \rangle - \langle f(0), g(0) \rangle)$ admits a limit when $t \rightarrow 0$. This shows that $t \mapsto \langle f(0), g(t) \rangle$ is differentiable at 0 if and only if $t \mapsto \langle f(t), g(t) \rangle$ is differentiable at 0, and in this case the above equation yields

$$\begin{aligned} \frac{d}{dt} \langle f(t), g(t) \rangle \Big|_{t=0} &= \lim_{t \rightarrow 0} \frac{\langle f(t), g(t) \rangle - \langle f(0), g(0) \rangle}{t}, \\ &= \lim_{t \rightarrow 0} \left(\left\langle \frac{f(t) - f(0)}{t}, g(t) \right\rangle + \frac{\langle f(0), g(t) \rangle - \langle f(0), g(0) \rangle}{t} \right), \\ &= \left\langle \frac{d}{dt}f(t) \Big|_{t=0}, g(0) \right\rangle + \frac{d}{dt} \langle f(0), g(t) \rangle \Big|_{t=0}. \end{aligned}$$

□

Note that the differentiability of $t \mapsto \langle f(0), g(t) \rangle$ at 0 does not necessarily imply the differentiability of g at 0. Take for example $n = 3$, $f(t) = \mathbf{1}/3$ for $t \in]-1, 1[$, and

$$g(t) = \begin{cases} -t\mathbf{e}_1 + t\frac{\mathbf{1}}{3}, & \text{if } t \in]-1, 0[; \\ -t\frac{\mathbf{1}}{3} + t\mathbf{e}_2, & \text{if } t \in [0, 1[. \end{cases}$$

Thus, the function $t \mapsto \langle f(0), g(t) \rangle = 0$ is differentiable at 0 but g is not. The preceding Lemma will be particularly useful in settings where it is desired to compute the derivative $\frac{d}{dt} \langle f(0), g(t) \rangle|_{t=0}$ without any explicit assumptions on the differentiability of $g(t)$ at 0. For example, this will come up when computing $\frac{d}{dt} \langle \mathbf{p}, D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \rangle|_{t=0}$, where $\mathbf{v} \in \mathbb{R}^{n-1}$ and $t \mapsto \tilde{\alpha}^t$ is smooth curve on $\text{int } \tilde{\Delta}_n$, with the only assumption that \tilde{L}_ℓ is twice differentiable at $\tilde{\alpha}^0 \in \text{int } \tilde{\Delta}_n$.

Lemma 7. *Let $\ell: \Delta_n \rightarrow [0, +\infty]^n$ be a proper loss. For any $\mathbf{p} \in \text{ri } \Delta_n$, it holds that*

$$\ell \text{ is continuous at } \mathbf{p} \xLeftrightarrow{(i)} \underline{L}_\ell \text{ is differentiable at } \mathbf{p} \xLeftrightarrow{(ii)} \partial[-\underline{L}_\ell](\mathbf{p}) = \{\nabla \underline{L}_\ell(\mathbf{p})\} = \{\ell(\mathbf{p})\}.$$

$\xLeftrightarrow{(i)}$. This equivalence has been shown before by [16].

[$\xLeftrightarrow{(ii)}$] Since $\underline{L}_\ell(\mathbf{p}) = -\sigma_{\mathcal{S}_\ell}(-\mathbf{p})$, for all $\mathbf{p} \in \text{ri } \Delta_n$, it follows that \underline{L}_ℓ is differentiable at \mathbf{p} if and only if $\partial[-\underline{L}_\ell](\mathbf{p}) = \partial\sigma_{\mathcal{S}_\ell}(-\mathbf{p}) = \{-\nabla\sigma_{\mathcal{S}_\ell}(-\mathbf{p})\} = \{\nabla \underline{L}_\ell(\mathbf{p})\}$ [7, Cor. D.2.1.4]. It remains to show that $\nabla \underline{L}_\ell(\mathbf{r}) = \ell(\mathbf{r})$ when \underline{L}_ℓ is differentiable at $\mathbf{r} \in \text{ri } \Delta_n$. Let $\alpha_x^t = \mathbf{r} + t\mathbf{e}_x$ and $\tilde{\alpha}_x^t = \Pi_n(\alpha_x^t)$, where $(\mathbf{e}_x)_{x \in [n]}$ is the standard basis of \mathbb{R}^n . For $x \in [n]$, the functions $f_x(t) := \alpha_x^t$ and $g_x(t) := \tilde{\ell}(\tilde{\alpha}_x^t)$ satisfy the conditions of Lemma 6. Therefore, $h_x(t) := \langle f_x(0), g_x(t) \rangle = \langle \mathbf{r}, \tilde{\ell}(\tilde{\alpha}_x^t) \rangle$ is differentiable at 0 and

$$\begin{aligned} \nabla \tilde{L}(\mathbf{r})\mathbf{e}_x &= \frac{d}{dt} \tilde{L}(\alpha_x^t) \Big|_{t=0} = \frac{d}{dt} \langle f_x(t), g_x(t) \rangle \Big|_{t=0}, \\ &= \left\langle \mathbf{e}_x, \tilde{\ell}(\tilde{\mathbf{r}}) \right\rangle + \frac{d}{dt} h_x(t) \Big|_{t=0}, \\ &= \tilde{\ell}_x(\tilde{\mathbf{r}}), \end{aligned}$$

where the last equality holds because h_x attains a minimum at 0 due to the properness of ℓ . The result being true for all $x \in [n]$ implies that $\nabla \tilde{L}(\tilde{\mathbf{r}}) = \tilde{\ell}(\tilde{\mathbf{r}}) = \ell(\mathbf{r})$. □

The next Lemma is a restatement of earlier results due to [14]. Our proof is more concise due to our definition of the Bayes risk in terms of the support function of the superprediction set.

Lemma 8 ([14]). *Let $\ell: \Delta_n \rightarrow [0, +\infty]^n$ be a proper loss whose Bayes risk is twice differentiable on $]0, +\infty[^n$ and let $X_{\mathbf{p}} = I_{\tilde{n}} - \mathbf{1}_{\tilde{n}}\tilde{\mathbf{p}}^\top$. The following holds*

- (i) $\forall \mathbf{p} \in \text{ri } \Delta_n, \langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_{\tilde{n}}^\top$.
- (ii) $\forall \tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n, D\tilde{\ell}(\tilde{\mathbf{p}}) = \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H\tilde{L}_\ell(\tilde{\mathbf{p}})$.
- (iii) $\forall \tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n, H\tilde{L}_{\log}(\tilde{\mathbf{p}}) = -(X_{\tilde{\mathbf{p}}})^{-1}(\text{diag}(\tilde{\mathbf{p}}))^{-1}$.

We show (i) and (ii). Let $\mathbf{p} \in \text{ri } \Delta_n$ and $f(\tilde{\mathbf{q}}) := \langle \mathbf{p}, \tilde{\ell}(\tilde{\mathbf{q}}) \rangle = \langle \mathbf{p}, \nabla L_\ell(\mathbf{q}) \rangle$, where the equality is due to Lemma 7. Since L_ℓ is twice differentiable $]0, +\infty[^n$, f is differentiable on $\text{int } \tilde{\Delta}_n$ and we have $Df(\tilde{\mathbf{q}}) = \langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle$. Since ℓ is proper, f reaches a minimum at $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$, and thus $\langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_{\tilde{n}}^\top$ (this shows (i)). On the other hand, we have $\nabla \tilde{L}_\ell(\tilde{\mathbf{p}}) = J_n^\top \nabla L_\ell(\mathbf{p}) = J_n^\top \tilde{\ell}(\tilde{\mathbf{p}})$. By differentiating and using the chain rule, we get $H\tilde{L}_\ell(\tilde{\mathbf{p}}) = [D\tilde{\ell}(\tilde{\mathbf{p}})]^\top J_n$. This means that $\forall i \in [\tilde{n}]$, $[H\tilde{L}_\ell(\tilde{\mathbf{p}})]_{\cdot, i} = \nabla \tilde{\ell}_i(\tilde{\mathbf{p}}) - \nabla \tilde{\ell}_n(\tilde{\mathbf{p}})$, and thus $\sum_{i=1}^{\tilde{n}} p_i [H\tilde{L}_\ell(\tilde{\mathbf{p}})]_{\cdot, i} = \sum_{i=1}^{\tilde{n}} p_i \nabla \tilde{\ell}_i(\tilde{\mathbf{p}}) - (1 - p_n) \nabla \tilde{\ell}_n(\tilde{\mathbf{p}})$. On the other hand, it follows from point (i) of the lemma that $\sum_{i=1}^{\tilde{n}} p_i \nabla \tilde{\ell}_i(\tilde{\mathbf{p}}) = \mathbf{0}_{\tilde{n}}$. Therefore, $[H\tilde{L}_\ell(\tilde{\mathbf{p}})]\tilde{\mathbf{p}} = -\nabla \tilde{\ell}_n(\tilde{\mathbf{p}})$ and, as a result, $\forall i \in [\tilde{n}]$, $[H\tilde{L}_\ell(\tilde{\mathbf{p}})]_{\cdot, i} - [H\tilde{L}_\ell(\tilde{\mathbf{p}})]\tilde{\mathbf{p}} = \nabla \tilde{\ell}_i(\tilde{\mathbf{p}})$. The last two equations can be combined as $D\tilde{\ell}(\tilde{\mathbf{p}}) = \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H\tilde{L}_\ell(\tilde{\mathbf{p}})$.

[We show (iii)] It follows from (ii), since $\forall i \in [\tilde{n}]$, $\nabla [\tilde{\ell}_{\log}]_i(\tilde{\mathbf{p}}) = \frac{1}{p_i} \mathbf{e}_i$, for $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$.

□

In the next lemma we state a new result for proper losses which will be crucial to prove a necessary condition for Φ -mixability (Theorem 14) — one of the main results of the paper.

Lemma 9. *Let $\ell: \Delta_n \rightarrow [0, +\infty]^n$ be a proper loss whose Bayes risk is twice differentiable on $]0, +\infty[^n$. For $\mathbf{v} \in \mathbb{R}^{n-1}$ and $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$,*

$$\left\langle \mathbf{p}, (D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}) \odot (D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}) \right\rangle = -\mathbf{v}^\top H\tilde{L}_\ell(\tilde{\mathbf{p}})[H\tilde{L}_{\log}(\tilde{\mathbf{p}})]^{-1}H\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{v}, \quad (5)$$

where $\mathbf{p} = \Pi_n(\tilde{\mathbf{p}})$ and L_{\log} is the Bayes risk of the log loss.

Furthermore, if $t \mapsto \tilde{\alpha}^t$ is a smooth curve in $\text{int } \tilde{\Delta}_n$ and satisfies $\tilde{\alpha}^0 = \tilde{\mathbf{p}}$ and $\frac{d}{dt}\tilde{\alpha}^t|_{t=0} = \mathbf{v}$, then $t \mapsto \langle \mathbf{p}, D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \rangle$ is differentiable at 0 and we have

$$\frac{d}{dt} \left\langle \mathbf{p}, D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \right\rangle \Big|_{t=0} = -\mathbf{v}^\top H\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{v}. \quad (6)$$

Proof. We know from Lemma 8 that for $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$, we have $D\tilde{\ell}(\tilde{\mathbf{p}}) = \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H\tilde{L}_\ell(\tilde{\mathbf{p}})$, where $X_{\tilde{\mathbf{p}}} = I_{n-1} - \mathbf{1}_{n-1}\tilde{\mathbf{p}}^\top$. Thus, we can write

$$\begin{aligned} \left\langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \odot D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle &= \mathbf{v}^\top (D\tilde{\ell}(\tilde{\mathbf{p}}))^\top \text{diag}(\mathbf{p}) D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}, \\ &= \mathbf{v}^\top (H\tilde{L}_\ell(\tilde{\mathbf{p}}))^\top [X_{\tilde{\mathbf{p}}}^\top, -\tilde{\mathbf{p}}] \text{diag}(\mathbf{p}) \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{v}. \end{aligned} \quad (7)$$

Observe that $[X_{\tilde{\mathbf{p}}}^\top, -\tilde{\mathbf{p}}] \text{diag}(\mathbf{p}) = [I_{n-1} - \tilde{\mathbf{p}}\mathbf{1}_{n-1}^\top, -\tilde{\mathbf{p}}] \text{diag}(\mathbf{p}) = [\text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top, -\tilde{\mathbf{p}}p_n]$. Thus,

$$\begin{aligned} [X_{\tilde{\mathbf{p}}}^\top, -\tilde{\mathbf{p}}] \text{diag}(\mathbf{p}) \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} &= [\text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top, -\tilde{\mathbf{p}}p_n] \begin{bmatrix} I_{n-1} - \mathbf{1}_{n-1}\tilde{\mathbf{p}}^\top \\ -\tilde{\mathbf{p}}^\top \end{bmatrix}, \\ &= \text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top + \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top(1 - p_n) + p_n\tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top, \\ &= \text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^\top, \\ &= \text{diag}(\tilde{\mathbf{p}}) X_{\tilde{\mathbf{p}}}, \\ &= -(H\tilde{L}_{\log}(\tilde{\mathbf{p}}))^{-1}, \end{aligned} \quad (8)$$

where the last equality is due to Lemma 8. The desired result follows by combining (7) and (8).

[We show (6)] Let $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$, we define $\tilde{\alpha}^t := \tilde{\mathbf{p}} + t\mathbf{v}$, $\alpha^t := \Pi_n(\tilde{\alpha}^t) = \mathbf{p} + tJ_n\mathbf{v}$, and $r(t) := \alpha^t / \|\alpha^t\|$, where $t \in \{s : \tilde{\mathbf{p}} + s\mathbf{v} \in \text{int } \tilde{\Delta}_n\}$. Since $t \mapsto r(t)$ is differentiable at 0 and $t \mapsto D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v}$ is continuous at 0, it follows from Lemma 3 that

$$\begin{aligned} \frac{d}{dt} \left\langle r(0), D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \right\rangle \Big|_{t=0} &= \frac{d}{dt} \left\langle r(t), D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \right\rangle \Big|_{t=0} - \left\langle \dot{r}(0), D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle, \\ &= - \left\langle \dot{r}(0), D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle, \end{aligned}$$

where the second equality holds since, according to Lemma 8, we have $\langle \alpha^t, D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \rangle = 0$. Since $r(0) = \mathbf{p} / \|\mathbf{p}\|$, $\dot{r}(0) = \|\mathbf{p}\|^{-1} (I_n - r(0)[r(0)]^\top) J_n \mathbf{v}$, and $J_n = \begin{bmatrix} I_{n-1} \\ -\mathbf{1}_{n-1}^\top \end{bmatrix}$, we get

$$\begin{aligned} \|\tilde{\mathbf{p}}\| \frac{d}{dt} \left\langle r(0), D\tilde{\ell}(\tilde{\alpha}^t)\mathbf{v} \right\rangle \Big|_{t=0} &= - \left\langle (I_n - r(0)[r(0)]^\top) J_n \mathbf{v}, D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle, \\ &= - \left\langle J_n \mathbf{v}, D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle, \\ &= - \left\langle J_n \mathbf{v}, \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v} \right\rangle, \\ &= -\mathbf{v}^\top H_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v}, \end{aligned} \tag{9}$$

where the passage to (9) is due to $r(0) = \mathbf{p} / \|\mathbf{p}\| \perp D\tilde{\ell}(\tilde{\mathbf{p}})$. In the last equality we used the fact that $J_n^\top \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} = [I_{n-1}, -\mathbf{1}_{n-1}] \begin{bmatrix} I_{n-1} - \mathbf{1}_{n-1}\tilde{\mathbf{p}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} = I_{n-1}$. \square

Proposition 10. Let $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy and $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ a closed admissible loss. If ℓ is Φ -mixable, then $\forall \mathfrak{l} \subseteq [k]$ with $|\mathfrak{l}| > 1$, ℓ is $\Phi_{\mathfrak{l}}$ -mixable and

$$\forall \mathbf{q} \in \text{rbd } \Delta_{\mathfrak{l}}, \forall \hat{\mathbf{q}} \in \text{ri } \Delta_{\mathfrak{l}}, \Phi'(\mathbf{q}; \hat{\mathbf{q}} - \mathbf{q}) = -\infty. \tag{10}$$

Given an entropy $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ and a loss $\ell: \mathcal{A} \rightarrow [0, +\infty]$, we define

$$\mathfrak{m}_{\Phi}(x, A, \mathbf{a}, \hat{\mathbf{q}}, \boldsymbol{\mu}) := \langle \boldsymbol{\mu}, \ell_x(A) \rangle + D_{\Phi}(\boldsymbol{\mu}, \hat{\mathbf{q}}) - \ell_x(\mathbf{a}),$$

where $x \in [n]$, $A \in \mathcal{A}^k$, $\mathbf{a} \in \mathcal{A}$, and $\mathbf{q}, \hat{\mathbf{q}} \in \Delta_k$. Reid et al. [9] showed that ℓ is Φ mixable if and only if

$$\widehat{\mathfrak{m}}_{\Phi} := \inf_{A \in \mathcal{A}^k} \sup_{\hat{\mathbf{q}} \in \Delta_k} \inf_{\mathbf{a}_* \in \mathcal{A}} \sup_{\boldsymbol{\mu} \in \Delta_k, x \in [n]} \mathfrak{m}_{\Phi}(x, A, \mathbf{a}_*, \hat{\mathbf{q}}, \boldsymbol{\mu}) \geq 0.$$

Proof of Proposition 10. [We show that ℓ is $\Phi_{\mathfrak{l}}$ -mixable] Let $\mathfrak{l} \subseteq [k]$, with $|\mathfrak{l}| > 1$, $A \in \mathcal{A}^k$, and $\mathbf{q} \in \Delta_{\mathfrak{l}}$. Since ℓ is Φ -mixable, the following holds

$$\exists \mathbf{a}_* \in \Delta_n, \forall x \in [n], \ell_x(\mathbf{a}_*) \leq \inf_{\hat{\mathbf{q}} \in \Delta_k} \langle \hat{\mathbf{q}}, \ell_x(A) \rangle + D_{\Phi}(\hat{\mathbf{q}}, \mathbf{q}), \tag{11}$$

$$\leq \inf_{\hat{\mathbf{q}} \in \Delta_{\mathfrak{l}}} \langle \hat{\mathbf{q}}, \ell_x(A) \rangle + D_{\Phi}(\hat{\mathbf{q}}, \mathbf{q}), \tag{12}$$

$$\begin{aligned} &= \inf_{\hat{\mathbf{q}} \in \Delta_{\mathfrak{l}}} \langle \Pi_{\mathfrak{l}} \hat{\mathbf{q}}, \Pi_{\mathfrak{l}} \ell_x(A) \rangle + D_{\Phi_{\mathfrak{l}}}(\Pi_{\mathfrak{l}} \hat{\mathbf{q}}, \Pi_{\mathfrak{l}} \mathbf{q}), \\ &= \inf_{\hat{\boldsymbol{\mu}} \in \Delta_{|\mathfrak{l}|}} \langle \hat{\boldsymbol{\mu}}, \ell_x(A \Pi_{\mathfrak{l}}^\top) \rangle + D_{\Phi_{\mathfrak{l}}}(\hat{\boldsymbol{\mu}}, \Pi_{\mathfrak{l}} \mathbf{q}), \end{aligned} \tag{13}$$

where in (11) we used the fact that $\Phi_{\mathfrak{l}}(\Pi_{\mathfrak{l}} \mathbf{q}) = \Phi(\mathbf{q})$, $\forall \mathbf{q} \in \Delta_{\mathfrak{l}}$. Given that $A \mapsto A \Pi_{\mathfrak{l}}^\top$ [resp. $\mathbf{q} \mapsto \Pi_{\mathfrak{l}} \mathbf{q}$] is onto from \mathcal{A}^k to $\mathcal{A}^{|\mathfrak{l}|}$ [resp. from $\Delta_{\mathfrak{l}}$ to $\Delta_{|\mathfrak{l}|}$], (13) implies that ℓ is $\Phi_{\mathfrak{l}}$ -mixable.

[We show (10)] Suppose that there exists $\hat{\mathbf{q}} \in \text{rbd } \Delta_k$ and $\mathbf{q} \in \text{ri } \Delta_k$ such that $|\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}})| < +\infty$. Let $f: [0, \epsilon] \rightarrow \mathbb{R}$ be defined by $f(\lambda) := \Phi(\hat{\mathbf{q}} + \lambda(\mathbf{q} - \hat{\mathbf{q}}))$, where $\epsilon > 0$ is such that $\hat{\mathbf{q}} + \epsilon(\mathbf{q} - \hat{\mathbf{q}}) \in \text{ri } \Delta_k$. The function f is closed and convex on $\text{dom } f = [0, \epsilon]$ and $\lim_{\lambda \downarrow 0} \frac{f(\lambda) - f(0)}{\lambda} = f'(0; 1) = \Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}})$ which is finite by assumption. Using this and the fact that $\lambda f'(0; 1) = f'(0; \lambda)$, we have $\lim_{\lambda \downarrow 0} \lambda^{-1}(f(\lambda) - f(0) - f'(0; \lambda)) = 0$. Substituting f by its expression in terms of Φ in the latter equality gives

$$\lim_{\lambda \downarrow 0} \lambda^{-1} D_{\Phi}(\hat{\mathbf{q}} + \lambda(\mathbf{q} - \hat{\mathbf{q}}), \hat{\mathbf{q}}) = 0. \tag{14}$$

Let $\eta > 0$ and $\theta^* \in [k]$ be such that $\hat{q}_{\theta^*} = 0$. Suppose that ℓ is an admissible, Φ -mixable loss. The fact that ℓ is admissible implies that there exists $(x_0, x_1, \mathbf{a}_0, \mathbf{a}_1) \in [n] \times [n] \times \mathcal{A} \times \mathcal{A}$ such that [9]

$$\mathbf{a}_1 \in \operatorname{argmin}\{\ell_{x_0}(\mathbf{a}) : \ell_{x_1}(\mathbf{a}) = \inf_{\hat{\mathbf{a}} \in \mathcal{A}} \ell_{x_1}(\hat{\mathbf{a}})\} \text{ and } \inf_{\hat{\mathbf{a}} \in \mathcal{A}} \ell_{x_0}(\mathbf{a}) = \ell_{x_0}(\mathbf{a}_0) < \ell_{x_0}(\mathbf{a}_1). \quad (15)$$

In particular, it holds that $\ell_{x_0}(\mathbf{a}_0) < \ell_{x_0}(\mathbf{a}_1)$. Fix $A \in \mathcal{A}^k$, such that $A_{\cdot, \theta^*} = \mathbf{a}_0$ and $A_{\cdot, \theta} = \mathbf{a}_1$ for $\theta \in [k] \setminus \{\theta^*\}$. Let

$$\mathbf{a}_* := \operatorname{argmax}_{\mathbf{a} \in \Delta_n} \inf_{\mu \in \Delta_k, x \in [n]} m_\Phi(x, A, \mathbf{a}, \hat{\mathbf{q}}, \mu),$$

with $\hat{\mathbf{q}} \in \operatorname{rbd} \Delta_k$ as in (14). Note that \mathbf{a}_* exists since ℓ is closed.

If \mathbf{a}_* is such that $\ell_{x_1}(\mathbf{a}_*) > \ell_{x_1}(\mathbf{a}_1)$, then taking $\mu = \hat{\mathbf{q}}$ puts all weights on experts predicting \mathbf{a}_1 , while $D_\Phi(\mu, \hat{\mathbf{q}}) = 0$. Therefore,

$$\widehat{m}_\Phi \leq \inf_{\mu \in \Delta_k, x \in [n]} m_\Phi(x, A, \mathbf{a}_*, \hat{\mathbf{q}}, \mu) \leq m_\Phi(x_1, A, \mathbf{a}, \hat{\mathbf{q}}, \hat{\mathbf{q}}) < 0.$$

This contradicts the Φ -mixability of ℓ . Therefore, $\ell_{x_1}(\mathbf{a}_*) = \ell_{x_1}(\mathbf{a}_1)$, which by (15) implies $\ell_{x_0}(\mathbf{a}_*) \geq \ell_{x_0}(\mathbf{a}_1)$. For $\mathbf{q}^\lambda = \hat{\mathbf{q}} + \lambda(\mathbf{q} - \hat{\mathbf{q}})$, with $\mathbf{q} \in \operatorname{ri} \Delta_k$ as in (11) and $\lambda \in [0, \epsilon]$,

$$\begin{aligned} \widehat{m}_\Phi &\leq \inf_{\mu \in \Delta_k, x \in [n]} m_\Phi(x, A, \mathbf{a}_*, \hat{\mathbf{q}}, \mu), \\ &\leq m_\Phi(x_0, A, \mathbf{a}, \hat{\mathbf{q}}, \mathbf{q}^\lambda), \\ &= \langle \mathbf{q}^\lambda, \ell_{x_0}(A) \rangle + D_\Phi(\mathbf{q}^\lambda, \hat{\mathbf{q}}) - \ell_{x_0}(\mathbf{a}_*), \\ &= (1 - \lambda q_{\theta^*}) \ell_{x_0}(\mathbf{a}_1) + \lambda q_{\theta^*} \ell_{x_0}(\mathbf{a}_0) + D_\Phi(\mathbf{q}^\lambda, \hat{\mathbf{q}}) - \ell_{x_0}(\mathbf{a}_*), \\ &\leq (1 - \lambda q_{\theta^*}) \ell_{x_0}(\mathbf{a}_*) + \lambda q_{\theta^*} \ell_{x_0}(\mathbf{a}_0) + D_\Phi(\mathbf{q}^\lambda, \hat{\mathbf{q}}) - \ell_{x_0}(\mathbf{a}_*), \\ &= \lambda q_{\theta^*} (\ell_{x_0}(\mathbf{a}_0) - \ell_{x_0}(\mathbf{a}_*)) + D_\Phi(\hat{\mathbf{q}} + \lambda(\mathbf{q} - \hat{\mathbf{q}}), \hat{\mathbf{q}}). \end{aligned}$$

Since $q_{\theta^*} > 0$ ($\mathbf{q} \in \operatorname{ri} \Delta_k$) and $\ell_{x_0}(\mathbf{a}_0) < \ell_{x_0}(\mathbf{a}_1) \leq \ell_{x_0}(\mathbf{a}_*)$, (11) implies that there exists $\lambda_* > 0$ small enough such that $\lambda_* q_{\theta^*} (\ell_{x_0}(\mathbf{a}_0) - \ell_{x_0}(\mathbf{a}_*)) + D_\Phi(\hat{\mathbf{q}} + \lambda_*(\mathbf{q} - \hat{\mathbf{q}}), \hat{\mathbf{q}}) < 0$. But this implies that $\widehat{m}_\Phi < 0$ which contradicts the Φ -mixability of ℓ . Therefore, $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}})$ is either equal to $+\infty$ or $-\infty$. The former case is not possible. In fact, since Φ is convex, it must have non-decreasing slopes; in particular, it holds that $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) \leq \Phi(\mathbf{q} - \hat{\mathbf{q}}) - \Phi(\hat{\mathbf{q}})$. Since Φ is finite on Δ_k (by definition of an entropy), we have $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) < +\infty$. Therefore, we have just shown that

$$\forall \hat{\mathbf{q}} \in \operatorname{rbd} \Delta_k, \forall \mathbf{q} \in \operatorname{ri} \Delta_k, \Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) = -\infty. \quad (16)$$

Now suppose that $(\hat{\mathbf{q}}, \mathbf{q}) \in \operatorname{rbd} \Delta_l \times \operatorname{ri} \Delta_l$ for $l \subseteq [k]$, with $|l| > 1$. Note that in this case, we have $(\Phi_l)'(\Pi_l \hat{\mathbf{q}}; \Pi_l(\mathbf{q} - \hat{\mathbf{q}})) = \Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}})$. We showed in the first step of this proof that under the assumptions of the proposition, ℓ must be Φ_l -mixable. Therefore, repeating the steps above that lead to (16) for Φ , $\hat{\mathbf{q}}$, and \mathbf{q} substituted by Φ_l , $\Pi_l \hat{\mathbf{q}} \in \operatorname{rbd} \Delta_{|l|}$, and $\Pi_l \mathbf{q} \in \operatorname{ri} \Delta_{|l|}$, we obtain $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) = \Phi'_l(\Pi_l \hat{\mathbf{q}}; \Pi_l(\mathbf{q} - \hat{\mathbf{q}})) = -\infty$. This shows (10). \square

Lemma 11. For $\eta > 0$, $S_\eta := \eta^{-1} S$ satisfies (10) for all $l \subseteq [k]$ such that $|l| > 1$, where S is the Shannon entropy.

Proof. Let $l \subseteq [k]$ such that $|l| > 1$. Let $(\hat{\mathbf{q}}, \mathbf{q}) \in \operatorname{rbd} \Delta_l \times \operatorname{ri} \Delta_l$ and $\mathbf{q}^\lambda := \hat{\mathbf{q}} + \lambda(\mathbf{q} - \hat{\mathbf{q}})$, for $\lambda \in]0, 1[$. Let $\mathcal{J} := \{j \in l : \hat{q}_j \neq 0\}$ and $\mathcal{K} := l \setminus \mathcal{J}$. We have

$$\begin{aligned} S(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) &= \lim_{\lambda \downarrow 0} \lambda^{-1} \left[\sum_{\theta \in \mathcal{J}} q_\theta^\lambda \log q_\theta^\lambda - \sum_{\theta' \in \mathcal{J}} \hat{q}_{\theta'} \log \hat{q}_{\theta'} \right], \\ &= \lim_{\lambda \downarrow 0} \lambda^{-1} \left[\sum_{\theta \in \mathcal{J}} (q_\theta^\lambda \log q_\theta^\lambda - \hat{q}_\theta \log \hat{q}_\theta) + \sum_{\theta' \in \mathcal{K}} q_{\theta'}^\lambda \log q_{\theta'}^\lambda \right]. \end{aligned} \quad (17)$$

Observe that the limit of either summation term inside the bracket in (17) is equal to zero. Thus, using l'Hopital's rule we get

$$\begin{aligned} S(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) &= \lim_{\lambda \downarrow 0} \left[\sum_{\theta \in \mathcal{J}} [(q_\theta - \hat{q}_\theta) \log q_\theta^\lambda + (q_\theta - \hat{q}_\theta)] + \sum_{\theta' \in \mathcal{K}} [q_{\theta'} \log q_{\theta'}^\lambda + q_{\theta'}] \right], \\ &= \sum_{\theta \in \mathcal{J}} (q_\theta - \hat{q}_\theta) \log \hat{q}_\theta + \sum_{\theta' \in \mathcal{K}} q_{\theta'} \left[\lim_{\lambda \downarrow 0} \log q_{\theta'}^\lambda \right], \end{aligned} \quad (18)$$

where in (18) we used the fact that $\sum_{\theta \in \mathcal{J}} (q_\theta - \hat{q}_\theta) + \sum_{\theta' \in \mathcal{R}} q_{\theta'} = 0$. Since for all $\theta' \in \mathcal{R}$, $\lim_{\lambda \downarrow 0} q_{\theta'}^\lambda = 0$, the right hand side of (6) is equal to $-\infty$. Therefore S satisfies (10). Since $S_\eta = \eta^{-1} S$, it is clear that S_η also satisfies (10). \square

Lemma 12. *Let $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy satisfying (10) for all $\mathfrak{l} \subseteq [k]$ such that $|\mathfrak{l}| > 1$. Then for all such \mathfrak{l} , it holds that*

$$\forall \mathbf{q} \in \Delta_{\mathfrak{l}}, \forall \boldsymbol{\mu} \in \Delta_k \setminus \Delta_{\mathfrak{l}}, D_\Phi(\boldsymbol{\mu}, \mathbf{q}) = +\infty.$$

Proof. Let $\boldsymbol{\mu} \in \Delta_k \setminus \Delta_{\mathfrak{l}}$ and $\mathcal{J} := \{\theta \in [k] : \mu_\theta \neq 0\} \cup \mathfrak{l}$. In this case, we have $\mathbf{q} \in \text{rbd } \Delta_{\mathcal{J}}$ and $\mathbf{q} + 2^{-1}(\boldsymbol{\mu} - \mathbf{q}) \in \text{ri } \Delta_{\mathcal{J}}$. Thus, since Φ satisfies (10) and $\Phi'(\mathbf{q}; \cdot)$ is 1-homogeneous [7, Prop. D.1.1.2], it follows that $2^{-1}\Phi'(\mathbf{q}; \boldsymbol{\mu} - \mathbf{q}) = \Phi'(\mathbf{q}; 2^{-1}(\boldsymbol{\mu} - \mathbf{q})) = -\infty$. Hence $D_\Phi(\boldsymbol{\mu}, \mathbf{q}) = +\infty$. \square

Lemma 13. *Let $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy satisfying (10) for all $\mathfrak{l} \subseteq [k]$ such that $|\mathfrak{l}| > 1$. If Φ satisfies (10), then $\partial\tilde{\Phi}(\tilde{\mathbf{q}}) = \emptyset, \forall \tilde{\mathbf{q}} \in \text{bd } \tilde{\Delta}_k$. Furthermore, $\forall \mathfrak{l} \subseteq [k]$ such that $|\mathfrak{l}| > 1$,*

$$\forall \mathbf{d} \in \mathbb{R}^k, \forall \mathbf{q} \in \text{ri } \Delta_{\mathfrak{l}}, \text{Mix}_\Phi(\mathbf{d}, \mathbf{q}) = \text{Mix}_{\Phi_{\mathfrak{l}}}(\Pi_{\mathfrak{l}}\mathbf{d}, \Pi_{\mathfrak{l}}\mathbf{q}).$$

Proof. Let $\boldsymbol{\mu} \in \text{rbd } \Delta_k$. Since Φ satisfies (10), it follows that $\forall \mathbf{q} \in \text{ri } \Delta_k, \tilde{\Phi}(\tilde{\boldsymbol{\mu}}; \tilde{\mathbf{q}} - \tilde{\boldsymbol{\mu}}) = \Phi'(\boldsymbol{\mu}; \mathbf{q} - \boldsymbol{\mu}) = -\infty$. Therefore, $\partial\tilde{\Phi}'(\tilde{\boldsymbol{\mu}}) = \emptyset$ [11, Thm. 23.4].

Let $\mathbf{d} \in \mathbb{R}^n, \mathfrak{l} \subseteq [k]$, with $|\mathfrak{l}| > 1$, and $\mathbf{q} \in \text{ri } \Delta_{\mathfrak{l}}$. Then

$$\begin{aligned} \text{Mix}_{\Phi_{\mathfrak{l}}}(\Pi_{\mathfrak{l}}\mathbf{d}, \Pi_{\mathfrak{l}}\mathbf{q}) &= \inf_{\boldsymbol{\pi} \in \Delta_{|\mathfrak{l}|}} \langle \boldsymbol{\pi}, \Pi_{\mathfrak{l}}\mathbf{d} \rangle + D_{\Phi_{\mathfrak{l}}}(\boldsymbol{\pi}, \Pi_{\mathfrak{l}}\mathbf{q}), \\ &= \inf_{\boldsymbol{\mu} \in \Delta_{\mathfrak{l}}} \langle \boldsymbol{\mu}, \mathbf{d} \rangle + D_\Phi(\boldsymbol{\mu}, \mathbf{q}), \\ &\leq \inf_{\boldsymbol{\mu} \in \Delta_k} \langle \boldsymbol{\mu}, \mathbf{d} \rangle + D_\Phi(\boldsymbol{\mu}, \mathbf{q}), \\ &= \text{Mix}_\Phi(\mathbf{d}, \mathbf{q}). \end{aligned} \tag{19}$$

To complete the proof, we need to show that (19) holds with equality. For this, it suffices to prove that $\forall \boldsymbol{\mu} \in \Delta_k \setminus \Delta_{\mathfrak{l}}, D_\Phi(\boldsymbol{\mu}, \mathbf{q}) = +\infty$. This follows from Corollary 12. \square

Lemma 14. *Let $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy satisfying (10) for all $\mathfrak{l} \subseteq [k]$ such that $|\mathfrak{l}| > 1$. Let $x \in [n], \mathbf{d} \in \mathbb{R}^k$, and $\mathbf{q} \in \Delta_k$. The infimum in*

$$\text{Mix}_\Phi(\mathbf{d}, \mathbf{q}) = \inf_{\boldsymbol{\mu} \in \Delta_k} \langle \boldsymbol{\mu}, \mathbf{d} \rangle + D_\Phi(\boldsymbol{\mu}, \mathbf{q}) \tag{20}$$

is attained at some $\mathbf{q}_ \in \Delta_k$. Furthermore, if $\mathbf{q} \in \text{ri } \Delta_k$ and \mathbf{q}_* is the infimum of (20) then for any $\mathbf{s}_q^* \in \text{argmax}\{\langle \mathbf{s}, \tilde{\mathbf{q}}_* - \tilde{\mathbf{q}} \rangle : \mathbf{s} \in \partial\tilde{\Phi}(\tilde{\mathbf{q}})\}$, we have*

$$\tilde{\mathbf{q}}_* \in \partial\tilde{\Phi}^*(\mathbf{s}_q^* - J_k^\top \mathbf{d}), \tag{21}$$

$$\text{Mix}_\Phi(\mathbf{d}, \mathbf{q}) = d_k + \tilde{\Phi}^*(\mathbf{s}_q^*) - \tilde{\Phi}^*(\mathbf{s}_q^* - J_k^\top \mathbf{d}). \tag{22}$$

Proof. Let $\mathbf{q} \in \text{ri } \Delta_k$. Since $\tilde{\mathbf{q}} \in \text{int dom } \tilde{\Phi} = \text{int } \tilde{\Delta}_k$, the function $\tilde{\boldsymbol{\mu}} \mapsto -\tilde{\Phi}'(\tilde{\mathbf{q}}; \tilde{\boldsymbol{\mu}} - \tilde{\mathbf{q}})$ is lower semicontinuous [11, Cor. 24.5.1]. Given that $\tilde{\boldsymbol{\mu}} \mapsto \langle \Pi_k(\tilde{\boldsymbol{\mu}}), \mathbf{d} \rangle + \tilde{\Phi}(\tilde{\boldsymbol{\mu}}) - \tilde{\Phi}(\tilde{\mathbf{q}})$ is a closed convex function, it is also lower semi-continuous. Therefore, the function

$$\tilde{\boldsymbol{\mu}} \mapsto \langle \Pi_k(\tilde{\boldsymbol{\mu}}), \mathbf{d} \rangle + \tilde{\Phi}(\tilde{\boldsymbol{\mu}}) - \tilde{\Phi}(\tilde{\mathbf{q}}) - \tilde{\Phi}'(\tilde{\mathbf{q}}; \tilde{\boldsymbol{\mu}} - \tilde{\mathbf{q}})$$

is lower semicontinuous, and thus attains its minimum on the compact set $\tilde{\Delta}_k$ at some point $\tilde{\mathbf{q}}_*$. Using the fact that $D_\Phi(\boldsymbol{\mu}, \mathbf{q}) = D_{\tilde{\Phi}}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{q}})$, we get that

$$\mathbf{q}_* := \Pi_k(\tilde{\mathbf{q}}_*) = \underset{\boldsymbol{\mu} \in \Delta_k}{\text{argmin}} \langle \boldsymbol{\mu}, \mathbf{d} \rangle + D_\Phi(\boldsymbol{\mu}, \mathbf{q}). \tag{23}$$

If $\mathbf{q} \in \text{rbd } \Delta_k$, then either \mathbf{q} is a vertex of Δ_k or there exists $\mathfrak{l} \subsetneq [k]$ such that $\mathbf{q} \in \text{ri } \Delta_{\mathfrak{l}}$. In the former case, it follows from (10) that $D_\Phi(\boldsymbol{\mu}, \mathbf{q}) = +\infty$ for all $\boldsymbol{\mu} \in \Delta_k \setminus \{\mathbf{q}\}$, and thus the infimum

of (20) is trivially attained at $\mu = q$. Now consider the alternative — $q \in \Delta_l$ with $|l| > 1$. Using Corollary 12, we have $D_\Phi(\mu, q) = +\infty$ for all $\mu \in \Delta_k \setminus \Delta_l$. Therefore,

$$\begin{aligned} \text{Mix}_\Phi(d, q) &= \inf_{\mu \in \Delta_l} \langle \mu, d \rangle + D_\Phi(\mu, q), \\ &= \inf_{\hat{\mu} \in \Delta_{|l|}} \langle \hat{\mu}, \Pi_l d \rangle + D_{\Phi_l}(\hat{\mu}, \Pi_l q), \end{aligned} \quad (24)$$

where $\Phi_l := \Phi \circ \Pi_l$. Since $\Pi_l q \in \text{ri } \Delta_{|l|}$, we can use the same argument as the previous paragraph with Φ and q replaced by Φ_l and $\Pi_l q$, respectively, to show that the infimum in (24) is attained at some $\hat{q}_* \in \Delta_{|l|}$. Thus, $q_* := \Pi_l^\top \hat{q}_* \in \Delta_k$ attains the infimum in (20).

Now we show the second part of the lemma. Let $q \in \text{ri } \Delta_k$ and q_* be the infimum of (20). Since $\tilde{\Phi}$ is convex and $\tilde{q} = \Pi_k(q) \in \text{int } \tilde{\Delta}_k = \text{int dom } \tilde{\Phi}$, we have $\partial \tilde{\Phi}(\tilde{q}) \neq \emptyset$ [11, Thm. 23.4]. This means that there exists $s_q^* \in \partial \tilde{\Phi}(\tilde{q})$ such that $\langle s_q^*, \tilde{q}_* - \tilde{q} \rangle = \tilde{\Phi}'(\tilde{q}; \tilde{q}_* - \tilde{q})$ [7, p.166]. We will now show that $s_q^* - J_k^\top d \in \partial \tilde{\Phi}(\tilde{q}_*)$, which will imply that $\tilde{q}_* \in \partial \tilde{\Phi}^*(s_q^* - J_k^\top d)$ (ibid., Cor. D.1.4.4). Let $q_* = \text{argmin}_{\mu \in \Delta_k} \langle \mu, d \rangle + D_\Phi(\mu, q)$. Thus, for all $\mu \in \Delta_k$,

$$\begin{aligned} &\langle \mu, d \rangle + \tilde{\Phi}(\tilde{\mu}) - \tilde{\Phi}(\tilde{q}) - \tilde{\Phi}'(\tilde{q}; \tilde{\mu} - \tilde{q}) \geq \langle q_*, d \rangle + \tilde{\Phi}(\tilde{q}_*) - \tilde{\Phi}(\tilde{q}) - \langle s_q^*, \tilde{q}_* - \tilde{q} \rangle, \\ \implies &\tilde{\Phi}(\tilde{\mu}) \geq \tilde{\Phi}(\tilde{q}_*) - \langle \tilde{\mu} - \tilde{q}_*, J_k^\top d \rangle + \langle s_q^*, \tilde{q} - \tilde{q}_* \rangle + \tilde{\Phi}'(\tilde{q}; \tilde{\mu} - \tilde{q}), \\ \implies &\tilde{\Phi}(\tilde{\mu}) \geq \tilde{\Phi}(\tilde{q}_*) - \langle \tilde{\mu} - \tilde{q}_*, J_k^\top d \rangle + \langle s_q^*, \tilde{q} - \tilde{q}_* \rangle + \langle s_q^*, \tilde{\mu} - \tilde{q} \rangle, \\ \implies &\tilde{\Phi}(\tilde{\mu}) \geq \tilde{\Phi}(\tilde{q}_*) + \langle \tilde{\mu} - \tilde{q}_*, s_q^* - J_k^\top d \rangle, \end{aligned}$$

where in the second line we used the fact that $\forall q \in \Delta_k, \langle q, d \rangle = \langle \tilde{q}, J_k^\top d \rangle + d_k$, and in third line we used the fact that $\forall s \in \partial \tilde{\Phi}(\tilde{q}), \langle s, \tilde{\mu} - \tilde{q} \rangle \leq \tilde{\Phi}'(\tilde{q}; \tilde{\mu} - \tilde{q})$ (ibid.). This shows that $s_q^* - J_k^\top d \in \partial \tilde{\Phi}(\tilde{q}_*)$.

Substituting $\tilde{\Phi}'(\tilde{q}; \tilde{q}_* - \tilde{q})$ by $\langle s_q^*, q_* - q \rangle$ in the expression for $\text{Mix}_\Phi(d, q)$, we get

$$\begin{aligned} \text{Mix}_\Phi(d, q) &= d_k + \langle \tilde{q}_*, J_k^\top d \rangle + \tilde{\Phi}(\tilde{q}_*) - \tilde{\Phi}(\tilde{q}) - \langle s_q^*, \tilde{q}_* - \tilde{q} \rangle, \\ &= d_k + \langle s_q^*, \tilde{q} \rangle - \tilde{\Phi}(\tilde{q}) - [\langle s_q^* - J_k^\top d, \tilde{q}_* \rangle - \tilde{\Phi}(\tilde{q}_*)], \\ &= d_k + \tilde{\Phi}^*(s_q^*) - \tilde{\Phi}^*(s_q^* - J_k^\top d), \end{aligned}$$

where in the last line we used the fact that $\tilde{\Phi}$ is a closed convex function, and thus $\forall \tilde{q} \in \tilde{\Delta}_k, s \in \partial \tilde{\Phi}(\tilde{q}) \implies \tilde{\Phi}^*(s) = \langle s, \tilde{q} \rangle - \tilde{\Phi}(\tilde{q})$ (ibid., Cor. E.1.4.4).

□

Lemma 15. *Let $q \in \Delta_k$. For any sequence (d_m) in $[0, +\infty]^k$ converging to $d \in [0, +\infty]^k$ coordinate-wise and any entropy $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ satisfying (10) for $l \subseteq [k]$ such that $|l| > 1$,*

$$\lim_{m \rightarrow \infty} \text{Mix}_\Phi(d_m, q) = \text{Mix}_\Phi(d, q). \quad (25)$$

Proof of Lemma 15. Let $q \in \Delta_k$ and $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy as in the statement of the Lemma. Let $(d_m) \subset \mathbb{R}^k$ such that $d_m \xrightarrow{m \rightarrow \infty} d \in \mathbb{R}^k$ in $[0, +\infty]^k$. Let $l := \{\theta \in [k] : d_\theta < +\infty\}$. If $l = \emptyset$ then the result holds trivially since, on the one hand, $\text{Mix}_\Phi(d, q) = +\infty$ and on the other hand $\text{Mix}_\Phi(d_m, q) \geq \min_{\theta \in [k]} d_{m,\theta} \xrightarrow{m \rightarrow \infty} +\infty$.

Assume now that $l \neq \emptyset$. Then

$$\text{Mix}_\Phi(d_m, q) = \inf_{\mu \in \Delta_k} \langle \mu, d_m \rangle + D_\Phi(\mu, q), \quad (26)$$

$$\leq \inf_{\hat{\mu} \in \Delta_l} \langle \hat{\mu}, d \rangle + D_\Phi(\hat{\mu}, q), \quad (27)$$

$$< +\infty, \quad (28)$$

where the last inequality stems from the fact that $\Pi_l d_m$ is a finite vector in $\mathbb{R}^{|l|}$. Therefore, (28) implies that the sequence $\alpha_m := \text{Mix}_\Phi(d_m, q)$ is bounded. We will show that (α_m) converges in \mathbb{R} and that its limit is exactly $\text{Mix}_\Phi(d, q)$. Let $(\hat{\alpha}_m)$ be any convergent subsequence of (α_m) , and let

($\hat{\mathbf{d}}_m$) be the corresponding subsequence of (\mathbf{d}_m). Consider the infimum in (100) with \mathbf{d}_m is replaced by $\hat{\mathbf{d}}_m$. From Lemma 14, this infimum is attained at some $\mathbf{q}_m \in \Delta_k$. Since Δ_k is compact, we may assume without loss of generality that \mathbf{q}_m converges to some $\bar{\mathbf{q}} \in \Delta_k$. Observe that $\bar{\mathbf{q}}$ must be $\Delta_{\mathfrak{I}}$; suppose that $\exists \theta_* \in \bar{\mathfrak{I}}$ such that $\bar{q}_{\theta_*} > 0$. Then

$$\begin{aligned}\hat{\alpha}_m &\geq \langle \mathbf{q}_m, \hat{\mathbf{d}}_m \rangle, \\ &\geq q_{m,\theta_*} \hat{d}_{m,\theta_*} \xrightarrow{m \rightarrow \infty} +\infty.\end{aligned}$$

This would contradict the fact that α_m is bounded, and thus $\bar{\mathbf{q}} \in \Delta_{\mathfrak{I}}$. Using this, we get

$$\begin{aligned}\text{Mix}_{\Phi}(\hat{\mathbf{d}}_m, \mathbf{q}) &= \langle \mathbf{q}_m, \hat{\mathbf{d}}_m \rangle + D_{\Phi}(\mathbf{q}_m, \mathbf{q}), \\ &\geq \langle \Pi_{\mathfrak{I}} \mathbf{q}_m, \Pi_{\mathfrak{I}} \hat{\mathbf{d}}_m \rangle + D_{\Phi}(\mathbf{q}_m, \mathbf{q}), \\ &\xrightarrow{m \rightarrow \infty} \langle \Pi_{\mathfrak{I}} \bar{\mathbf{q}}, \Pi_{\mathfrak{I}} \mathbf{d} \rangle + D_{\Phi}(\bar{\mathbf{q}}, \mathbf{q}), \\ &= \langle \bar{\mathbf{q}}, \mathbf{d} \rangle + D_{\Phi}(\bar{\mathbf{q}}, \mathbf{q}), \tag{29} \\ &\geq \inf_{\hat{\boldsymbol{\mu}} \in \Delta_{\mathfrak{I}}} \langle \hat{\boldsymbol{\mu}}, \mathbf{d} \rangle + D_{\Phi}(\hat{\boldsymbol{\mu}}, \mathbf{q}). \tag{30}\end{aligned}$$

where in (29) we use the fact that $\bar{\mathbf{q}} \in \Delta_{\mathfrak{I}}$. Combining (30) with (27) shows that $\hat{\alpha}_m$ converges to $\text{Mix}_{\Phi}(\mathbf{d}, \mathbf{q}) = \inf_{\hat{\boldsymbol{\mu}} \in \Delta_{\mathfrak{I}}} \langle \hat{\boldsymbol{\mu}}, \mathbf{d} \rangle + D_{\Phi}(\hat{\boldsymbol{\mu}}, \mathbf{q})$. Since ($\hat{\alpha}_m$) was any convergent subsequence of (α_m) (which is bounded), the result follows. \square

C Proofs of Results in the Main Body

C.1 Proof of Theorem 4

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}_{ℓ} , as ℓ .

Proof. We will construct a proper support loss $\underline{\ell}$ of ℓ .

Let $\mathbf{p} \in \text{ri } \Delta_n$ ($-\mathbf{p} \in \text{int dom } \sigma_{\mathcal{S}_{\ell}}$). Since the support function of a non-empty set is closed and convex, we have $\sigma_{\mathcal{S}_{\ell}}^{**} = \sigma_{\mathcal{S}_{\ell}}$ [7, Prop. C.2.1.2]. Pick any $\mathbf{v} \in \partial \sigma_{\mathcal{S}_{\ell}}^{**}(-\mathbf{p}) = \partial \sigma_{\mathcal{S}_{\ell}}^{*}(-\mathbf{p}) \neq \emptyset$. Since $\sigma_{\mathcal{S}_{\ell}}^{*} = \iota_{\mathcal{S}_{\ell}}$ [11], we can apply Proposition 1-(iv) with f replaced by $\sigma_{\mathcal{S}_{\ell}}^{*}$ to obtain $\langle -\mathbf{p}, \mathbf{v} \rangle = \sigma_{\mathcal{S}_{\ell}}^{*}(-\mathbf{p}) + \iota_{\mathcal{S}_{\ell}}(\mathbf{v})$. The fact that $\langle -\mathbf{p}, \mathbf{v} \rangle$ and $\sigma_{\mathcal{S}_{\ell}}^{*}(-\mathbf{p})$ are both finite implies that $\iota_{\mathcal{S}_{\ell}}(\mathbf{v}) = 0$. Therefore, $\mathbf{v} \in \mathcal{S}_{\ell}$ and $\langle \mathbf{p}, \mathbf{v} \rangle = -\sigma_{\mathcal{S}_{\ell}}^{*}(-\mathbf{p}) = \underline{L}_{\ell}(\mathbf{p})$. Define $\underline{\ell}(\mathbf{p}) := \mathbf{v} \in \mathcal{S}_{\ell}$.

Now let $\mathbf{p} \in \text{rbd } \Delta_n$ and $\mathbf{q} := \mathbf{1}_n/n$. Since the \underline{L}_{ℓ} is a closed concave function and $\mathbf{q} \in \text{int dom } \underline{L}_{\ell}$, it follows that $\underline{L}_{\ell}(\mathbf{p} + m^{-1}(\mathbf{q} - \mathbf{p})) \xrightarrow{m \rightarrow \infty} \underline{L}_{\ell}(\mathbf{p})$ [7, Prop. B.1.2.5]. Note that $\mathbf{q}_m := \mathbf{p} + m^{-1}(\mathbf{q} - \mathbf{p}) \in \text{ri } \Delta_n, \forall m \in \mathbb{N}$. Now let $v_{x,m} := \ell_x(\mathbf{q}_m)$, where $\ell(\mathbf{q}_m)$ is as constructed in the previous paragraph. If $(v_{1,m})$ is bounded [resp. unbounded], we can extract a subsequence $(v_{1,\varphi_1(m)})$ which converges [resp. diverges to $+\infty$], where $\varphi_1: \mathbb{N} \rightarrow \mathbb{N}$ is an increasing function. By repeating this process for $(v_{2,\varphi_1(m)})$ and so on, we can construct an increasing function $\varphi := \varphi_n \circ \dots \circ \varphi_1: \mathbb{N} \rightarrow \mathbb{N}$, such that $\mathbf{v}_m := [v_{x,\varphi(m)}]_{x \in [n]}^{\top}$ has a well defined (coordinate-wise) limit in $[0, +\infty]^n$. Define $\underline{\ell}(\mathbf{p}) := \lim_{m \rightarrow \infty} \mathbf{v}_m$. By continuity of the inner product, we have

$$\begin{aligned}\langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle &= \lim_{m \rightarrow \infty} \langle \mathbf{q}_{\varphi(m)}, \mathbf{v}_m \rangle = \lim_{m \rightarrow \infty} \langle \mathbf{q}_{\varphi(m)}, \underline{\ell}(\mathbf{q}_{\varphi(m)}) \rangle, \\ &= \lim_{m \rightarrow \infty} \underline{L}_{\ell}(\mathbf{q}_{\varphi(m)}) = \underline{L}_{\ell}(\mathbf{p}).\end{aligned}$$

By construction, $\forall m \in \mathbb{N}, \mathbf{p}_m := \mathbf{q}_{\varphi(m)} \in \text{ri } \Delta_n$ and $\underline{\ell}(\mathbf{p}_m) = \mathbf{v}_m \xrightarrow{m \rightarrow \infty} \underline{\ell}(\mathbf{p})$. Therefore, $\underline{\ell}$ is support loss of ℓ .

It remains to show that it is proper; that is $\forall \mathbf{p} \in \Delta_n, \forall \mathbf{q} \in \Delta_n, \langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \underline{\ell}(\mathbf{q}) \rangle$. Let $\mathbf{q} \in \text{ri } \Delta_n$. We just showed that $\forall \mathbf{p} \in \Delta_n, \langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle = \underline{L}_{\ell}(\mathbf{p})$ and that $\underline{\ell}(\mathbf{q}) \in \mathcal{S}_{\ell}$. Using the fact that $\underline{L}_{\ell}(\mathbf{p}) = \inf_{\mathbf{z} \in \mathcal{S}_{\ell}} \langle \mathbf{p}, \mathbf{z} \rangle$, we obtain $\langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \underline{\ell}(\mathbf{q}) \rangle$.

Now let $\mathbf{q} \in \text{rbd } \Delta_k$. Since $\underline{\ell}$ is a support loss, we know that there exists a sequence $(\mathbf{q}_m) \subset \text{ri } \Delta_n$ such that $\underline{\ell}(\mathbf{q}_m) \xrightarrow{m \rightarrow \infty} \underline{\ell}(\mathbf{q})$. But as we established in the previous paragraph, $\langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \underline{\ell}(\mathbf{q}_m) \rangle$. By passing to the limit $m \rightarrow \infty$, we obtain $\langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \underline{\ell}(\mathbf{q}) \rangle$. Therefore $\underline{\ell}$ is a proper loss with Bayes risk \underline{L}_{ℓ} . \square

C.2 Proofs of Theorem 5 and Proposition 12

For a set \mathcal{C} , we denote $\text{co } \mathcal{C}$ and $\overline{\text{co}} \mathcal{C}$ its *convex hull* and *closed convex hull*, respectively.

Definition 16 ([7]). Let \mathcal{C} be non-empty convex set in \mathbb{R}^n . We say that $\mathbf{u} \in \mathcal{C}$ is an *extreme point* of \mathcal{C} if there are no two different points \mathbf{u}_1 and \mathbf{u}_2 in \mathcal{C} and $\lambda \in]0, 1[$ such that $\mathbf{u} = \lambda \mathbf{u}_1 + (1 - \lambda) \mathbf{u}_2$.

We denote the set of extreme points of a set \mathcal{C} by $\text{ext } \mathcal{C}$.

Lemma 17. Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a closed loss. Then $\text{ext } \overline{\text{co}} \mathcal{S}_\ell \subseteq \mathcal{S}_\ell$.

Proof. Since $\text{co } \mathcal{S}_\ell \subseteq \mathbb{R}^n$ is connected, $\text{co } \mathcal{S}_\ell = \{\mathbf{v} + \sum_{k=1}^n \alpha_k \ell(\mathbf{a}_k) : (\mathbf{a}_k)_{k \in [n]}, \alpha, \mathbf{v}) \in \mathcal{A}^n \times \Delta_n \times [0, +\infty]^n\}$ [7, Prop. A.1.3.7].

We claim that $\overline{\text{co}} \mathcal{S}_\ell = \text{co } \mathcal{S}_\ell$. Let $(\mathbf{z}_m) := (\mathbf{v}_m + \sum_{k=1}^n \alpha_{m,k} \ell(\mathbf{a}_{m,k}))$ be a convergent sequence in $[0, +\infty]^n$, where (α_m) , $([\mathbf{a}_{m,k}]_{k \in [n]})$ and (\mathbf{v}_m) are sequences in Δ_n , \mathcal{A}^n , and $[0, +\infty]^n$, respectively. Since Δ_n is compact, we may assume, by extracting a subsequence if necessary, that $\alpha_m \xrightarrow{m \rightarrow \infty} \alpha^* \in \Delta_n$. Let $\mathcal{K} := \{k \in [n] : \alpha_k^* \neq 0\}$. Since \mathbf{z}_m converges, $([\ell(\mathbf{a}_{m,k})]_{k \in \mathcal{K}}, \mathbf{v}_m)$ is a bounded sequence in $[0, +\infty]^{|\mathcal{K}|+n}$. Since ℓ is closed, we may assume, by extraction a subsequence if necessary, that $\forall k \in \mathcal{K}$, $\ell(\mathbf{a}_{m,k}) \xrightarrow{m \rightarrow \infty} \ell(\mathbf{a}_k^*)$ and $\mathbf{v}_m \xrightarrow{m \rightarrow \infty} \mathbf{v}^*$, where $[\mathbf{a}_k^*]_{k \in \mathcal{K}} \in \mathcal{A}^{|\mathcal{K}|}$ and $\mathbf{v}^* \in [0, +\infty]^n$. Consequently,

$$\begin{aligned} \mathbf{v}^* + \sum_{k=1}^n \alpha_k^* \ell(\mathbf{a}_k^*) &= \lim_{m \rightarrow \infty} \left[\mathbf{v}_{m,k} + \sum_{k \in \mathcal{K}} \alpha_{m,k} \ell(\mathbf{a}_{m,k}) \right], \\ &\leq \lim_{m \rightarrow \infty} \mathbf{z}_m, \end{aligned}$$

where the last inequality is coordinate-wise. Therefore, there exists $\mathbf{v}' \in [0, +\infty]^n$ such that $\lim_{m \rightarrow \infty} \mathbf{z}_m = \mathbf{v}' + \mathbf{v}^* + \sum_{k=1}^n \alpha_k^* \ell(\mathbf{a}_k^*) \in \text{co } \mathcal{S}_\ell$. This shows that $\overline{\text{co}} \mathcal{S}_\ell \subset \text{co } \mathcal{S}_\ell$, and thus $\overline{\text{co}} \mathcal{S}_\ell = \text{co } \mathcal{S}_\ell$ which proves our first claim.

By definition of an extreme point, $\text{ext } \overline{\text{co}} \mathcal{S}_\ell \subseteq \overline{\text{co}} \mathcal{S}_\ell$. Let $\mathbf{e} \in \text{ext } \overline{\text{co}} \mathcal{S}_\ell$ and $(\mathbf{a}_{k \in [n]}, \alpha, \mathbf{v}) \in \mathcal{A}^n \times \Delta_n \times [0, +\infty]^n$ such that $\mathbf{e} = \sum_{k=1}^n \alpha_k \ell(\mathbf{a}_k) + \mathbf{v}$. If there exists $i, j \in [n]$ such that $\alpha_i \alpha_j \neq 0$ or $\alpha_i v_j \neq 0$ then \mathbf{e} would violate the definition of an extreme point. Therefore, the only possible extreme points are of the form $\{\ell(\mathbf{a}) : \mathbf{a} \in \text{dom } \ell\} = \mathcal{S}_\ell$. \square

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

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

Proof. Let $\mathbf{p} \in \text{ri } \Delta_n$ and suppose that \underline{L}_ℓ is differentiable at \mathbf{p} . In this case, $\sigma_{\mathcal{S}_\ell}$ is differentiable at $-\mathbf{p}$, which implies [7, Cor. D.2.1.4]

$$\mathcal{F}(\mathbf{p}) := \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p}) = \{\nabla \sigma_{\mathcal{S}_\ell}(-\mathbf{p})\}. \quad (31)$$

On the other hand, the fact that $\sigma_{\mathcal{S}_\ell} = \sigma_{\overline{\text{co}} \mathcal{S}_\ell}$ [7, Prop. C.2.2.1], implies $\mathcal{F}(\mathbf{p}) = \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p}) = \partial \sigma_{\overline{\text{co}} \mathcal{S}_\ell}(-\mathbf{p})$. The latter being an *exposed face* of $\overline{\text{co}} \mathcal{S}_\ell$ implies that every extreme point of $\mathcal{F}(\mathbf{p})$ is also an extreme point of $\overline{\text{co}} \mathcal{S}_\ell$ [7, Prop. A.2.3.7, Prop. A.2.4.3]. Therefore, from (31), $\underline{\ell}(\mathbf{p}) = \nabla \sigma_{\mathcal{S}_\ell}(-\mathbf{p})$ is the only extreme point of $\mathcal{F}(\mathbf{p}) \subset \overline{\text{co}} \mathcal{S}_\ell$. From Lemma 17, there exists $\mathbf{a}_* \in \mathcal{A}$ such that $\ell(\mathbf{a}_*) = \underline{\ell}(\mathbf{p})$. In this paragraph, we showed the following

$$\forall \mathbf{p} \in \text{ri } \Delta_n, \exists \mathbf{a}_* \in \text{dom } \ell, \ell(\mathbf{a}_*) = \underline{\ell}(\mathbf{p}). \quad (32)$$

For the rest of this proof we will assume that \underline{L}_ℓ is differentiable on $]0, +\infty]^n$. Let $\mathbf{p} \in \text{ri } \Delta_n \cap \text{dom } \underline{\ell}$. Since $\underline{\ell}$ is a support loss, there exists (\mathbf{p}_m) in $\text{ri } \Delta_n$ such that $(\underline{\ell}(\mathbf{p}_m))_m$ converges to $\underline{\ell}(\mathbf{p})$. From (32) it holds that $\forall \mathbf{p}_m \in \text{ri } \Delta_n, \exists \mathbf{a}_m \in \mathcal{A}, \ell(\mathbf{a}_m) = \underline{\ell}(\mathbf{p}_m)$. Since $(\ell(\mathbf{a}_m))_m$ converges and ℓ is closed, there exists $\mathbf{a}_* \in \mathcal{A}$ such that $\ell(\mathbf{a}_*) = \lim_{m \rightarrow \infty} \ell(\mathbf{a}_m) = \underline{\ell}(\mathbf{p})$.

Now let $\mathbf{a} \in \text{dom } \ell$ and $f(\mathbf{p}, x) := \underline{\ell}_x(\mathbf{p}) - \ell_x(\mathbf{a})$. Since $\ell(\mathbf{a}) \in \mathcal{S}_\ell$ and $\underline{\ell}$ is proper, we have for all $\mathbf{p} \in \text{ri } \Delta_n, \mathbb{E}_{x \sim \mathbf{p}}[f(\mathbf{p}, x)] \leq 0$ and $-\infty < f(\mathbf{p}, x), \forall x \in [n]$. Therefore, Lemma 5 implies that for

all $m \in \mathbb{N} \setminus \{0\}$ there exists $\mathbf{p}_m \in \text{ri } \Delta_n$, such that $\forall x \in [n], \ell_x(\mathbf{p}_m) \leq \ell_x(\mathbf{a}) + 1/m$. On one hand, since $(\ell(\mathbf{p}_m))$ is bounded (from the previous inequality), we may assume by extracting a subsequence if necessary, that $(\ell(\mathbf{p}_m))_m$ converges. On the other hand, since $\mathbf{p}_m \in \text{ri } \Delta_n$, (32) implies that there exists $\mathbf{a}_m \in \text{dom } \ell$ such that $\ell(\mathbf{p}_m) = \ell(\mathbf{a}_m)$. Since ℓ is closed and $(\ell(\mathbf{a}_m))_m$ converges, there exists $\mathbf{a}_* \in \mathcal{A}$, such that $\ell(\mathbf{a}_*) = \lim_{m \rightarrow \infty} \ell(\mathbf{a}_m) = \lim_{m \rightarrow \infty} \ell(\mathbf{p}_m) \leq \ell(\mathbf{a})$. But since ℓ is admissible, the latter component-wise inequality implies that $\ell(\mathbf{a}_*) = \ell(\mathbf{a}) = \lim_{m \rightarrow \infty} \ell(\mathbf{p}_m)$. \square

Lemma 18. *Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss satisfying Assumption 1. If \underline{L}_ℓ is not differentiable at \mathbf{p} then there exist $\mathbf{a}_0, \mathbf{a}_1 \in \text{dom } \ell$, such that $\ell(\mathbf{a}_0) \neq \ell(\mathbf{a}_1)$ and $\underline{L}_\ell(\mathbf{p}) = \langle \mathbf{p}, \ell(\mathbf{a}_0) \rangle = \langle \mathbf{p}, \ell(\mathbf{a}_1) \rangle$.*

Proof. Suppose \underline{L}_ℓ is not differentiable at $\mathbf{p} \in \text{ri } \Delta_n$. Then from the definition of the Bayes risk, $\sigma_{\mathcal{S}_\ell}$ is not differentiable at $-\mathbf{p}$. This implies that $\mathcal{F}(\mathbf{p}) := \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p})$ has more than one element [7, Cor. D.2.1.4]. Since $\sigma_{\mathcal{S}_\ell} = \sigma_{\overline{\text{co}} \mathcal{S}_\ell}$ (ibid., Prop. C.2.2.1), $\mathcal{F}(\mathbf{p}) = \partial \sigma_{\overline{\text{co}} \mathcal{S}_\ell}(-\mathbf{p})$ is a subset of $\overline{\text{co}} \mathcal{S}_\ell$ and every extreme point of $\mathcal{F}(\mathbf{p})$ is also an extreme point of $\overline{\text{co}} \mathcal{S}_\ell$ (ibid., Prop. A.2.3.7). Thus, from Lemma 17, we have $\text{ext } \mathcal{F}(\mathbf{p}) \subset \mathcal{S}_\ell$. On the other hand, since $-\mathbf{p} \in \text{int dom } \sigma_{\mathcal{S}_\ell}$, $\mathcal{F}(\mathbf{p})$ is a compact, convex set [11, Thm. 23.4], and thus $\mathcal{F}(\mathbf{p}) = \text{co}(\text{ext } \mathcal{F}(\mathbf{p}))$ [7, Thm. A.2.3.4]. Hence, the fact that $\mathcal{F}(\mathbf{p})$ has more than one element implies there exists $\mathbf{a}_0, \mathbf{a}_1 \in \mathcal{A}$ such that $\ell(\mathbf{a}_0), \ell(\mathbf{a}_1) \in \text{ext } \mathcal{F}(\mathbf{p}) \subseteq \mathcal{F}(\mathbf{p})$ and $\ell(\mathbf{a}_0) \neq \ell(\mathbf{a}_1)$. Since $\mathcal{F}(\mathbf{p}) = \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p})$, Proposition 1-(iv) and the fact that $\sigma_{\mathcal{S}_\ell}^* = \iota_{\mathcal{S}_\ell}$ imply that $\underline{L}_\ell(\mathbf{p}) = \langle \mathbf{p}, \ell(\mathbf{a}_0) \rangle = \langle \mathbf{p}, \ell(\mathbf{a}_1) \rangle$. \square

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

Proof. Let $\mathbf{l} = \{1, 2\}$. Since ℓ is Φ -mixable, it must be $\Phi_{\mathbf{l}}$ -mixable, where $\Phi_{\mathbf{l}} := \Phi_{\mathbf{l}} \circ \Pi_{\mathbf{l}}^T: \mathbb{R}^2 \rightarrow \mathbb{R} \cup \{+\infty\}$ (Proposition 10). Let $\tilde{\Psi} := \Phi_{\mathbf{l}}$.

For $w \in]0, +\infty[$ and $z \in \text{int dom } \tilde{\Psi}^* = \mathbb{R}$ (see appendix E), we define $(\tilde{\Psi}^*)'_\infty(w) := \lim_{t \rightarrow +\infty} [\tilde{\Psi}^*(z + tw) - \tilde{\Psi}^*(z)]/t$. The value of $(\tilde{\Psi}^*)'_\infty(w)$ does not depend on the choice of z , and it holds that $(\tilde{\Psi}^*)'_\infty(w) = \sigma_{\text{dom } \tilde{\Psi}}(w)$ and $(\tilde{\Psi}^*)'_\infty(-w) = \sigma_{\text{dom } \tilde{\Psi}}(-w)$ [7, Prop. C.1.2.2]. In our case, we have $\text{dom } \tilde{\Psi} = [0, 1]$ (by definition of $\tilde{\Psi}$), which implies that $\sigma_{\text{dom } \tilde{\Psi}}(1) = 1$ and $\sigma_{\text{dom } \tilde{\Psi}}(-1) = 0$. Therefore, $(\tilde{\Psi}^*)'_\infty(1) + (\tilde{\Psi}^*)'_\infty(-1) = 1$. As a result $\tilde{\Psi}^*$ cannot be affine. For all $\delta > 0$, let $g_\delta: \mathbb{R} \times \{-1, 0, +1\} \rightarrow \mathbb{R}$ be defined by

$$g_\delta(s, u) := [\tilde{\Psi}^*(s + \delta(u + 1)/2) - \tilde{\Psi}^*(s + \delta(u - 1)/2)]/\delta.$$

Since $\tilde{\Psi}^*$ is convex it must have non-decreasing slopes (ibid., p.13). Combining this with the fact that $\tilde{\Psi}^*$ is not affine implies that

$$\exists s_\delta^* \in \mathbb{R}, g_\delta(s_\delta^*, -1) < g_\delta(s_\delta^*, +1). \quad (33)$$

The fact that $\tilde{\Psi}^*$ has non-decreasing slopes also implies that

$$g_\delta(s_\delta^*, +1) = [\tilde{\Psi}^*(s_\delta^* + \delta) - \tilde{\Psi}^*(s_\delta^*)]/\delta \leq \lim_{t \rightarrow \infty} [\tilde{\Psi}^*(s_\delta^* + t) - \tilde{\Psi}^*(s_\delta^*)]/t = (\tilde{\Psi}^*)'_\infty(1) = 1.$$

Similarly, we have $0 = -(\tilde{\Psi}^*)'_\infty(-1) \leq g_\delta(s_\delta^*, -1)$. Let $\tilde{\mu} \in \partial \tilde{\Psi}^*(s_\delta^*)$. Since $\tilde{\Psi}$ is a closed convex function the following equivalence holds $\tilde{\mu} \in \partial \tilde{\Psi}^*(s_\delta^*) \iff s_\delta^* \in \partial \tilde{\Psi}(\tilde{\mu})$ (ibid., Cor. D.1.4.4). Thus, if $\tilde{\mu} \in \{0, 1\} = \text{bd } \tilde{\Delta}_2$, then $\partial \tilde{\Psi}(\tilde{\mu}) \neq \emptyset$, which is not possible since ℓ is Ψ -mixable (Lemma 13).

[We show $\underline{L}_\ell \in C^1([0, +\infty]^n)$] We will now show that \underline{L}_ℓ is continuously differentiable on $]0, +\infty[^n$. Since \underline{L}_ℓ is 1-homogeneous, it suffices to check the differentiability on $\text{ri } \Delta_n$. Suppose \underline{L}_ℓ is not differentiable at $\mathbf{p} \in \text{ri } \Delta_n$. From Lemma 18, there exists $\mathbf{a}_0, \mathbf{a}_1 \in \mathcal{A}$ such that $\ell(\mathbf{a}_0), \ell(\mathbf{a}_1) \in \partial \sigma_{\mathcal{S}_\ell}(-\mathbf{p})$ and $\ell(\mathbf{a}_0) \neq \ell(\mathbf{a}_1)$. Let $A := [\mathbf{a}_0, \mathbf{a}_1] \in \mathbb{R}^{n \times 2}$, $\delta := \min\{|\ell_x(\mathbf{a}_0) - \ell_x(\mathbf{a}_1)| : x \in [n], |\ell_x(\mathbf{a}_0) - \ell_x(\mathbf{a}_1)| > 0\}$, and $s_\delta^* \in \mathbb{R}$ as in (33). We denote $g^- := g_\delta(s_\delta^*, -1)$ and $g^+ := g_\delta(s_\delta^*, +1) \in]0, 1]$. Let $\tilde{\mu} \in \partial \tilde{\Psi}^*(s_\delta^*) \in \text{int } \tilde{\Delta}_2$ and $\boldsymbol{\mu} = \Pi_2(\tilde{\mu}) \in \text{ri } \Delta_2$. From the fact that ℓ is Ψ -mixable, $J_2^T \ell_x(A) = \ell_x(\mathbf{a}_0) - \ell_x(\mathbf{a}_1)$, and (8), there must exist $\mathbf{a}_* \in \mathcal{A}$ such that for all $x \in [n]$,

$$\begin{aligned} \ell_x(\mathbf{a}_*) &\leq \text{Mix}_\Psi(\ell_x(A), \boldsymbol{\mu}), \\ &= \ell_x(\mathbf{a}_1) + \tilde{\Psi}^*(s_\delta^*) - \tilde{\Psi}^*(s_\delta^* - \ell_x(\mathbf{a}_0) + \ell_x(\mathbf{a}_1)), \end{aligned}$$

and by letting sgn be the *sign* function

$$\leq \ell_x(a_1) + g_\delta(s_\delta^*, -\text{sgn}[\ell_x(a_0) - \ell_x(a_1)])[\ell_x(a_0) - \ell_x(a_1)], \quad (34)$$

where in (34) we used the fact that $\tilde{\Psi}^*$ has non-decreasing slopes and the definition of δ . When $\ell_x(a_0) \leq \ell_x(a_1)$, (34) becomes $\ell_x(a_*) \leq (1 - g^+)\ell_x(a_1) + g^+\ell_x(a_0)$. Otherwise, we have $\ell_x(a_*) \leq (1 - g^-)\ell_x(a_1) + g^-\ell_x(a_0) < (1 - g^+)\ell_x(a_1) + g^+\ell_x(a_0)$. Since ℓ is admissible, there must exist at least one $x \in [n]$ such that $\ell_x(a_0) > \ell_x(a_1)$. Combining this with the fact that $p_x > 0, \forall x \in [n]$ ($\mathbf{p} \in \text{ri } \Delta_n$), implies that $\langle \mathbf{p}, \ell(a_*) \rangle < \langle \mathbf{p}, (1 - g^+)\ell(a_1) + g^+\ell(a_0) \rangle = \underline{L}_\ell(\mathbf{p})$. This contradicts the fact that $\ell(a_*) \in \mathcal{S}_\ell$. Therefore, \underline{L}_ℓ must be differentiable at \mathbf{p} . As argued earlier, this implies that \underline{L}_ℓ must be differentiable on $]0, +\infty[^n$. Combining this with the fact that \underline{L}_ℓ is concave on $]0, +\infty[^n$, implies that \underline{L}_ℓ is continuously differentiable on $]0, +\infty[^n$ (ibid., Rmk. D.6.2.6).

[We show $\tilde{\Phi}^* \in C^1(\mathbb{R}^{k-1})$] Suppose that $\tilde{\Phi}^*$ is not differentiable at some $\mathbf{s}^* \in \mathbb{R}^{k-1}$. Then there exists $\mathbf{d} \in \mathbb{R}^{k-1} \setminus \{\mathbf{0}_k\}$ such that $-(\tilde{\Phi}^*)'(\mathbf{s}^*; -\mathbf{d}) < (\tilde{\Phi}^*)'(\mathbf{s}^*; \mathbf{d})$. Since $\mathbf{s}^* \in \text{int dom } \tilde{\Phi}^*$, $(\tilde{\Phi}^*)'(\mathbf{s}^*; \cdot)$ is finite and convex [7, Prop. D.1.1.2], and thus it is continuous on $\text{dom } \tilde{\Phi}^* = \mathbb{R}^{k-1}$ (ibid., Rmk. B.3.1.3). Consequently, there exists $\delta^* > 0$ such that

$$\forall \hat{\mathbf{d}} \in \mathbb{R}^{k-1}, \|\hat{\mathbf{d}} - \mathbf{d}\| \leq \delta^* \implies -(\tilde{\Phi}^*)'(\mathbf{s}^*; -\hat{\mathbf{d}}) < (\tilde{\Phi}^*)'(\mathbf{s}^*; \hat{\mathbf{d}}) \quad (35)$$

Let $g: \{-1, 1\} \rightarrow \mathbb{R}$ be such that

$$g(u) := \sup_{\|\hat{\mathbf{d}} - \mathbf{d}\| \leq \delta^*} u \cdot (\tilde{\Phi}^*)'(\mathbf{s}^*; u\hat{\mathbf{d}}).$$

Note that since $\tilde{\Phi}^*$ has increasing slopes ($\tilde{\Phi}^*$ is convex), $g(1) \leq \sup_{\|\hat{\mathbf{d}} - \mathbf{d}\| \leq \delta^*} (\tilde{\Phi}^*)'_\infty(\hat{\mathbf{d}}) = \sup_{\|\hat{\mathbf{d}} - \mathbf{d}\| \leq \delta^*} \sigma_{\text{dom } \tilde{\Phi}}(\hat{\mathbf{d}}) \leq 1$, where the last inequality holds because $\tilde{\Delta}_k \subset \mathcal{B}(\mathbf{0}_k, 1)$, and thus $\sigma_{\text{dom } \tilde{\Phi}}(\hat{\mathbf{d}}) = \sigma_{\tilde{\Delta}_k}(\hat{\mathbf{d}}) \leq \sigma_{\mathcal{B}(\mathbf{0}_k, 1)}(\hat{\mathbf{d}}) = 1$. Let $\Delta g := g(1) - g(-1)$. From (35), it is clear that $\Delta g > 0$.

Suppose that \underline{L}_ℓ is twice differentiable on $]0, +\infty[^n$ and let $\tilde{\ell}$ be a support loss of ℓ . By definition of a support loss, $\forall \mathbf{p} \in \text{ri } \Delta_k, \tilde{\ell}(\tilde{\mathbf{p}}) = \underline{\ell}(\mathbf{p}) = \nabla \underline{L}_\ell(\mathbf{p})$ (where $\tilde{\ell} := \underline{\ell} \circ \Pi_n$). Thus, since \underline{L}_ℓ is twice differentiable on $]0, +\infty[^n$, $\tilde{\ell}$ is differentiable on $\text{int } \tilde{\Delta}_n$. Furthermore, $\tilde{\ell}$ is continuous on $\text{ri } \Delta_k$ given that $\underline{L}_\ell \in C^1(]0, +\infty[^n)$ as shown in the first part of this proof. We may assume without loss of generality that ℓ is not a constant function. Thus, from Theorem 5, $\tilde{\ell}$ is not a constant function either. Consequently, the mean value theorem applied to $\tilde{\ell}$ (see e.g. [12, Thm. 5.10]) between any two points in $\text{ri } \Delta_n$ with distinct images under $\tilde{\ell}$, implies that there exists $(\tilde{\mathbf{p}}_*, \mathbf{v}_*) \in \text{int } \tilde{\Delta}_n \times \mathbb{R}^{n-1}$, such that $D\tilde{\ell}(\tilde{\mathbf{p}}_*)\mathbf{v}_* \neq \mathbf{0}_n$. For the rest of the proof let $(\tilde{\mathbf{p}}, \mathbf{v}) := (\tilde{\mathbf{p}}_*, \mathbf{v}_*)$ and define $\mathfrak{I} := \{x \in [n] : D\tilde{\ell}_x(\tilde{\mathbf{p}})\mathbf{v} \neq 0\}$. From Lemma 8, we have $\langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_n^T$, which implies that there exists $x \in \mathfrak{I}, D\tilde{\ell}_x(\tilde{\mathbf{p}})\mathbf{v} > 0$. Thus, the set

$$\mathfrak{K} := \left\{x \in \mathfrak{I} : D\tilde{\ell}_x(\tilde{\mathbf{p}})\mathbf{v} > 0\right\} \quad (36)$$

is non-empty. From this and the fact that $\mathbf{p} \in \text{ri } \Delta_n$, it follows that

$$\sum_{x' \in \mathfrak{K}} p_{x'} D\tilde{\ell}_{x'}(\tilde{\mathbf{p}})\mathbf{v} > 0. \quad (37)$$

Let $\tilde{\mathbf{p}}^t := \tilde{\mathbf{p}} + t\mathbf{v}$. From Taylor's Theorem (see e.g. [?, §151]) applied to the function $t \mapsto \tilde{\ell}(\tilde{\mathbf{p}}^t)$, there exists $\epsilon^* > 0$ and functions $\delta_x: [-\epsilon^*, \epsilon^*] \rightarrow \mathbb{R}^n, x \in [n]$, such that $\lim_{t \rightarrow 0} t^{-1}\delta_x(t) = 0$ and

$$\forall |t| \leq \epsilon^*, \quad \underline{\ell}_x(\mathbf{p}^t) = \underline{\ell}_x(\mathbf{p}) + t D\tilde{\ell}_x(\mathbf{p})\mathbf{v} + \delta_x(t). \quad (38)$$

For $x \in [n]$, let $\mathbf{d}_x \in \mathbb{R}^{\tilde{k}}$ and suppose that $\|\mathbf{d}_x - \mathbf{d}\| \leq \delta^*$ (we will define \mathbf{d}_x explicitly later). By shrinking ϵ^* if necessary, we may assume that

$$\forall x \in \mathcal{I}, \forall \theta \in [k], \forall |t| \leq \epsilon^*, \quad d_\theta t^{-1} \delta_x(t) \leq \frac{\delta^* |\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}|}{\sqrt{n}}, \quad (39)$$

$$\forall x \notin \mathcal{I}, \quad \tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \left[\delta_x\left(\epsilon^* \frac{d_\theta}{\|\mathbf{d}\|}\right)\right]_{\theta \in [\tilde{k}]}\right) \leq \epsilon^* \frac{\Delta g}{4\|\mathbf{d}\|} \sum_{x' \in \mathcal{K}} p_{x'} \mathbf{D}_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v}, \quad (40)$$

$$\begin{aligned} \forall x \in [n], \quad \tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right) &\leq -(\tilde{\Phi}^*)'\left(\mathbf{s}^*; -\epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right) \\ &\quad + \epsilon^* \frac{\Delta g}{4\|\mathbf{d}\|} \sum_{x' \in \mathcal{K}} p_{x'} \mathbf{D}_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v}, \end{aligned} \quad (41)$$

where (41) is satisfied for small enough ϵ^* because of (37) and the fact that

$$\frac{1}{\epsilon} \left(\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \epsilon \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right) \right) \xrightarrow{\epsilon \rightarrow 0} -(\tilde{\Phi}^*)'\left(\mathbf{s}^*; -\frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right),$$

and (40) is also satisfied for small enough ϵ^* because $\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \left[\delta_x\left(\epsilon \frac{d_\theta}{\|\mathbf{d}\|}\right)\right]_{\theta \in [\tilde{k}]}\right) = O\left(\max_{\theta \in [\tilde{k}]} \left|\delta_x\left(\epsilon \frac{d_\theta}{\|\mathbf{d}\|}\right)\right|\right) = o(\epsilon)$, where the first equality is due to the fact that $(\lambda, \mathbf{z}) \mapsto \frac{1}{\lambda} \left(\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*(\mathbf{s}^* - \lambda \mathbf{z}) \right)$ is uniformly bounded on compact subsets of $\mathbb{R} \times \mathbb{R}^{\tilde{k}}$ (by continuity of the directional derivative $(\tilde{\Phi}^*)'(\mathbf{s}^*; \cdot)$).

If $\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v} \leq 0$, then by the positive homogeneity of the directional derivative, the definition of the function g , and (41), we get

$$\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right) \leq \epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} g(1) + \epsilon^* \frac{\Delta g}{4\|\mathbf{d}\|} \sum_{x' \in \mathcal{K}} p_{x'} \mathbf{D}_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v}. \quad (42)$$

On the other hand, if $\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v} > 0$, then from the monotonicity of the slopes of $\tilde{\Phi}^*$, the positive homogeneity of the directional derivative, and the definition of the function g , it follows that

$$\begin{aligned} \frac{1}{\epsilon^*} \left(\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right) \right) &\leq -(\tilde{\Phi}^*)'\left(\mathbf{s}^*; -\frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x\right), \\ &= -\frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} (\tilde{\Phi}^*)'(\mathbf{s}^*; -\mathbf{d}_x), \\ &\leq \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} g(-1), \\ &= \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} (-\Delta g + g(1)). \end{aligned} \quad (43)$$

Let $\lambda_\theta := \epsilon^* \frac{d_\theta}{\|\mathbf{d}\|}$, for $\theta \in [\tilde{k}]$. From Theorem 5, there exists $[\mathbf{a}_\theta]_{\theta \in [k]} \in \mathcal{A}^k$, such that

$$\ell(\mathbf{a}_k) = \underline{\ell}(\mathbf{p}) \quad \text{and} \quad \forall \theta \in [\tilde{k}], \quad \ell(\mathbf{a}_\theta) = \underline{\ell}(\mathbf{p}^{\lambda_\theta}) = \underline{\ell}(\mathbf{p}) + \epsilon^* \frac{d_\theta}{\|\mathbf{d}\|} \mathbf{D}_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v} + \delta\left(\epsilon^* \frac{d_\theta}{\|\mathbf{d}\|}\right), \quad (44)$$

where $[\delta(\cdot)]_x := \delta_x(\cdot)$ for $x \in [n]$.

From the fact that ℓ is Φ -mixable, it follows that there exists $\mathbf{a}_* \in \mathcal{A}$ such that for all $x \in [n]$,

$$\ell_x(\mathbf{a}_*) \leq \text{Mix}_\Phi(\ell_x(\mathbf{a}_{1:k}), \boldsymbol{\mu}) = \ell_x(\mathbf{a}_k) + \tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - J_k^\top \ell_x(\mathbf{a}_{1:k})\right). \quad (45)$$

For $x \in [n]$, we now define $\mathbf{d}_x \in \mathbb{R}^{\tilde{k}}$ explicitly as

$$\forall \theta \in [\tilde{k}], \quad d_{x,\theta} := \begin{cases} d_\theta + \frac{\|\mathbf{d}\|}{\epsilon^* [\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}]} \delta_x\left(\epsilon^* \frac{d_\theta}{\|\mathbf{d}\|}\right), & \text{if } x \in \mathcal{I} \\ d_\theta, & \text{otherwise.} \end{cases}$$

From (39), we have $\|\mathbf{d}_x - \mathbf{d}\| \leq \delta^*, \forall x \in [n]$. Furthermore, from (44) and the fact that for all $x \in [n]$, $J_k^\top \ell_x(\mathbf{a}_{1:k}) = [\ell_x(\mathbf{a}_\theta) - \ell_x(\mathbf{a}_k)]_{\theta \in [\tilde{k}]}$, we have

$$J_k^\top \ell_x(\mathbf{a}_{1:k}) = \begin{cases} \epsilon^* \frac{\mathbf{D}_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\mathbf{d}\|} \mathbf{d}_x, & \text{if } x \in \mathcal{I}; \\ \left[\delta_x\left(\epsilon^* \frac{d_\theta}{\|\mathbf{d}\|}\right)\right]_{\theta \in [\tilde{k}]}, & \text{otherwise.} \end{cases} \quad (46)$$

Using this, together with (42) and (43), we get $\forall x \in \mathcal{I}$,

$$\begin{aligned} \tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*(\mathbf{s}^* - J_k^\top \ell_x(\mathbf{a}_{1:k})) &= \tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \epsilon^* \frac{D_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\tilde{\mathbf{d}}\|} \mathbf{d}_x\right), \\ &\leq \epsilon^* \frac{D_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\tilde{\mathbf{d}}\|} g(1) - \epsilon^* \Delta g \frac{D_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v}}{\|\tilde{\mathbf{d}}\|} \mathbb{1}_{\{D_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v} > 0\}} \\ &\quad + \epsilon^* \frac{\Delta g}{4\|\tilde{\mathbf{d}}\|} \mathbb{1}_{\{D_{\tilde{\ell}_x}(\tilde{\mathbf{p}})\mathbf{v} \leq 0\}} \sum_{x' \in \mathcal{R}} p_{x'} D_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v}. \end{aligned} \quad (47)$$

Combining (45), (46), and (47) yields

$$\begin{aligned} \langle \mathbf{p}, \ell(\mathbf{a}_*) \rangle &\leq \langle \mathbf{p}, \ell(\mathbf{a}_k) \rangle + \frac{\epsilon^*}{\|\tilde{\mathbf{d}}\|} \langle \mathbf{p}, D_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v} \rangle g(1) - \frac{3\epsilon^* \Delta g}{4\|\tilde{\mathbf{d}}\|} \sum_{x' \in \mathcal{R}} p_{x'} D_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v} \\ &\quad + \sum_{x \notin \mathcal{I}} p_x \left(\tilde{\Phi}^*(\mathbf{s}^*) - \tilde{\Phi}^*\left(\mathbf{s}^* - \left[\delta_x \left(\epsilon^* \frac{d_\theta}{\|\tilde{\mathbf{d}}\|} \right) \right]_{\theta \in [\tilde{k}]}\right) \right), \end{aligned}$$

using (40) and the fact that $\langle \mathbf{p}, D_{\tilde{\ell}}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_n^\top$ (see Lemma 8), we get

$$\leq \langle \mathbf{p}, \ell(\mathbf{a}_k) \rangle - \frac{\epsilon^* \Delta g}{2\|\tilde{\mathbf{d}}\|} \sum_{x' \in \mathcal{R}} p_{x'} D_{\tilde{\ell}_{x'}}(\tilde{\mathbf{p}})\mathbf{v}, \quad (48)$$

$$< \langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle, \quad (49)$$

where in (49) we used (37) and the fact that $\underline{\ell}(\mathbf{p}) = \ell(\mathbf{a}_k)$ (see (45)). Equation 49 shows that $\ell(\mathbf{a}^*) \notin \mathcal{S}_\ell$, which is a contradiction. \square

C.3 Proof of Theorem 7

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

Proof. Let $\eta := \underline{\eta}_\ell$. We will show that $\exp(-\eta \mathcal{S}_\ell)$ is convex, which will imply that ℓ is η -mixable [5].

Since $\underline{\eta}_\ell = \inf_{\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n} (\lambda_{\max}([\mathbf{H}\tilde{\underline{L}}_{\log}(\tilde{\mathbf{p}})]^{-1} \mathbf{H}\tilde{\underline{L}}_\ell(\tilde{\mathbf{p}})))^{-1} > 0$, $\eta \underline{L}_\ell - \underline{L}_{\log}$ is convex on $\text{ri } \Delta_n$ [14, Thm. 10]. Let $\mathbf{p} \in \text{ri } \Delta_n$ and define

$$\Lambda(\mathbf{r}) := \underline{L}_{\log}(\mathbf{r}) + \langle \mathbf{r}, \eta \underline{\ell}(\mathbf{p}) - \ell_{\log}(\mathbf{p}) \rangle, \quad \mathbf{r} \in \text{ri } \Delta_n.$$

Since Λ is equal to \underline{L}_{\log} plus an affine function, it follows that $\eta \underline{L}_\ell - \Lambda$ is also convex on $\text{ri } \Delta_n$. On the one hand, since $\underline{\ell}$ and ℓ_{\log} are proper losses, we have $\langle \mathbf{p}, \underline{\ell}(\mathbf{p}) \rangle = \underline{L}_\ell(\mathbf{p})$ and $\langle \mathbf{p}, \ell_{\log}(\mathbf{p}) \rangle = \underline{L}_{\log}(\mathbf{p})$ which implies that

$$\eta \underline{L}_\ell(\mathbf{p}) - \Lambda(\mathbf{p}) = 0. \quad (50)$$

On the other hand, since \underline{L}_ℓ and \underline{L}_{\log} are differentiable we have $\underline{\ell}(\mathbf{p}) = \nabla \underline{L}_\ell(\mathbf{p})$ and $\nabla \underline{L}_{\log}(\mathbf{p}) = \ell_{\log}(\mathbf{p})$, which yields $\eta \nabla \underline{L}_\ell(\mathbf{p}) - \nabla \Lambda(\mathbf{p}) = \mathbf{0}_n$. This implies that $\eta \underline{L}_\ell - \Lambda$ attains a minimum at \mathbf{p} [7, Thm. D.2.2.1]. Combining this fact with (50) gives $\eta \underline{L}_\ell(\mathbf{r}) \geq \Lambda(\mathbf{r})$, $\forall \mathbf{r} \in \text{ri } \Delta_n$, or equivalently $-\eta \underline{L}_\ell \leq -\Lambda$. By Proposition 1-(iii), this implies

$$[-\eta \underline{L}_\ell]^* \geq [-\Lambda]^*. \quad (51)$$

Using Proposition 1-(ii), we get $[-\Lambda]^*(\mathbf{s}) = [-\underline{L}_{\log}]^*(\mathbf{s} - \ell_{\log}(\mathbf{p}) + \eta \underline{\ell}(\mathbf{p}))$ for $\mathbf{s} \in \mathbb{R}^n$. Since $-\eta \underline{L}_\ell(\mathbf{u}) = -\underline{L}_\ell(\eta \mathbf{u}) = \sigma_{\mathcal{S}_\ell}(-\eta \mathbf{u})$ and $\sigma_{\mathcal{S}_\ell} = \iota_{\mathcal{S}_\ell}$, Proposition 1-(v) implies $[-\eta \underline{L}_\ell]^*(\mathbf{s}) = \iota_{\mathcal{S}_\ell}(-\mathbf{s}/\eta)$. Similarly, we have $[-\underline{L}_{\log}]^*(\mathbf{s}) = \iota_{\mathcal{S}_{\log}}(-\mathbf{s})$. Therefore, (51) implies

$$\forall \mathbf{s} \in \mathbb{R}^n, \quad \iota_{\mathcal{S}_\ell}(-\mathbf{s}/\eta) \geq \iota_{\mathcal{S}_{\log}}(-\mathbf{s} + \ell_{\log}(\mathbf{p}) - \eta \underline{\ell}(\mathbf{p})).$$

This inequality implies that if $\mathbf{s} \in -\eta \mathcal{S}_\ell$, then $\mathbf{s} \in -\mathcal{S}_{\log} + \ell_{\log}(\mathbf{p}) - \eta \underline{\ell}(\mathbf{p})$. In particular, if $\mathbf{u} \in e^{-\eta \mathcal{S}_\ell}$ then

$$\mathbf{u} \in e^{-\mathcal{S}_{\log} + \ell_{\log}(\mathbf{p}) - \eta \underline{\ell}(\mathbf{p})} \subseteq \mathcal{H}_{\tau(\mathbf{p}), 1} = \{\mathbf{v} \in \mathbb{R}^n : \langle \mathbf{v}, \mathbf{p} \odot e^{\eta \underline{\ell}(\mathbf{p})} \rangle \leq 1\}. \quad (52)$$

To see the set inclusion in (52), consider $\mathbf{s} \in -\mathcal{S}_{\log} + \ell_{\log}(\mathbf{p}) - \eta\ell(\mathbf{p})$, then by definition of the superprediction set \mathcal{S}_{\log} there exists $\mathbf{r} \in \Delta_n$ and $\mathbf{v} \in [0, +\infty[^n$, such that $\mathbf{s} = \log \mathbf{r} - \log \mathbf{p} - \eta\ell(\mathbf{p}) - \mathbf{v}$. Thus,

$$\langle e^{\mathbf{s}}, \mathbf{p} \odot e^{\eta\ell(\mathbf{p})} \rangle = \langle \mathbf{r}, e^{-\mathbf{v}} \rangle \leq 1, \quad (53)$$

where the inequality is true because $\mathbf{r} \in \Delta_n$ and $\mathbf{v} \in [0, +\infty[^n$. The above argument shows that $e^{-\eta\mathcal{S}_{\ell}} \subseteq \mathcal{H}_{\tau(\mathbf{p}),1}$, where $\tau(\mathbf{p}) := \mathbf{p} \odot e^{\eta\ell(\mathbf{p})}$. Furthermore, $e^{-\eta\mathcal{S}_{\ell}} \subseteq \mathcal{H}_{\tau(\mathbf{p}),1} \cap]0, +\infty[^n$, since all elements of $e^{-\eta\mathcal{S}_{\ell}}$ have non-negative, finite components. The latter set inclusion still holds for $\hat{\mathbf{p}} \in \text{ri } \Delta_n$. In fact, from the definition of a support loss, there exists a sequence (\mathbf{p}_m) in $\text{ri } \Delta_n$ converging to $\hat{\mathbf{p}}$ such that $\ell(\mathbf{p}_m) \xrightarrow{m \rightarrow \infty} \ell(\hat{\mathbf{p}})$. Equation 53 implies that for $\mathbf{u} \in e^{-\eta\mathcal{S}_{\ell}}$, $\langle \mathbf{u}, \mathbf{p}_m \odot e^{\eta\ell(\mathbf{p}_m)} \rangle \leq 1$. Since the inner product is continuous, by passage to the limit, we obtain $\langle \mathbf{u}, \hat{\mathbf{p}} \odot e^{\eta\ell(\hat{\mathbf{p}})} \rangle \leq 1$. Therefore,

$$e^{-\eta\mathcal{S}_{\ell}} \subseteq \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{\tau(\mathbf{p}),1} \cap]0, +\infty[^n. \quad (54)$$

Now suppose $\mathbf{u} \in \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{\tau(\mathbf{p}),1} \cap]0, +\infty[^n$; that is, for all $\mathbf{p} \in \Delta_n$,

$$\begin{aligned} 1 \geq \langle \mathbf{u}, \mathbf{p} \odot e^{\eta\ell(\mathbf{p})} \rangle &= \langle \mathbf{p}, \mathbf{u} \odot e^{\eta\ell(\mathbf{p})} \rangle = \langle \mathbf{p}, e^{\eta\ell(\mathbf{p}) + \log \mathbf{u}} \rangle, \\ &\geq e^{\langle \mathbf{p}, \eta\ell(\mathbf{p}) \rangle + \langle \mathbf{p}, \log \mathbf{u} \rangle}, \end{aligned} \quad (55)$$

where the first equality is obtained merely by expanding the expression of the inner product, and the second inequality is simply Jensen's Inequality. Since $\mathbf{u} \mapsto e^{\mathbf{u}}$ is strictly convex, the Jensen's inequality in (55) is strict unless $\exists(c, \mathbf{p}) \in \mathbb{R} \times \Delta_n$, such that

$$\eta\ell(\mathbf{p}) + \log \mathbf{u} = c\mathbf{1}_n. \quad (56)$$

By substituting (56) into (55), we get $1 \geq \exp(c)$, and thus $c \leq 0$. Furthermore, (56) together with the fact that $\mathbf{u} \in]0, +\infty[^n$ imply that $\mathbf{p} \in \text{dom } \ell$, and thus there exists $\mathbf{a} \in \text{dom } \ell$ such that $\ell(\mathbf{a}) = \ell(\mathbf{p})$ (Theorem 5). Using this and rearranging (56), we get $\mathbf{u} = \exp(-\eta\ell(\mathbf{a}) + c\mathbf{1})$. Since $c \leq 0$, this means that $\mathbf{u} \in \exp(-\eta\mathcal{S}_{\ell})$. Suppose now that (56) does not hold. In this case, (55) must be a strict inequality for all $\mathbf{p} \in \Delta_n$. By applying the log on both side of (55),

$$\forall \mathbf{p} \in \Delta_n, \eta\ell(\mathbf{p}) + \langle \mathbf{p}, \log \mathbf{u} \rangle = \langle \mathbf{p}, \eta\ell(\mathbf{p}) \rangle + \langle \mathbf{p}, \log \mathbf{u} \rangle < 0. \quad (57)$$

Since $\mathbf{p} \mapsto \underline{L}_{\ell}(\mathbf{p}) = -\sigma_{\mathcal{S}_{\ell}}(-\mathbf{p})$ is a closed concave function, the map $g: \mathbf{p} \mapsto \eta\underline{L}_{\ell}(\mathbf{p}) + \langle \mathbf{p}, \log \mathbf{u} \rangle$ is also closed and concave, and thus upper semi-continuous. Since Δ_n is compact, the function g must attain its maximum in Δ_n . Due to (57) this maximum is negative; there exists $c_1 > 0$ such that

$$\forall \mathbf{p} \in \Delta_n, \langle \mathbf{p}, \eta\ell(\mathbf{p}) \rangle - \langle \mathbf{p}, -\log \mathbf{u} \rangle \leq -c_1. \quad (58)$$

Let $f(\mathbf{p}, x) := \eta\underline{L}_{\ell}(\mathbf{p}) + \log u_x + c_1$, for $x \in [n]$. It follows from (58) that for all $\mathbf{p} \in \Delta_n$, $\mathbb{E}_{x \sim \mathbf{p}} f(\mathbf{p}, x) \leq 0$ and $\forall x \in [n], -\infty < f(\mathbf{p}, x)$. Thus, Lemma 6 applied to f with $\epsilon = c_1/2$, implies that there exists $\mathbf{p}_* \in \text{ri } \Delta_n$, such that $\eta\ell(\mathbf{p}_*) \leq -\log \mathbf{u} - c_1/2 \leq -\log \mathbf{u}$. From this inequality, $\mathbf{p}_* \in \text{dom } \ell$, and therefore, there exists $\mathbf{a}_* \in \text{dom } \ell$ such that $\ell(\mathbf{a}_*) = \ell(\mathbf{p}_*)$ (Theorem 5). This shows that $\eta\ell(\mathbf{a}_*) \leq -\log \mathbf{u}$, which implies that $\mathbf{u} \in \exp(-\eta\mathcal{S}_{\ell})$. Therefore, $\bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{\tau(\mathbf{p}),1} \cap]0, +\infty[^n \subseteq e^{-\eta\mathcal{S}_{\ell}}$. Combining this with (54) shows that $e^{-\eta\mathcal{S}_{\ell}} = \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{\tau(\mathbf{p}),1} \cap]0, +\infty[^n$. Since $e^{-\eta\mathcal{S}_{\ell}}$ is the intersection of convex set, it is a itself convex set. Since $\text{dom } \ell = \mathcal{A}$ by assumption, it follows that $\mathcal{S}_{\ell} = \mathcal{S}_{\ell}^{\infty}$, and thus $e^{-\eta\mathcal{S}_{\ell}^{\infty}}$ is convex. This last fact implies that ℓ is η -mixable [5]. \square

C.4 Proof of Theorem 10

We start by the following characterization of Δ -differentiability (this was defined on page 5 of the main body of the paper).

Lemma 19. *Let $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy. Then Φ is Δ -differentiable if and only if $\forall \mathbf{l} \subseteq [k]$ such that $|\mathbf{l}| > 1$, $\tilde{\Phi}_{\mathbf{l}} := \Phi \circ \Pi_k \circ [\Pi_k^{\tilde{\mathbf{l}}}]^{\top}$ is differentiable on $\text{int } \tilde{\Delta}_{|\mathbf{l}|}$.*

Proof. This is a direct consequence of Proposition B.4.2.1 in [7], since 1) $\tilde{\Phi}_l$ is convex; and 2)

$$\begin{aligned}\tilde{\Phi}'(\tilde{\mathbf{u}}; \tilde{\mathbf{v}} - \tilde{\mathbf{u}}) &= \tilde{\Phi}'([\Pi_l^{\tilde{k}}]^\top \tilde{\mathbf{u}}; [\Pi_l^{\tilde{k}}]^\top (\mathbf{v} - \mathbf{u})), \\ &= \Phi'(\Pi_k [\Pi_l^{\tilde{k}}]^\top \tilde{\mathbf{u}}; \Pi_k [\Pi_l^{\tilde{k}}]^\top (\mathbf{v} - \mathbf{u})),\end{aligned}$$

for all $\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \in \text{int } \tilde{\Delta}_{|l|}$ and $\tilde{\Phi} := \Phi \circ \Pi_k$. \square

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

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

where \mathbf{e}_θ is the θ th basis element of \mathbb{R}^k .

Proof. For all $l \subseteq [k]$ such that $|l| > 1$, let $\tilde{\Phi} := \Phi \circ \Pi_k$ and $\tilde{\Phi}_l := \tilde{\Phi} \circ [\Pi_l^{\tilde{k}}]^\top$. From Lemma 14 the infimum involved in the definition of the expert distribution \mathbf{q}^t in Algorithm 2 is indeed attained. It remains to verify that this minimum is unique. This will become clear in what follows.

Let $l^0 = [k]$ and $\mathcal{I}^t := \{\theta \in [k] : \ell_{x^t}(\mathbf{a}_\theta^t) < +\infty\}$, $t \in [T]$. For $t \in [T]$, we define the non-increasing sequence of subsets (l^t) of $[k]$ defined by $l^t := \mathcal{I}^t \cap l^{t-1}$. We show by induction that $\mathbf{q}^t \in \Delta_{l^t}$ and

$$\nabla \tilde{\Phi}_{l^t}(\Pi_{l^t}^{\tilde{k}} \tilde{\mathbf{q}}^t) = \Pi_{l^t}^{\tilde{k}} \left(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) - \sum_{s=1}^t J_k^\top \ell_{x^s}(A^s) \right), \quad (59)$$

where $A^s := [\mathbf{a}_\theta^s] \in \mathcal{A}^k$, $s \in \mathbb{N}$. Suppose that (59) holds true up to some $t \geq 1$. We will now show that it holds for $t+1$. To simplify expressions, we denote $\tilde{\mathbf{x}}_l := \Pi_l^{\tilde{k}} \tilde{\mathbf{x}} \in \mathbb{R}^l$ for $\tilde{\mathbf{x}} \in \mathbb{R}^{\tilde{k}}$, and $\mathbf{z}^t := \ell_{x^t}(A^t)$, $t \in [T]$. From the definition of \mathbf{q}^t in Algorithm 2, we have

$$\mathbf{q}^{t+1} \in \mathcal{M} := \underset{\mu \in \Delta_k}{\text{Argmin}} \langle \mu, \mathbf{z}^{t+1} \rangle + D_\Phi(\mu, \mathbf{q}^t).$$

Using the definition of \mathcal{I}^{t+1} ,

$$\begin{aligned}\mathcal{M} &= \underset{\mu \in \Delta_{l^{t+1}}}{\text{Argmin}} \langle \mu, \mathbf{z}^{t+1} \rangle + D_\Phi(\mu, \mathbf{q}^t), \\ &= \underset{\mu \in \Delta_{l^{t+1}}}{\text{Argmin}} \langle \mu, \mathbf{z}^{t+1} \rangle + \tilde{\Phi}_{l^t}(\tilde{\mu}_{l^t}) - \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \tilde{\Phi}'_{l^t}(\tilde{\mathbf{q}}_{l^t}^t; \tilde{\mu}_{l^t} - \tilde{\mathbf{q}}_{l^t}^t).\end{aligned}$$

Now using the facts that $\mathbf{q}^t \in \Delta_{l^t}$, $\mu \in \Delta_{l^{t+1}} \subseteq \Delta_{l^t}$, Φ is Δ -differentiable, and Lemma 19, we have

$$\mathcal{M} = \underset{\mu \in \Delta_{l^{t+1}}}{\text{Argmin}} \langle \mu, \mathbf{z}^{t+1} \rangle + \tilde{\Phi}_{l^{t+1}}(\tilde{\mu}_{l^{t+1}}) - \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \langle \tilde{\mu}_{l^t} - \tilde{\mathbf{q}}_{l^t}^t, \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) \rangle.$$

Using the facts that $\langle \mu, \mathbf{z}^{t+1} \rangle = z_k^{t+1} + \langle \tilde{\mu}_{l^{t+1}}, \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1} \rangle$, for $\tilde{\mu} \in \tilde{\Delta}_{l^{t+1}}$, and $\langle \tilde{\mu}_{l^t}, \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) \rangle \mathcal{M} = \langle \tilde{\mu}_{l^{t+1}}, \Pi_{l^{t+1}}^{\tilde{k}} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) \rangle$ (since $\mu \in \Delta_{l^{t+1}}$)

$$\begin{aligned}\mathcal{M} &= \underset{\mu \in \Delta_{l^{t+1}}}{\text{Argmin}} \langle \tilde{\mu}_{l^{t+1}}, -\Pi_{l^{t+1}}^{\tilde{k}} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) + \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1} \rangle + \tilde{\Phi}_{l^{t+1}}(\tilde{\mu}_{l^{t+1}}) \\ &\quad + \langle \tilde{\mathbf{q}}_{l^t}^t, \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) \rangle - \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t),\end{aligned}$$

and since the last two terms are independent of μ ,

$$\mathcal{M} = \underset{\mu \in \Delta_{l^{t+1}}}{\text{Argmin}} \langle \tilde{\mu}_{l^{t+1}}, -\Pi_{l^{t+1}}^{\tilde{k}} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) + \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1} \rangle + \tilde{\Phi}_{l^{t+1}}(\tilde{\mu}_{l^{t+1}}).$$

Now using Fenchel duality property in Proposition 1-(iv),

$$\mathcal{M} = \{\mu \in \Delta_{l^{t+1}} : \Pi_{l^{t+1}}^{\tilde{k}} \circ \Pi_k(\mu) = \tilde{\mu}_{l^{t+1}} \in \partial \tilde{\Phi}_{l^{t+1}}^*(\Pi_{l^{t+1}}^{\tilde{k}} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1})\}.$$

Finally, due to Lemma 10 and Proposition 12, $\tilde{\Phi}_{l^{t+1}}^*$ is differentiable on $\mathbb{R}^{|l^{t+1}|-1}$, and thus

$$\mathcal{M} = \{\Pi_k \circ [\Pi_{l^{t+1}}^{\tilde{k}}]^\top \circ \nabla \tilde{\Phi}_{l^{t+1}}^*(\Pi_{l^{t+1}}^{\tilde{k}} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1})\}. \quad (60)$$

From (60), we obtain

$$\nabla \tilde{\Phi}_{l^{t+1}}(\Pi_{l^{t+1}}^{\tilde{k}} \tilde{\mathbf{q}}^{t+1}) = \Pi_{l^{t+1}}^{l^t} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1}. \quad (61)$$

Thus using the induction assumption and the fact that $\Pi_{l^{t+1}}^{l^t} \Pi_{l^t}^{\tilde{k}} = \Pi_{l^{t+1}}^{\tilde{k}}$ (since $l^{t+1} \subseteq l^t$), the result follows, i.e. (59) is true for all $t \in [T]$. Furthermore, $\mathbf{q}^{t+1} \in \Delta_{l^{t+1}}$, since $\Pi_{l^{t+1}}^{\tilde{k}} \tilde{\mathbf{q}}^{t+1} \in \text{dom } \tilde{\Phi}_{l^{t+1}} \subseteq \tilde{\Delta}_{|l^{t+1}|}$. Using the same arguments as above, one arrives at

$$\begin{aligned} \text{Mix}_{\Phi}(\mathbf{q}^t, \mathbf{z}^{t+1}) &= z_k^{t+1} + \inf_{\mu \in \Delta_{l^{t+1}}} \langle \tilde{\mu}_{l^{t+1}}, -\Pi_{l^{t+1}}^{l^t} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) + \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1} \rangle + \tilde{\Phi}_{l^{t+1}}(\tilde{\mu}_{l^{t+1}}) \\ &\quad + \langle \tilde{\mathbf{q}}_{l^t}^t, \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) \rangle - \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t). \end{aligned}$$

Using the Fenchel duality property Proposition 1-(vi) and (60),

$$= z_k^{t+1} + \tilde{\Phi}_{l^t}^*(\nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t)) - \tilde{\Phi}_{l^{t+1}}^*(\Pi_{l^{t+1}}^{l^t} \nabla \tilde{\Phi}_{l^t}(\tilde{\mathbf{q}}_{l^t}^t) - \Pi_{l^{t+1}}^{\tilde{k}} J_k^\top \mathbf{z}^{t+1}). \quad (62)$$

On the other hand, Φ -mixability implies that there exists $\mathbf{a}_*^t \in \mathcal{A}^t$, such that for all $x^t \in [n]$,

$$\forall t \in [T], \ell_{x^t}(\mathbf{a}_*^t) \leq \text{Mix}_{\Phi}(\mathbf{q}^{t-1}, \mathbf{z}^t),$$

Summing this inequality for $t = 1, \dots, T$ yields,

$$\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) \leq \sum_{t=1}^T \text{Mix}_{\Phi}(\mathbf{q}^{t-1}, \mathbf{z}^t),$$

and thus using (62) and (61) yields

$$\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) \leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \tilde{\Phi}^*(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0)) - \tilde{\Phi}_{l^T}^*(\Pi_{l^T}^{l^{T-1}} \nabla \tilde{\Phi}_{l^{T-1}}(\tilde{\mathbf{q}}_{l^{T-1}}^{T-1}) - \Pi_{l^T}^{\tilde{k}} J_k^\top \mathbf{z}^T).$$

Finally, using (59) together with the fact that $\Pi_{l^T}^{l^{T-1}} \Pi_{l^{T-1}}^{\tilde{k}} = \Pi_{l^T}^{\tilde{k}}$

$$\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) \leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \tilde{\Phi}^*(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0)) - \tilde{\Phi}_{l^T}^* \left(\Pi_{l^T}^{\tilde{k}} \left(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) - \sum_{t=1}^T J_k^\top \ell_{x^t}(\mathbf{A}^t) \right) \right).$$

Using the definition of the Fenchel dual and Proposition 1-(vi) again, the above inequality becomes

$$\begin{aligned} \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) &\leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \langle \tilde{\mathbf{q}}^0, \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \rangle - \tilde{\Phi}(\tilde{\mathbf{q}}^0) \\ &\quad - \sup_{\pi \in \Delta_{|l^T|}} \left[\left\langle \tilde{\pi}, \Pi_{l^T}^{\tilde{k}} \left(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) - \sum_{t=1}^T J_k^\top \ell_{x^t}(\mathbf{A}^t) \right) \right\rangle - \tilde{\Phi}_{l^T}(\tilde{\pi}) \right], \\ &= \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \langle \tilde{\mathbf{q}}^0, \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \rangle - \tilde{\Phi}(\tilde{\mathbf{q}}^0) \\ &\quad + \inf_{\mu \in \Delta_{l^T}} \left[\left\langle \tilde{\mu}, \sum_{t=1}^T J_k^\top \ell_{x^t}(\mathbf{A}^t) - \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \right\rangle + \tilde{\Phi}(\tilde{\mu}) \right]. \quad (63) \end{aligned}$$

Using the fact that $\forall \theta \in [k] \setminus l^T$, $\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_\theta^t) = +\infty$ (by definition of (l^t)), the right hand side of (63) becomes

$$\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \langle \tilde{\mathbf{q}}^0, \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \rangle - \tilde{\Phi}(\tilde{\mathbf{q}}^0) + \inf_{\mu \in \Delta_k} \left[\left\langle \tilde{\mu}, \sum_{t=1}^T J_k^\top \ell_{x^t}(\mathbf{A}^t) - \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \right\rangle + \tilde{\Phi}(\tilde{\mu}) \right].$$

Thus, we get

$$\begin{aligned} \forall \mu \in \Delta_k, \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) &\leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \left\langle \tilde{\mu}, \sum_{t=1}^T J_k^\top \ell_{x^t}(\mathbf{A}^t) \right\rangle \\ &\quad + \tilde{\Phi}(\tilde{\mu}) - \tilde{\Phi}(\tilde{\mathbf{q}}^0) - \langle \tilde{\mu} - \tilde{\mathbf{q}}^0, \nabla \tilde{\Phi}(\tilde{\mathbf{q}}^0) \rangle. \end{aligned}$$

Using the facts that $\sum_{t=1}^T \ell_{x^t}(\mathbf{a}_k^t) + \left\langle \tilde{\boldsymbol{\mu}}, \sum_{t=1}^T J_k^\top \ell_{x^t}(A^t) \right\rangle = \left\langle \boldsymbol{\mu}, \sum_{t=1}^T \ell_{x^t}(A^t) \right\rangle$ and the definition of the divergence,

$$\forall \boldsymbol{\mu} \in \Delta_k, \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) \leq \left\langle \boldsymbol{\mu}, \sum_{t=1}^T \ell_{x^t}(A^t) \right\rangle + D_\Phi(\boldsymbol{\mu}, \mathbf{q}^0),$$

which for $\boldsymbol{\mu} = \mathbf{e}_\theta$ implies

$$\forall \theta \in [k], \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) \leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_\theta^t) + D_\Phi(\mathbf{e}_\theta, \mathbf{q}^0). \quad (64)$$

When instead of Φ -mixability, we have (η, Φ) -mixability, the last term in (64) becomes $\frac{D_\Phi(\mathbf{e}_\theta, \mathbf{q}^0)}{\eta}$ and the desired result follows. \square

C.5 Proof of Theorem 11

We require the following result:

Proposition 20. *For the Shannon entropy S , it holds that $\tilde{S}^*(\mathbf{v}) = \log(\langle \exp(\mathbf{v}), \mathbf{1}_{\bar{k}} \rangle + 1)$, $\forall \mathbf{v} \in \mathbb{R}^{k-1}$, and $S^*(\mathbf{z}) = \log \langle \exp(\mathbf{z}), \mathbf{1}_k \rangle$, $\forall \mathbf{z} \in \mathbb{R}^k$.*

Proof. Given $\mathbf{v} \in \mathbb{R}^{k-1}$, we first derive the expression of the Fenchel dual $\tilde{S}^*(\mathbf{v}) := \sup_{\tilde{\mathbf{q}} \in \tilde{\Delta}_k} \langle \tilde{\mathbf{q}}, \mathbf{v} \rangle - \tilde{S}(\tilde{\mathbf{q}})$. Setting the gradient of $\tilde{\mathbf{q}} \mapsto \langle \tilde{\mathbf{q}}, \mathbf{v} \rangle - \tilde{S}(\tilde{\mathbf{q}})$ to $\mathbf{0}_{\bar{k}}$ gives $\mathbf{v} = \nabla \tilde{S}(\tilde{\mathbf{q}})$. For $\mathbf{q} \in]0, +\infty[^k$, we have $\nabla S(\mathbf{q}) = \log \mathbf{q} + \mathbf{1}_k$, and from appendix A we know that $\nabla \tilde{S}(\tilde{\mathbf{q}}) = J_k^\top \nabla S(\mathbf{q})$. Therefore,

$$\mathbf{v} = \nabla \tilde{S}(\tilde{\mathbf{q}}) \implies \mathbf{v} = J_k^\top \nabla S(\mathbf{q}) \implies \mathbf{v} = \log \frac{\tilde{\mathbf{q}}}{q_k},$$

where the right most equality is equivalent to $\tilde{\mathbf{q}}/q_k = \exp(\mathbf{v})$. Since $\langle \tilde{\mathbf{q}}, \mathbf{1}_{\bar{k}} \rangle = 1 - q_k$, we get $q_k = (\langle \exp(\mathbf{v}), \mathbf{1}_{\bar{k}} \rangle + 1)^{-1}$. Therefore, the supremum in the definition of $\tilde{S}^*(\mathbf{v})$ is attained at $\tilde{\mathbf{q}}_* = \exp(\mathbf{v})(\langle \exp(\mathbf{v}), \mathbf{1}_{\bar{k}} \rangle + 1)^{-1}$. Hence $\tilde{S}^*(\mathbf{v}) = \langle \tilde{\mathbf{q}}_*, \mathbf{v} \rangle - \langle \tilde{\mathbf{q}}_*, \log \tilde{\mathbf{q}}_* \rangle = \log(\langle \exp(\mathbf{v}), \mathbf{1}_{\bar{k}} \rangle + 1)$. Finally, using (2) we get $S^*(\mathbf{z}) = \log \langle \exp(\mathbf{z}), \mathbf{1}_k \rangle$, for $\mathbf{z} \in \mathbb{R}^k$. \square

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

Proof.

Claim 1. *For all $\mathbf{q} \in \Delta_k$, $A := \mathbf{a}_{1:k} \in \mathbb{R}^k$, and $x \in [n]$*

$$-\eta^{-1} \log \langle \exp(-\eta \ell_x(A)), \mathbf{q} \rangle = \text{Mix}_S^\eta(\ell_x(A), \mathbf{q}). \quad (65)$$

Let $\mathbf{q} \in \text{ri } \Delta_k$. From Proposition 20, the Shannon entropy is such that S^ is differentiable on \mathbb{R}^k , and thus it follows from Lemma 14 ((21)-(22)) that for any $\mathbf{d} \in [0, +\infty[^k$*

$$\text{Mix}_S(\mathbf{d}, \mathbf{q}) = S^*(\nabla S(\mathbf{q})) - S^*(\nabla S(\mathbf{q}) - \mathbf{d}). \quad (66)$$

By definition of S , $\nabla S(\mathbf{q}) = \log \mathbf{q} + \mathbf{1}_k$, and due to Proposition 20, $S^(\mathbf{z}) = \log \langle \exp \mathbf{z}, \mathbf{1}_k \rangle$, $\mathbf{z} \in \mathbb{R}^k$. Therefore,*

$$\nabla S(\mathbf{q}) - \eta \mathbf{d} = \log(\exp(-\eta \mathbf{d}) \odot \mathbf{q}) + \mathbf{1}_k. \quad (67)$$

On the other hand, from [9] we also have

$$\text{Mix}_S^\eta(\mathbf{d}, \mathbf{q}) = \eta^{-1} \text{Mix}_S(\eta \mathbf{d}, \mathbf{q}), \quad \eta > 0. \quad (68)$$

Combining (66)-(68), yields

$$-\eta^{-1} \log \langle \exp(-\eta \mathbf{d}), \mathbf{q} \rangle = \text{Mix}_S^\eta(\mathbf{d}, \mathbf{q}). \quad (69)$$

Suppose now that $\mathbf{q} \in \text{ri } \Delta_{\mathbf{l}}$ for $\mathbf{l} \subseteq [k]$ such that $|\mathbf{l}| > 1$. By repeating the argument above for $S_{\mathbf{l}} := S \circ \Pi_{\mathbf{l}}^T$, we get

$$\begin{aligned} \forall \mathbf{d} \in [0, +\infty]^n, \text{Mix}_{S_{\mathbf{l}}}^{\eta}(\Pi_{\mathbf{l}}\mathbf{d}, \Pi_{\mathbf{l}}\mathbf{q}) &= -\eta^{-1} \log \langle \exp(-\eta \Pi_{\mathbf{l}}\mathbf{d}), \Pi_{\mathbf{l}}\mathbf{q} \rangle, \\ &= -\eta^{-1} \log \langle \exp(-\eta \mathbf{d}), \mathbf{q} \rangle. \end{aligned} \quad (70)$$

Fix $x \in [n]$ and let $\hat{\mathbf{d}} := \ell_x(A) \in [0, +\infty]^k$. Let $(\hat{\mathbf{d}}_m) \subset [0, +\infty]^k$ be any sequence converging to $\hat{\mathbf{d}}$. Lemma 15, $\text{Mix}_S^{\eta}(\hat{\mathbf{d}}_m, \mathbf{q}) \xrightarrow{m \rightarrow \infty} \text{Mix}_S^{\eta}(\hat{\mathbf{d}}, \mathbf{q})$. Using this with (70) gives

$$\begin{aligned} -\eta^{-1} \log \langle \exp(-\eta \ell_x(A)), \mathbf{q} \rangle &= \lim_{m \rightarrow \infty} -\eta^{-1} \log \langle \exp(-\eta \hat{\mathbf{d}}_m), \mathbf{q} \rangle, \\ &= \lim_{m \rightarrow \infty} \text{Mix}_S^{\eta}(\hat{\mathbf{d}}_m, \mathbf{q}), \\ &= \text{Mix}_S^{\eta}(\hat{\mathbf{d}}, \mathbf{q}) = \text{Mix}_S^{\eta}(\ell_x(A), \mathbf{q}). \end{aligned} \quad (71)$$

It remains to check the case where \mathbf{q} is a vertex; Without loss of generality assume that $\mathbf{q} = \mathbf{e}_1$ and let $\mu \in \Delta_k \setminus \{\mathbf{e}_1\}$. Then there exists $\mathbf{l}_* \subset [k]$, such that $(\mathbf{e}_1, \mu) \in (\text{rbd } \Delta_{\mathbf{l}_*}) \times (\text{ri } \Delta_{\mathbf{l}_*})$ and by Lemma 11, $S'(\mathbf{e}_1; \mu - \mathbf{e}_1) = -\infty$. Therefore, $\forall \mathbf{q} \in \Delta_k \setminus \{\mathbf{e}_1\}$, $D_{S_{\eta}}(\mathbf{q}, \mathbf{e}_1) = +\infty$, which implies

$$\begin{aligned} \forall x \in [n], \text{Mix}_S^{\eta}(\ell_x(A), \mathbf{e}_1) &= \inf_{\mathbf{q} \in \Delta_k} \langle \mathbf{q}, \ell_x(A) \rangle + D_{S_{\eta}}(\mathbf{q}, \mathbf{e}_1), \\ &= \langle \mathbf{e}_1, \ell_x(A) \rangle + D_{S_{\eta}}(\mathbf{e}_1, \mathbf{e}_1), \\ &= \langle \mathbf{e}_1, \ell_x(A) \rangle, \\ &= \ell_x(\mathbf{a}_1) = -\eta^{-1} \log \langle \exp(-\eta \ell_x(A)), \mathbf{e}_1 \rangle. \end{aligned} \quad (72)$$

Combining (72) and (71) proves the claim in (65). The desired equivalence follows trivially from the definitions of η -mixability and (η, S) -mixability. \square

C.6 Proof of Theorem 13

We need the following lemma to show Theorem 13.

Lemma 21. *Let Φ be as in Theorem 13. Then $\eta_{\ell}\Phi - S$ is convex on Δ_k only if Φ satisfies (10).*

Proof. Let $\hat{\mathbf{q}} \in \text{rbd } \Delta_k$. Suppose that there exists $\mathbf{q} \in \text{ri } \Delta_k$ such that $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) > -\infty$. Since Φ is convex, it must have non-decreasing slopes; in particular, it holds that $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) \leq \Phi(\mathbf{q}) - \Phi(\hat{\mathbf{q}})$. Therefore, since Φ is finite on Δ_k (by definition of an entropy), we have $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) < +\infty$. Since by assumption $\eta_{\ell}\Phi - S$ is convex and finite on the simplex, we can use the same argument to show that $[\eta_{\ell}\Phi - S]'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) = \eta_{\ell}\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) - S'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) < +\infty$. This is a contradiction since $S'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) = -\infty$ (Lemma 11). Therefore, it must hold that $\Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}}) = -\infty$.

Suppose now that $(\hat{\mathbf{q}}, \mathbf{q}) \in (\text{rbd } \Delta_{\mathbf{l}}) \times (\text{ri } \Delta_{\mathbf{l}})$ for $\mathbf{l} \subseteq [k]$, with $|\mathbf{l}| > 1$. Let $\Phi_{\mathbf{l}} := \Phi \circ \Pi_{\mathbf{l}}^T$ and $S_{\mathbf{l}} := S \circ \Pi_{\mathbf{l}}^T$. Since $\eta_{\ell}\Phi - S$ is convex on Δ_k and $\Pi_{\mathbf{l}}$ is a linear function, $\eta_{\ell}\Phi_{\mathbf{l}} - S_{\mathbf{l}}$ is convex on $\Delta_{|\mathbf{l}|}$. Repeating the steps above for Φ and S substituted by $\Phi_{\mathbf{l}}$ and $S_{\mathbf{l}}$, respectively, we get that $(\Phi_{\mathbf{l}})'(\Pi_{\mathbf{l}}\hat{\mathbf{q}}; \Pi_{\mathbf{l}}\mathbf{q} - \Pi_{\mathbf{l}}\hat{\mathbf{q}}) = -\infty$. Since $(\Phi_{\mathbf{l}})'(\Pi_{\mathbf{l}}\hat{\mathbf{q}}; \Pi_{\mathbf{l}}\mathbf{q} - \Pi_{\mathbf{l}}\hat{\mathbf{q}}) = \Phi'(\hat{\mathbf{q}}; \mathbf{q} - \hat{\mathbf{q}})$ the proof is completed. \square

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

Proof. Assume $\eta_{\ell}\Phi - S$ is convex on Δ_k . For this to hold, it is necessary that $\eta_{\ell} > 0$ since $-S$ is strictly concave. Let $\eta := \eta_{\ell}$ and $S_{\eta} := \eta^{-1}S$. Then $\tilde{S}_{\eta} = \eta^{-1}\tilde{S}$ and $\tilde{\Phi} - \tilde{S}_{\eta} = (\Phi - S_{\eta}) \circ \Pi_k$ is convex on $\tilde{\Delta}_k$, since $\Phi - S_{\eta}$ is convex on Δ_k and Π_k is affine.

Let $x \in [n]$, $A := [\mathbf{a}_{\theta}]_{\theta \in [k]}$, and $\mathbf{q} \in \Delta_k$. Suppose that $\mathbf{q} \in \text{ri } \Delta_k$ and let $\mathbf{s}_{\mathbf{q}}^* \in \partial \tilde{\Phi}(\tilde{\mathbf{q}})$ be as in Proposition 14. Note that if $\ell_x(\mathbf{a}_{\theta}) = +\infty, \forall \theta \in [k]$, then the Φ -mixability condition (8) is trivially satisfied. Suppose, without loss of generality, that $\ell_x(\mathbf{a}_k) < +\infty$. Let $(\mathbf{d}_m) \subset$

$[0, +\infty]^k$ be any sequence such that $\mathbf{d}_m \xrightarrow{m \rightarrow \infty} \mathbf{d} := \ell_x(A) \in [0, +\infty]^k$. From Lemmas 11 and 15, $\text{Mix}_\Psi(\mathbf{d}_m, \mathbf{q}) \xrightarrow{m \rightarrow \infty} \text{Mix}_\Psi(\mathbf{d}, \mathbf{q})$ for $\Psi \in \{\Phi, S_\eta\}$.

Let $\tilde{\Upsilon}_q : \mathbb{R}^{k-1} \rightarrow \mathbb{R} \cup \{+\infty\}$ be defined by

$$\tilde{\Upsilon}_q(\tilde{\boldsymbol{\mu}}) := \tilde{S}_\eta(\tilde{\boldsymbol{\mu}}) + \langle \tilde{\boldsymbol{\mu}}, \mathbf{s}_q^* - \nabla \tilde{S}_\eta(\tilde{\mathbf{q}}) \rangle - \tilde{\Phi}^*(\mathbf{s}_q^*) + \tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}})),$$

and it's Fenchel dual follows from Proposition 1 (i+ii):

$$\tilde{\Upsilon}_q^*(\mathbf{v}) = \tilde{S}_\eta^*(\mathbf{v} - \mathbf{s}_q^* + \nabla \tilde{S}_\eta(\tilde{\mathbf{q}})) + \tilde{\Phi}^*(\mathbf{s}_q^*) - \tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}})),$$

After substituting \mathbf{v} by $\mathbf{s}_q^* - J_k^\top \mathbf{d}$ in the expression of $\tilde{\Upsilon}_q^*$ and rearranging, we get

$$\tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}})) - \tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}}) - J_k^\top \mathbf{d}_m) = \tilde{\Phi}^*(\mathbf{s}_q^*) - \tilde{\Upsilon}_q^*(\mathbf{s}_q^* - J_k^\top \mathbf{d}_m). \quad (73)$$

Since $\mathbf{s}_q^* \in \partial \tilde{\Phi}(\tilde{\mathbf{q}})$ and $\tilde{\Phi}$ is a closed convex function, combining Proposition 1-(iv) and the fact that $\tilde{\Phi}^{**} = \tilde{\Phi}$ [7, Cor. E.1.3.6] yields $\langle \tilde{\mathbf{q}}, \mathbf{s}_q^* \rangle - \tilde{\Phi}^*(\mathbf{s}_q^*) = \tilde{\Phi}(\tilde{\mathbf{q}})$. Thus, after substituting $\tilde{\boldsymbol{\mu}}$ by $\tilde{\mathbf{q}}$ in the expression of $\tilde{\Upsilon}_q$, we get

$$\tilde{\Phi}(\tilde{\mathbf{q}}) = \tilde{\Upsilon}_q(\tilde{\mathbf{q}}). \quad (74)$$

On the other hand, $\tilde{\Phi} - \tilde{\Upsilon}_q$ is convex on $\tilde{\Delta}_k$, since $\tilde{\Upsilon}_q$ is equal to \tilde{S}_η plus an affine function. Thus, $\partial[\tilde{\Phi} - \tilde{\Upsilon}_q](\tilde{\mathbf{q}}) + \partial \tilde{\Upsilon}_q(\tilde{\mathbf{q}}) = \partial \tilde{\Phi}(\tilde{\mathbf{q}})$, since $\tilde{\Phi}$ and $\tilde{\Upsilon}_q$ are both convex (ibid., Thm. D.4.1.1). Since $\tilde{\Upsilon}_q$ is differentiable at $\tilde{\mathbf{q}}$, we have $\partial \tilde{\Upsilon}_q(\tilde{\mathbf{q}}) = \{\nabla \tilde{\Upsilon}_q(\tilde{\mathbf{q}})\} = \{\mathbf{s}_q^*\}$. Furthermore, since $\mathbf{s}_q^* \in \partial \tilde{\Phi}(\tilde{\mathbf{q}})$, then $\mathbf{0}_{\tilde{k}} \in \partial \tilde{\Phi}(\tilde{\mathbf{q}}) - \partial \tilde{\Upsilon}_q(\tilde{\mathbf{q}}) = \partial[\tilde{\Phi} - \tilde{\Upsilon}_q](\tilde{\mathbf{q}})$. Hence, $\tilde{\Phi} - \tilde{\Upsilon}_q$ attains a minimum at $\tilde{\mathbf{q}}$ (ibid., Thm. D.2.2.1). Due to this and (74), $\tilde{\Phi} \geq \tilde{\Upsilon}_q$, which implies that $\tilde{\Phi}^* \leq \tilde{\Upsilon}_q^*$ (Proposition 1-(iii)). Using this in (73) gives for all $m \in \mathbb{N}$

$$\begin{aligned} \tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}})) - \tilde{S}_\eta^*(\nabla \tilde{S}_\eta(\tilde{\mathbf{q}}) - J_k^\top \mathbf{d}_m) &\leq \tilde{\Phi}^*(\mathbf{s}_q^*) - \tilde{\Phi}^*(\mathbf{s}_q^* - J_k^\top \mathbf{d}_m), \\ \implies \text{Mix}_{S_\eta}^\eta(\mathbf{d}_m, \mathbf{q}) &\leq \text{Mix}_\Phi(\mathbf{d}_m, \mathbf{q}), \end{aligned}$$

where the implication is obtained by adding $[\mathbf{d}_m]_k$ on both sides of the first inequality and using Proposition 14.

Suppose now that $\mathbf{q} \in \text{ri } \Delta_\ell$, with $|\ell| > 1$, and let $\Phi_\ell := \Phi \circ \Pi_\ell^\top$ and $S_\ell := S \circ \Pi_\ell^\top$. Note that since $\eta_\ell \Phi - S$ is convex on Δ_k and Π_ℓ is a linear function, $\eta_\ell \Phi_\ell - S_\ell$ is convex on $\Delta_{|\ell|}$. Repeating the steps above for Φ, S, \mathbf{q} , and A substituted by $\Phi_\ell, S_\ell, \Pi_\ell \mathbf{q}$, and $A \Pi_\ell^\top$, respectively, yields

$$\begin{aligned} \text{Mix}_{S_\ell}^\eta(\Pi_\ell \mathbf{d}_m, \Pi_\ell \mathbf{q}) &\leq \text{Mix}_{\Phi_\ell}(\Pi_\ell \mathbf{d}_m, \Pi_\ell \mathbf{q}), \\ \implies \text{Mix}_S^\eta(\mathbf{d}_m, \mathbf{q}) &\leq \text{Mix}_\Phi(\mathbf{d}_m, \mathbf{q}), \\ \implies \text{Mix}_S^\eta(\ell_x(A), \mathbf{q}) &\leq \text{Mix}_\Phi(\ell_x(A), \mathbf{q}), \end{aligned} \quad (75)$$

where the first implication follows from Lemma 13, since S_η and Φ both satisfy (10) (see Lemmas 11 and 21), and (75) is obtained by passage to the limit $m \rightarrow \infty$. Since $\eta = \eta_\ell > 0$, ℓ is η -mixable, which implies that ℓ is S_η -mixable (Theorem 11). Therefore, there exists $\mathbf{a}_* \in \mathcal{A}$, such that

$$\ell_x(\mathbf{a}_*) \leq \text{Mix}_S^\eta(\ell_x(A), \mathbf{q}) \leq \text{Mix}_\Phi(\ell_x(A), \mathbf{q}). \quad (76)$$

To complete the proof (that is, to show that ℓ is Φ -mixable), it remains to consider the case where \mathbf{q} is a vertex of Δ_k . Without loss of generality assume that $\mathbf{q} = \mathbf{e}_1$ and let $\boldsymbol{\mu} \in \Delta_k \setminus \{\mathbf{e}_1\}$. Thus, there exists $\ell_* \subseteq [k]$, with $|\ell_*| > 1$, such that $(\mathbf{e}_1, \boldsymbol{\mu}) \in (\text{rbd } \Delta_{\ell_*}) \times (\text{ri } \Delta_{\ell_*})$, and Lemma 21 implies that $\Phi'(\mathbf{e}_1; \boldsymbol{\mu} - \mathbf{e}_1) = -\infty$. Therefore, $\forall \mathbf{q} \in \Delta_k \setminus \{\mathbf{e}_1\}$, $D_\Phi(\mathbf{q}, \mathbf{e}_1) = +\infty$, which implies

$$\begin{aligned} \forall x \in [n], \text{Mix}_\Phi(\ell_x(A), \mathbf{e}_1) &= \inf_{\mathbf{q} \in \Delta_k} \langle \mathbf{q}, \ell_x(A) \rangle + D_\Phi(\mathbf{q}, \mathbf{e}_1), \\ &= \langle \mathbf{e}_1, \ell_x(A) \rangle + D_\Phi(\mathbf{e}_1, \mathbf{e}_1) = \langle \mathbf{e}_1, \ell_x(A) \rangle, \\ &= \ell_x(\mathbf{a}_1). \end{aligned} \quad (77)$$

The Φ -mixability condition (8) is trivially satisfied in this case. Combining (76) and (77) shows that ℓ is Φ -mixable. \square

C.7 Proof of Theorem 14

The following Lemma gives necessary regularity conditions on the entropy Φ under the assumptions of Theorem 14.

Lemma 22. *Let Φ and ℓ be as in Theorem 14. Then the following holds*

- (i) $\tilde{\Phi}$ is strictly concave on $\text{int } \tilde{\Delta}_k$.
- (ii) $\tilde{\Phi}^*$ is continuously differentiable on \mathbb{R}^{k-1} .
- (iii) $\tilde{\Phi}^*$ is twice differentiable on \mathbb{R}^{k-1} and $\forall \tilde{\mathbf{q}} \in \text{int } \tilde{\Delta}_k$, $\mathbf{H}\tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}})) = (\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}))^{-1}$.
- (iv) For the Shannon entropy, we have $(\mathbf{H}\tilde{S}(\tilde{\mathbf{q}}))^{-1} = \mathbf{H}\tilde{S}^*(\nabla\tilde{S}(\tilde{\mathbf{q}})) = \text{diag } \tilde{\mathbf{q}} - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top$.

Proof. Since ℓ is Φ -mixable and \underline{L}_ℓ is twice differentiable on $]0, +\infty[^n$, $\tilde{\Phi}^*$ is continuously differentiable on \mathbb{R}^{n-1} (Proposition 12). Therefore, $\tilde{\Phi}$ is strictly convex on $\text{ri } \tilde{\Delta}_k$ [7, Thm. E.4.1.2].

The differentiability of $\tilde{\Phi}$ and $\tilde{\Phi}^*$ implies $\nabla\tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}})) = \tilde{\mathbf{q}}$ (ibid.). Since $\tilde{\Phi}$ is twice differentiable on $\text{int } \tilde{\Delta}_k$ (by assumption), the latter equation implies that $\tilde{\Phi}^*$ is twice differentiable on $\nabla\tilde{\Phi}(\text{int } \tilde{\Delta}_k)$. Using the chain rule, we get $\mathbf{H}\tilde{\Phi}^*(\nabla\tilde{\Phi}(\mathbf{u}))\mathbf{H}\tilde{\Phi}(\mathbf{u}) = I_{\tilde{k}}$. Multiplying both sides of the equation by $(\mathbf{H}\tilde{\Phi}(\mathbf{u}))^{-1}$ from the right gives the expression in (iii). Note that $\mathbf{H}\tilde{\Phi}(\cdot)$ is in fact invertible on $\text{int } \tilde{\Delta}_k$ since $\tilde{\Phi}$ is strictly convex on $\text{int } \tilde{\Delta}_k$. It remains to show that $\nabla\tilde{\Phi}(\text{int } \tilde{\Delta}_k) = \mathbb{R}^{k-1}$. This set equality follows from 1) $[\tilde{\mathbf{q}} \in \partial\tilde{\Phi}^*(s) \iff s \in \partial\tilde{\Phi}(\tilde{\mathbf{q}})]$ (ibid., Cor. E.1.4.4); 2) $\text{dom } \tilde{\Phi}^* = \mathbb{R}^{k-1}$; and 3) $\forall \tilde{\mathbf{q}} \in \text{bd } \tilde{\Delta}_k, \partial\tilde{\Phi}(\tilde{\mathbf{q}}) = \emptyset$ (Lemma 13).

For the Shannon entropy, we have $\tilde{S}^*(\mathbf{v}) = \log(\langle \exp(\mathbf{v}), \mathbf{1}_{\tilde{k}} \rangle + 1)$ (Proposition 20) and $\nabla\tilde{S}(\tilde{\mathbf{q}}) = \log \frac{\tilde{\mathbf{q}}}{q_k}$, for $(\mathbf{v}, \tilde{\mathbf{q}}) \in \mathbb{R}^{k-1} \times \tilde{\Delta}_k$. Thus $(\mathbf{H}\tilde{S}(\tilde{\mathbf{q}}))^{-1} = \mathbf{H}\tilde{S}^*(\nabla\tilde{S}(\tilde{\mathbf{q}})) = \text{diag } \tilde{\mathbf{q}} - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top$. \square

To show Theorem 14, we analyze a particular parameterized curve defined in the next lemma.

Lemma 23. *Let $\ell: \Delta_n \rightarrow [0, +\infty]^n$ be a proper loss whose Bayes risk \underline{L}_ℓ is twice differentiable on $]0, +\infty[^n$, and let Φ be an entropy such that $\tilde{\Phi}$ and $\tilde{\Phi}^*$ are twice differentiable on $\text{int } \tilde{\Delta}_k$ and \mathbb{R}^{k-1} , respectively. For $(\tilde{\mathbf{p}}, \tilde{\mathbf{q}}, V) \in \text{int } \tilde{\Delta}_n \times \text{int } \tilde{\Delta}_k \times \mathbb{R}^{\tilde{n} \times \tilde{k}}$, let $\beta: \mathbb{R} \rightarrow \mathbb{R}^n$ be the curve defined by*

$$\forall x \in [n], \quad \beta_x(t) = \tilde{\ell}_x(\tilde{\mathbf{p}}) + \tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}})) - \tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}}) - J_k^\top \tilde{\ell}_x(\tilde{P}^t)), \quad (78)$$

where $\tilde{P}^t = [\tilde{\mathbf{p}}\mathbf{1}_k^\top + tV, \tilde{\mathbf{p}}] \in \mathbb{R}^{\tilde{n} \times k}$ and $t \in \{s \in \mathbb{R} : \forall j \in [\tilde{k}], \tilde{\mathbf{p}} + sV_{\cdot,j} \in \text{int } \tilde{\Delta}_n\}$. Then

$$\begin{aligned} \beta(0) &= \tilde{\ell}(\tilde{\mathbf{p}}), \\ \dot{\beta}(0) &= \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}})V\tilde{\mathbf{q}}, \\ \frac{d}{dt} \left\langle \mathbf{p}, \dot{\beta}(t) \right\rangle \Big|_{t=0} &= - \sum_{j=1}^{k-1} q_j V_{\cdot,j}^\top \mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})V_{\cdot,j} - \text{tr}(\text{diag } (\mathbf{p}) \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}})V(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}))^{-1}(\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}})V)^\top). \end{aligned} \quad (79)$$

Proof. Since $\tilde{P}^t = [\tilde{\mathbf{p}}\mathbf{1}_k^\top + tV, \tilde{\mathbf{p}}] \in \mathbb{R}^{\tilde{n} \times k}$, $\tilde{P}^0 = \tilde{\mathbf{p}}\mathbf{1}_k^\top$ and $\tilde{\ell}_x(\tilde{P}^0) = \tilde{\ell}_x(\tilde{\mathbf{p}})\mathbf{1}_k$. As a result, $J_k^\top \tilde{\ell}_x(\tilde{P}^0) = \mathbf{0}_{\tilde{k}}$, and thus $\beta_x(0) = \tilde{\ell}_x(\tilde{\mathbf{p}}) + \tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}})) - \tilde{\Phi}^*(\nabla\tilde{\Phi}(\tilde{\mathbf{q}}) - \mathbf{0}_{\tilde{k}}) = \tilde{\ell}_x(\tilde{\mathbf{p}})$. This shows that $\beta(0) = \tilde{\ell}(\tilde{\mathbf{p}})$. Let $\gamma_x(t) := \nabla\tilde{\Phi}(\tilde{\mathbf{q}}) - J_k^\top \tilde{\ell}_x(\tilde{P}^t)$. For $j \in [k-1]$,

$$\begin{aligned} \frac{d}{dt} [\gamma_x(t)]_j &= \frac{d}{dt} \left([\nabla\tilde{\Phi}(\tilde{\mathbf{q}})]_j - [J_k^\top \tilde{\ell}_x(\tilde{P}^t)]_j \right), \\ &= - \frac{d}{dt} \left(\tilde{\ell}_x(\tilde{P}_{\cdot,j}^t) - \tilde{\ell}_x(\tilde{P}_{\cdot,k}^t) \right), \\ &= - \frac{d}{dt} \left(\tilde{\ell}_x(\tilde{\mathbf{p}} + tV_{\cdot,j}) - \tilde{\ell}_x(\tilde{\mathbf{p}}) \right), \quad \left(\text{since } \frac{d}{dt} \tilde{\ell}_x(\tilde{P}_{\cdot,k}^t) = \frac{d}{dt} \tilde{\ell}_x(\tilde{\mathbf{p}}) = 0 \right) \\ &= -\mathbf{D}\tilde{\ell}_x(\tilde{P}_{\cdot,j}^t)V_{\cdot,j}. \end{aligned}$$

From the definition of \tilde{P}^t , $\tilde{P}_{\cdot,j}^0 = \tilde{\mathbf{p}}$, $\forall j \in [\tilde{k}]$, and therefore, $\dot{\gamma}_x(0) = -(\mathbf{D}\tilde{\ell}_x(\tilde{\mathbf{p}})V)^\top$. By differentiating β_x in (78) and using the chain rule, $\dot{\beta}_x(t) = -(\dot{\gamma}_x(t))^\top \nabla \tilde{\Phi}^*(\gamma_x(t))$. By setting $t = 0$, $\dot{\beta}_x(0) = -(\dot{\gamma}_x(0))^\top \nabla \tilde{\Phi}^*(\nabla \tilde{\Phi}(\tilde{\mathbf{q}})) = \mathbf{D}\tilde{\ell}_x(\tilde{\mathbf{p}})V\tilde{\mathbf{q}}$. Thus, $\dot{\beta}(0) = \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}})V\tilde{\mathbf{q}}$. Furthermore,

$$\begin{aligned}
\left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) \rangle \right|_{t=0} &= \left. \frac{d}{dt} \sum_{x=1}^n p_x \left(\sum_{j=1}^{k-1} \mathbf{D}\tilde{\ell}_x(\tilde{P}_{\cdot,j}^t) V_{\cdot,j} [\nabla \tilde{\Phi}^*(\gamma_x(t))]_j \right) \right|_{t=0}, \\
&= \sum_{j=1}^{k-1} \left. \frac{d}{dt} \left(\sum_{x=1}^n p_x \mathbf{D}\tilde{\ell}_x(\tilde{P}_{\cdot,j}^t) V_{\cdot,j} [\nabla \tilde{\Phi}^*(\gamma_x(t))]_j \right) \right|_{t=0}, \\
&= \sum_{j=1}^{k-1} \left(\left. \frac{d}{dt} \langle \mathbf{p}, \mathbf{D}\tilde{\ell}(\tilde{P}_{\cdot,j}^t) V_{\cdot,j} q_j \rangle \right|_{t=0} + \sum_{x=1}^n p_x \mathbf{D}\tilde{\ell}_x(\tilde{\mathbf{p}}) V_{\cdot,j} \left. \frac{d}{dt} [\nabla \tilde{\Phi}^*(\gamma_x(t))]_j \right|_{t=0} \right), \\
&= - \sum_{j=1}^{k-1} q_j V_{\cdot,j}^\top \mathbf{H}\tilde{\ell}_\ell(\tilde{\mathbf{p}}) V_{\cdot,j} - \sum_{x=1}^n \sum_{i=1}^{k-1} p_x \mathbf{D}\tilde{\ell}_x(\tilde{\mathbf{p}}) V_{\cdot,j} [\mathbf{H}\tilde{\Phi}^*(\nabla \tilde{\Phi}(\tilde{\mathbf{q}}))]_{j,i} \mathbf{D}\tilde{\ell}_x(\tilde{\mathbf{p}}) V_{\cdot,i}, \\
&= - \sum_{j=1}^{k-1} q_j V_{\cdot,j}^\top \mathbf{H}\tilde{\ell}_\ell(\tilde{\mathbf{p}}) V_{\cdot,j} - \text{tr}(\text{diag}(\mathbf{p}) \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) V \mathbf{H}\tilde{\Phi}^*(\nabla \tilde{\Phi}(\tilde{\mathbf{q}})) (\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) V)^\top), \\
&= - \sum_{j=1}^{k-1} q_j V_{\cdot,j}^\top \mathbf{H}\tilde{\ell}_\ell(\tilde{\mathbf{p}}) V_{\cdot,j} - \text{tr}(\text{diag}(\mathbf{p}) \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) V (\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}))^{-1} (\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) V)^\top),
\end{aligned}$$

where in the third equality we used Lemma 6, in the fourth equality we used Lemma 9, and in the sixth equality we used Lemma 22-(iii). \square

In next lemma, we state a necessary condition for Φ -mixability in terms of the parameterized curve β defined in Lemma 23.

Lemma 24. *Let ℓ , Φ , and β be as in Lemma 23. If $\exists(\tilde{\mathbf{p}}, \tilde{\mathbf{q}}, V) \in \text{int } \tilde{\Delta}_n \times \text{int } \tilde{\Delta}_k \times \mathbb{R}^{\tilde{n} \times \tilde{k}}$ such that the curve $\gamma(t) := \tilde{\ell}(\tilde{\mathbf{p}} + tV\tilde{\mathbf{q}})$ satisfies $\left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0} < 0$, then ℓ is not Φ -mixable. In particular, $\exists P \in \text{ri } \Delta_n^k$, such that $[\text{Mix}_\Phi(\ell_x(P), \mathbf{q})]_{x \in [n]}^\top$ lies outside \mathcal{S}_ℓ .*

Proof. First note that for any triplet $(\tilde{\mathbf{p}}, \tilde{\mathbf{q}}, V) \in \text{int } \tilde{\Delta}_n \times \text{int } \tilde{\Delta}_k \times \mathbb{R}^{\tilde{n} \times \tilde{k}}$, the map $t \mapsto \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle$ is differentiable at 0. This follows from Lemmas 6 and 23. Let $r(t) := \Pi_n(\tilde{\mathbf{p}} + tV\tilde{\mathbf{q}})$ and $\delta(t) := \langle r(t), \beta(t) - \gamma(t) \rangle$. Then

$$\dot{\delta}(t) = \langle r(t), \dot{\beta}(t) - \dot{\gamma}(t) \rangle + \langle V\tilde{\mathbf{q}}, \beta(t) - \gamma(t) \rangle.$$

Since $t \mapsto \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle$ is differentiable at 0, it follows from Lemma 6 that $t \mapsto \dot{\delta}(t)$ is also differentiable at 0, and thus

$$\begin{aligned}
\ddot{\delta}(0) &= \left. \frac{d}{dt} \langle r(t), \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0} + \langle J_n V \tilde{\mathbf{q}}, \dot{\beta}(0) - \dot{\gamma}(0) \rangle, \\
&= \langle \dot{r}(0), \dot{\beta}(0) - \dot{\gamma}(0) \rangle + \left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0}, \tag{80}
\end{aligned}$$

$$\begin{aligned}
&= \langle J_n V \tilde{\mathbf{q}}, \dot{\beta}(0) - \dot{\gamma}(0) \rangle + \left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0}, \\
&= \left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0} < 0, \tag{81}
\end{aligned}$$

where (80) and (81) hold because $\dot{\beta}(0) = D\tilde{\ell}(\tilde{\mathbf{p}})V\tilde{\mathbf{q}} = \dot{\gamma}(0)$ (see Lemma 23). According to Taylor's theorem (see e.g. [?, §151]), there exists $\epsilon > 0$ and $h : [-\epsilon, \epsilon] \rightarrow \mathbb{R}$ such that

$$\forall |t| \leq \epsilon, \delta(t) = \delta(0) + t\dot{\delta}(0) + \frac{t^2}{2}\ddot{\delta}(0) + h(t)t^2, \quad (82)$$

and $\lim_{t \rightarrow 0} h(t) = 0$. From Lemma 23, $\beta(0) = \gamma(0) = 0$ and $\dot{\beta}(0) = \dot{\gamma}(0)$. Therefore, $\delta(0) = \dot{\delta}(0) = 0$ and (82) becomes $\delta(t) = \frac{t^2}{2}\ddot{\delta}(0) + h(t)t^2$. Due to (81) and the fact that $\lim_{t \rightarrow 0} h(t) = 0$, we can choose $\epsilon_* > 0$ small enough such that $\delta(\epsilon_*) = \frac{\epsilon_*^2}{2}\ddot{\delta}(0) + h(\epsilon_*)\epsilon_*^2 < 0$. This means that $\langle \Pi_n(\tilde{\mathbf{p}} + \epsilon_* V\tilde{\mathbf{q}}), \beta(\epsilon_*) \rangle < \langle \Pi_n(\tilde{\mathbf{p}} + \epsilon_* V\tilde{\mathbf{q}}), \tilde{\ell}(\tilde{\mathbf{p}} + \epsilon_* V\tilde{\mathbf{q}}) \rangle = \langle \Pi_n(\tilde{\mathbf{p}} + \epsilon_* V\tilde{\mathbf{q}}), \ell(\Pi_n(\tilde{\mathbf{p}} + \epsilon_* V\tilde{\mathbf{q}})) \rangle$. Therefore, $\beta(\epsilon_*)$ must lie outside the superprediction set. Thus, the mixability condition (8) does not hold for $P^{\epsilon_*} = \Pi_n[\tilde{\mathbf{p}}\mathbf{1}_k^\top + \epsilon_* V, \tilde{\mathbf{p}}] \in \text{ri } \Delta_n^k$. This completes the proof. \square

Theorem 14 Let $\ell : \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss such that \underline{L}_ℓ 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 $\text{int } \tilde{\Delta}_k$. Then ℓ is Φ -mixable only if $\eta_\ell \Phi - S$ is convex on Δ_k .

Proof. We will prove the contrapositive; suppose that $\eta_\ell \Phi - S$ is not convex on Δ_k and we show that ℓ cannot be Φ -mixable. Note first that from Lemma 22-(iii), $\tilde{\Phi}^*$ is twice differentiable on \mathbb{R}^{k-1} . Thus Lemmas 23 and 24 apply. Let $\underline{\ell}$ be a proper support loss of ℓ and suppose that $\eta_\ell \Phi - S$ is not convex on Δ_k . This implies that $\eta_\ell \tilde{\Phi} - \tilde{S}$ is not convex on $\text{int } \tilde{\Delta}_k$, and by Lemma 3 there exists $\tilde{\mathbf{q}}_* \in \text{int } \tilde{\Delta}_k$, such that $1 > \eta_\ell \lambda_{\min}(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}_*)(\mathbf{H}\tilde{S}(\tilde{\mathbf{q}}_*))^{-1})$. From this and the definition of η_ℓ , there exists $\tilde{\mathbf{p}}_* \in \text{int } \tilde{\Delta}_n$ such that

$$1 > \frac{\lambda_{\min}(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}_*)(\mathbf{H}\tilde{S}(\tilde{\mathbf{q}}_*))^{-1})}{\lambda_{\max}([\mathbf{H}\tilde{L}_{\log}(\tilde{\mathbf{p}}_*)]^{-1}\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}}_*))} = \frac{\lambda_{\min}(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}_*)(\text{diag}(\tilde{\mathbf{q}}_*) - \tilde{\mathbf{q}}_*\tilde{\mathbf{q}}_*^\top))}{\lambda_{\max}([\mathbf{H}\tilde{L}_{\log}(\tilde{\mathbf{p}}_*)]^{-1}\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}}_*))}, \quad (83)$$

where the equality is due to Lemma 22-(iv). For the rest of this proof let $(\tilde{\mathbf{p}}, \tilde{\mathbf{q}}) = (\tilde{\mathbf{p}}^*, \tilde{\mathbf{q}}^*)$. By assumption, \tilde{L}_ℓ twice differentiable and concave on $\text{int } \tilde{\Delta}_n$, and thus $-\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})$ is symmetric positive semi-definite. Therefore, there exists a symmetric positive semi-definite matrix $\Lambda_{\tilde{\mathbf{p}}}$ such that $\Lambda_{\tilde{\mathbf{p}}}\Lambda_{\tilde{\mathbf{p}}} = -\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})$. From Lemma 22-(i), $\tilde{\Phi}$ is strictly convex on $\text{int } \tilde{\Delta}_k$, and so there exists a symmetric positive definite matrix $K_{\tilde{\mathbf{q}}}$ such that $K_{\tilde{\mathbf{q}}}K_{\tilde{\mathbf{q}}} = \mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}})$. Let $\mathbf{w} \in \mathbb{R}^{n-1}$ be the unit norm eigenvector of $[\mathbf{H}\tilde{L}_{\log}(\tilde{\mathbf{p}})]^{-1}\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})$ associated with $\lambda_*^\ell := \lambda_{\max}([\mathbf{H}\tilde{L}_{\log}(\tilde{\mathbf{p}})]^{-1}\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}}))$. Suppose that $c_\ell := \mathbf{w}^\top \mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} = 0$. Since $\mathbf{w}^\top \Lambda_{\tilde{\mathbf{p}}}\Lambda_{\tilde{\mathbf{p}}}\mathbf{w} = -c_\ell = 0$, it follows from the positive semi-definiteness of $\Lambda_{\tilde{\mathbf{p}}}$ that $\Lambda_{\tilde{\mathbf{p}}}\mathbf{w} = \mathbf{0}_{\tilde{n}}$, and thus $\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} = -\Lambda_{\tilde{\mathbf{p}}}\Lambda_{\tilde{\mathbf{p}}}\mathbf{w} = \mathbf{0}_{\tilde{n}}$. This implies that $\lambda_*^\ell = 0$, which is not possible due to (83). Therefore, $\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} \neq \mathbf{0}_{\tilde{n}}$. Furthermore, the negative semi-definiteness of $\mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})$ implies that

$$c_\ell = \mathbf{w}^\top \mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} < 0. \quad (84)$$

Let $\mathbf{v} \in \mathbb{R}^{k-1}$ be the unit norm eigenvector of $K_{\tilde{\mathbf{q}}}(\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top)K_{\tilde{\mathbf{q}}}$ associated with $\lambda_*^\Phi := \lambda_{\min}(K_{\tilde{\mathbf{q}}}(\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top)K_{\tilde{\mathbf{q}}}) = \lambda_{\min}(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}})(\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top))$, where the equality is due to Lemma 2. Let $\hat{\mathbf{v}} := K_{\tilde{\mathbf{q}}}\mathbf{v}$.

We will show that for $V = \mathbf{w}\hat{\mathbf{v}}^\top$, the parametrized curve β defined in Lemma 23 satisfies $\left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \right|_{t=0} < 0$, where $\gamma(t) = \tilde{\ell}(\tilde{\mathbf{p}} + tV\tilde{\mathbf{q}})$. According to Lemma 24 this would imply that there exists $P \in \text{ri } \Delta_n^k$, such that $[\text{Mix}_\Phi(\underline{\ell}_x(P), \mathbf{q})]_{x \in [n]}^\top$ lies outside \mathcal{S}_ℓ . From Theorem 5, we know that there exists $A_* \in \mathcal{A}^k$, such that $\ell_x(A_*) = \underline{\ell}_x(P), \forall x \in [n]$. Therefore, $[\text{Mix}_\Phi(\ell_x(A_*), \mathbf{q})]_{x \in [n]}^\top = [\text{Mix}_\Phi(\underline{\ell}_x(P), \mathbf{q})]_{x \in [n]}^\top \notin \mathcal{S}_\ell$, and thus ℓ is not Φ -mixable.

From Lemma 23 (Equation 79) and the fact that $V_{\cdot,j} = \hat{\mathbf{v}}_j \mathbf{w}$, for $j \in [\tilde{k}]$, we can write

$$\begin{aligned} \left. \frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) \rangle \right|_{t=0} &= - \sum_{j=1}^{k-1} q_j \hat{\mathbf{v}}_j^2 \mathbf{w}^\top \mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} - \text{tr}(\text{diag}(\mathbf{p}) D\tilde{\ell}(\tilde{\mathbf{p}})V(\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}))^{-1}(D\tilde{\ell}(\tilde{\mathbf{p}})V)^\top), \\ &= - \langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \odot \hat{\mathbf{v}} \rangle \mathbf{w}^\top \mathbf{H}\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{w} - (\hat{\mathbf{v}}^\top (\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}}))^{-1} \hat{\mathbf{v}}) \langle \mathbf{p}, [D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{w}] \odot [(D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{w})], \end{aligned}$$

where the second equality is obtained by noting that 1) $(\hat{\mathbf{v}}^\top (\mathbf{H}\tilde{\Phi}(\mathbf{q}))^{-1} \hat{\mathbf{v}})$ is a scalar quantity and can be factorized out; and 2) $\text{tr}(\text{diag}(\mathbf{p}) \mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) \mathbf{w} (\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) \mathbf{w})^\top) = \langle \mathbf{p}, (\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) \mathbf{w}) \odot (\mathbf{D}\tilde{\ell}(\tilde{\mathbf{p}}) \mathbf{w}) \rangle$.

On the other hand, from Lemma 9, $\frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \Big|_{t=0} = -\langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \rangle^2 \mathbf{w}^\top \mathbf{H}\tilde{\mathcal{L}}_\ell(\tilde{\mathbf{q}}) \mathbf{w}$. Using (5) and the definition of c_ℓ , we get

$$\begin{aligned}
\frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \Big|_{t=0} &= [-\langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \odot \hat{\mathbf{v}} \rangle + \langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \rangle^2] c_\ell + \\
&\quad (\hat{\mathbf{v}}^\top (\mathbf{H}\tilde{\Phi}(\mathbf{q}))^{-1} \hat{\mathbf{v}}) (\mathbf{w}^\top (\mathbf{H}\tilde{\mathcal{L}}_\ell(\tilde{\mathbf{p}})) (\mathbf{H}\tilde{\mathcal{L}}_{\log}(\tilde{\mathbf{p}}))^{-1} \mathbf{H}\tilde{\mathcal{L}}_\ell(\mathbf{p}) \mathbf{w}), \\
&= -c_\ell [\langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \odot \hat{\mathbf{v}} \rangle - \langle \tilde{\mathbf{q}}, \hat{\mathbf{v}} \rangle^2 - \lambda_*^\ell (\hat{\mathbf{v}}^\top (\mathbf{H}\tilde{\Phi}(\mathbf{q}))^{-1} \hat{\mathbf{v}})], \\
&= -c_\ell [\hat{\mathbf{v}}^\top (\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top) \hat{\mathbf{v}} - \lambda_*^\ell (\hat{\mathbf{v}}^\top (\mathbf{H}\tilde{\Phi}(\mathbf{q}))^{-1} \hat{\mathbf{v}})], \\
&= -c_\ell [\hat{\mathbf{v}}^\top (\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top) \hat{\mathbf{v}} - \lambda_*^\ell (\mathbf{v}^\top K_q (K_q K_q)^{-1} K_q \mathbf{v})], \\
&= -c_\ell [\mathbf{v}^\top K_q (\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top) K_q \mathbf{v} - \lambda_*^\ell], \\
&= -c_\ell [\lambda_*^\Phi - \lambda_*^\ell], \\
&= -c_\ell [\lambda_{\min}(\mathbf{H}\tilde{\Phi}(\mathbf{q})(\text{diag}(\tilde{\mathbf{q}}) - \tilde{\mathbf{q}}\tilde{\mathbf{q}}^\top)) - \lambda_{\max}(\mathbf{H}\tilde{\mathcal{L}}_\ell(\tilde{\mathbf{p}})(\mathbf{H}\tilde{\mathcal{L}}_{\log}(\tilde{\mathbf{p}}))^{-1})],
\end{aligned} \tag{85}$$

where in (85) we used the fact that $\mathbf{v}^\top \mathbf{v} = 1$. The last equality combined with (83) and (84) shows that $\frac{d}{dt} \langle \mathbf{p}, \dot{\beta}(t) - \dot{\gamma}(t) \rangle \Big|_{t=0} < 0$, which completes the proof. \square

C.8 Proof of Lemma 15

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 $\eta_\ell \in \mathfrak{H}_\ell$.*

Proof. Suppose $\mathfrak{H}_\ell \neq \emptyset$. Let $\mathbf{q} \in \Delta_k$, $A := \mathbf{a}_{1:k} \in \mathcal{A}^k$. By definition of η_ℓ there exists $(\eta_m) \subset [0, +\infty[$ such that ℓ is η_m -mixable and $\eta_m \xrightarrow{m \rightarrow \infty} \eta_\ell$. Therefore, $\forall m \in \mathbb{N}$, $\exists \mathbf{a}_m \in \mathcal{A}$ such that

$$\forall x \in [n], \ell_x(\mathbf{a}_m) \leq -\eta_m^{-1} \log \langle \mathbf{q}, \exp(-\eta_m(\ell_x(A))) \rangle < +\infty, \tag{86}$$

where the right-most inequality follows from the fact $\text{dom } \ell = \mathcal{A}$. Therefore, the sequence $(\ell(\mathbf{a}_m)) \subset [0, +\infty]^n$ is bounded, and thus admits a convergent subsequence. If we let \mathbf{s} be the limit of this subsequence, then from (86) it follows that

$$\forall x \in [n], \mathbf{s} \leq -\eta_\ell^{-1} \log \langle \mathbf{q}, \exp(-\eta_\ell(\ell_x(A))) \rangle, \tag{87}$$

On the other hand, since ℓ is closed (by Assumption 1), it follows that there exists $\mathbf{a}_* \in \mathcal{A}$ such that $\ell(\mathbf{a}_*) = \mathbf{s}$. Combining this with (87) implies that ℓ is η_ℓ -mixable, and thus $\eta_\ell \in \mathfrak{H}_\ell$. \square

C.9 Proof of Theorem 17

Theorem 17 *Let ℓ and Φ be as in Theorem 16. Then*

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

Proof. From Theorem 16, ℓ is Φ_η -mixable if and only if $\underline{\eta}_\ell \Phi_\eta - \mathbf{S} = \eta^{-1} \underline{\eta}_\ell \Phi - \mathbf{S}$ is convex on Δ_k . When this is the case, Lemma 3 implies that

$$1 \leq \eta^{-1} \underline{\eta}_\ell \left(\inf_{\tilde{\mathbf{q}} \in \text{int } \tilde{\Delta}_k} \lambda_{\min}[\mathbf{H}\tilde{\Phi}(\tilde{\mathbf{q}})[\mathbf{H}\tilde{\mathcal{S}}(\tilde{\mathbf{q}})]^{-1}] \right), \tag{88}$$

where we used the facts that $\mathbf{H}(\eta^{-1} \underline{\eta}_\ell \tilde{\Phi}) = \eta^{-1} \underline{\eta}_\ell \mathbf{H}\tilde{\Phi}$, $\lambda_{\min}(\cdot)$ is linear, and $\eta^{-1} \underline{\eta}_\ell$ is independent of $\tilde{\mathbf{q}} \in \text{int } \tilde{\Delta}_k$. Inequality 88 shows that the largest η such that ℓ is Φ_η -mixable is given by η_ℓ^Φ in (11). \square

C.10 Proof of Theorem 18

Theorem 18 Let $S, \Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$, where S is the Shannon entropy and Φ 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 $\underline{\ell}$ twice differentiable on $]0, +\infty[^n$, is Φ -mixable only if $R_\ell^S \leq R_\ell^\Phi$.

Proof. Suppose ℓ is Φ -mixable. Then from Theorem 16, $\eta_\ell \Phi - S$ is convex on Δ_k , and thus $\underline{\eta}_\ell = \eta_\ell^S > 0$ (Corollary 17). Furthermore, $\underline{\eta}_\ell \tilde{\Phi} - \tilde{S} = [\underline{\eta}_\ell \Phi - S] \circ \Pi_k$ is convex on $\text{int } \tilde{\Delta}_k$, since Π_k is an affine function. It follows from Lemma 3 and Corollary 17 that

$$\eta_\ell^\Phi = \underline{\eta}_\ell \inf_{\tilde{q} \in \text{int } \tilde{\Delta}_k} \lambda_{\min}(\text{H}\tilde{\Phi}(\tilde{q})(\text{H}\tilde{S}(\tilde{q}))^{-1}) \geq 1 > 0.$$

Let $\mu \in \text{ri } \Delta_k$ and $\theta_* := \arg\max_\theta D_S(e_\theta, \mu)$. By definition of an entropy and the fact that the directional derivatives $\Phi'(\mu; \cdot)$ and $S'(\mu; \cdot)$ are finite on Δ_k [7, Prop. D.1.1.2], it holds that $D_\Phi(e_{\theta_*}, \mu), D_S(e_{\theta_*}, \mu) \in]0, +\infty[$. Therefore, there exists $\alpha > 0$ such that $\alpha^{-1} D_\Phi(e_{\theta_*}, \mu) = D_S(e_{\theta_*}, \mu)$. If we let $\Psi := \alpha^{-1} \Phi$, we get

$$D_\Psi(e_{\theta_*}, \mu) = D_S(e_{\theta_*}, \mu). \quad (89)$$

Let $d_\Psi(\tilde{q}) := \tilde{\Psi}(\tilde{q}) - \tilde{\Psi}(\tilde{\mu}) - \langle \tilde{q} - \tilde{\mu}, \nabla \tilde{\Psi}(\tilde{\mu}) \rangle$. Observe that

$$d_\Psi(\tilde{q}) = \Psi(q) - \Psi(\mu) - \langle q - \mu, \nabla \Psi(\mu) \rangle = D_\Psi(q, \mu).$$

We define d_S similarly. Suppose that $\eta_\ell^\Psi > \eta_\ell^S = \underline{\eta}_\ell$. Then, from Corollary 17, $\forall \tilde{q} \in \text{int } \tilde{\Delta}_k$, $\lambda_{\min}(\text{H}\tilde{\Psi}(\tilde{q})(\text{H}\tilde{S}(\tilde{q}))^{-1}) > 1$. This implies that $\forall \tilde{q} \in \text{int } \tilde{\Delta}_k$, $\lambda_{\min}(\text{H}d_\Psi(\tilde{q})(\text{H}d_S(\tilde{q}))^{-1}) > 1$, and from Lemma 3, $d_\Psi - d_S$ must be strictly convex on $\text{int } \tilde{\Delta}_k$. We also have $\nabla d_\Psi(\tilde{\mu}) - \nabla d_S(\tilde{\mu}) = 0$ and $d_\Psi(\tilde{\mu}) - d_S(\tilde{\mu}) = 0$. Therefore, $d_\Psi - d_S$ attains a strict minimum at $\tilde{\mu}$ (ibid., Thm. D.2.2.1); that is, $d_\Psi(\tilde{q}) > d_S(\tilde{q})$, $\forall \tilde{q} \in \tilde{\Delta}_k \setminus \{\tilde{\mu}\}$. In particular, for $\tilde{q} = \Pi_k(e_{\theta_*})$, we get $D_\Psi(e_{\theta_*}, \mu) = d_\Psi(\tilde{q}) > d_S(\tilde{q}) = D_S(e_{\theta_*}, \mu)$, which contradicts (89). Therefore, $\eta_\ell^\Psi \leq \eta_\ell^S$, and thus

$$\begin{aligned} R_\ell^S(\mu) &= \max_\theta D_S(e_\theta, \mu) / \eta_\ell^S = D_S(e_{\theta_*}, \mu) / \eta_\ell^S, \\ &\leq D_\Psi(e_{\theta_*}, \mu) / \eta_\ell^\Psi, \end{aligned} \quad (90)$$

$$\begin{aligned} &\leq \max_\theta D_\Psi(e_\theta, \mu) / \eta_\ell^\Psi, \\ &= R_\ell^\Psi(\mu), \end{aligned} \quad (91)$$

where (90) is due to $D_\Psi(e_{\theta_*}, \mu) = D_S(e_{\theta_*}, \mu)$ and $\eta_\ell^\Psi \leq \eta_\ell^S$. Equation 91, implies that $R_\ell^S(\mu) \leq R_\ell^\Psi(\mu)$, since $R_\ell^\Psi(\mu) = R_\ell^{\alpha\Phi}(\mu) = R_\ell^\Phi(\mu)$ [9]. Therefore,

$$\forall \mu \in \text{ri } \Delta_k, R_\ell^S(\mu) \leq R_\ell^\Phi(\mu). \quad (92)$$

It remains to consider the case where μ is in the relative boundary of Δ_k . Let $\mu \in \text{rbd } \Delta_k$. There exists $l_0 \subsetneq [k]$ such that $\mu \in \Delta_{l_0}$. Let $\theta^* \in [k] \setminus l_0$ and $l := l_0 \cup \{\theta^*\}$. It holds that $\mu \in \text{rbd } \Delta_l$ and $\mu + 2^{-1}(e_{\theta^*} - \mu) \in \text{ri } \Delta_l$. Since ℓ is Φ -mixable, it follows from Proposition 10 and the 1-homogeneity of $\Phi'(\mu; \cdot)$ [7, Prop. D.1.1.2] that

$$\Phi'(\mu; e_{\theta^*} - \mu) = 2\Phi'(\mu; [\mu + 2^{-1}(e_{\theta^*} - \mu)] - \mu) = -\infty.$$

Hence,

$$\begin{aligned} R_\ell^\Phi(\mu) &= \max_{\theta \in [k]} D_\Phi(e_\theta, \mu), \\ &\geq D_\Phi(e_{\theta^*}, \mu) = \Phi(e_{\theta^*}) - \Phi(\mu) - \Phi'(\mu; e_{\theta^*} - \mu) = +\infty. \end{aligned} \quad (93)$$

Inequality 93 also applies to S , since ℓ is $(\underline{\eta}_\ell^{-1} S)$ -mixable. From (93) and (92), we conclude that $\forall \mu \in \Delta_k, R_\ell^S(\mu) \leq R_\ell^\Phi(\mu)$. \square

C.11 Proof of Theorem 19

Theorem 19 Let $\Phi : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be a Δ -differentiable entropy. Let $\ell : \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss such that \underline{L}_ℓ 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 ℓ is (η, Φ) -mixable then for initial distribution $\mathbf{q}^0 = \operatorname{argmin}_{\mathbf{q} \in \Delta_k} \max_{\theta \in [k]} D_\Phi(\mathbf{e}_\theta, \mathbf{q})$ and any sequence $(x^t, \mathbf{a}_{1:k}^t)_{t=1}^T$, the AGAA achieves the regret

$$\forall \theta \in [k], \quad \text{Loss}_{AGAA}^\ell(T) - \text{Loss}_\theta^\ell(T) \leq R_\ell^\Phi + \sum_{t=1}^{T-1} (v_\theta^t - \langle \mathbf{v}^t, \mathbf{q}^t \rangle).$$

Proof. Recall that $\Phi_t(\mathbf{w}) := \Phi(\mathbf{w}) - \langle \mathbf{w}, \beta^t - \theta^t \rangle$, where $\theta^t = -\eta \sum_{s=1}^{t-1} \ell_{x^s}(A^s)$. From Theorem 16 and since Φ_t is equal to Φ plus an affine function, it is clear that if ℓ is (η, Φ) -mixable then ℓ is (η, Φ_t) -mixable. Thus, for all $(A^t, \mathbf{q}^{t-1}) \in \mathcal{A}^k \times \Delta_k$, there exists $\mathbf{a}_*^t \in \mathcal{A}$ such that for any outcome $x^t \in [n]$

$$\ell_{x^t}(\mathbf{a}_*^t) \leq \eta^{-1} [\Phi_t^*(\nabla \Phi_t(\mathbf{q}^{t-1})) - \Phi_t^*(\nabla \Phi_t(\mathbf{q}^{t-1}) - \eta \ell_{x^t}(A^t))].$$

Summing over t from 1 to T , we get

$$\begin{aligned} \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) &\leq \eta^{-1} [\Phi_1^*(\nabla \Phi_1(\mathbf{q}^0)) - \Phi_T^*(\nabla \Phi_T(\mathbf{q}^{T-1}) - \eta \ell_{x^T}(A^T))] \\ &\quad + \eta^{-1} \sum_{t=1}^{T-1} [\Phi_{t+1}^*(\nabla \Phi_{t+1}(\mathbf{q}^t)) - \Phi_t^*(\nabla \Phi_t(\mathbf{q}^{t-1}) - \eta \ell_{x^t}(A^t))]. \end{aligned} \quad (94)$$

ODue to the properties of the entropic dual [9] and the definition of Φ_t , the following holds for all $t \in [T]$ and $\mathbf{z} \in \mathbb{R}^k$,

$$\nabla \Phi_t(\mathbf{q}^{t-1}) = -\eta \sum_{s=1}^{t-1} \ell_{x^s}(A^s), \quad (95)$$

$$\Phi_t^*(\mathbf{z}) = \Phi^*(\mathbf{z} + \nabla \Phi(\mathbf{q}^{t-1}) + \eta \sum_{s=1}^{t-1} \ell_{x^s}(A^s)), \quad (96)$$

$$\nabla \Phi(\mathbf{q}^t) = \nabla \Phi(\mathbf{q}^{t-1}) - \eta \ell_{x^t}(A^t) - \eta \mathbf{v}^t. \quad (97)$$

Using (95)-(96), we get for all $0 \leq t < T$, $\Phi_{t+1}^*(\nabla \Phi_{t+1}(\mathbf{q}^t)) = \Phi^*(\nabla \Phi(\mathbf{q}^t))$, and in particular $\Phi_1^*(\nabla \Phi_1(\mathbf{q}^0)) = \Phi^*(\nabla \Phi(\mathbf{q}^0))$. Similarly, using (95)-(97), gives $\Phi_t^*(\nabla \Phi_t(\mathbf{q}^{t-1}) - \eta \ell_{x^t}(A^t)) = \Phi^*(\nabla \Phi(\mathbf{q}^t) + \eta \mathbf{v}^t)$ for all $1 \leq t \leq T$. Substituting back into (94) yields

$$\begin{aligned} \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_*^t) &\leq \eta^{-1} [\Phi^*(\nabla \Phi(\mathbf{q}^0)) - \Phi^*(\nabla \Phi(\mathbf{q}^T) + \eta \mathbf{v}^T)] \\ &\quad + \eta^{-1} \sum_{t=1}^{T-1} [\Phi^*(\nabla \Phi(\mathbf{q}^t)) - \Phi^*(\nabla \Phi(\mathbf{q}^t) + \eta \mathbf{v}^t)], \end{aligned} \quad (98)$$

To conclude the proof, we note that since Φ is convex it holds that

$$\Phi^*(\nabla \Phi(\mathbf{q}^t)) - \Phi^*(\nabla \Phi(\mathbf{q}^t) + \eta \mathbf{v}^t) \leq -\eta \langle \mathbf{v}^t, \nabla \Phi^*(\nabla \Phi(\mathbf{q}^t)) \rangle = -\eta \langle \mathbf{v}^t, \mathbf{q}^t \rangle, \quad (99)$$

which allows us to bound the sum on the right hand side of 98. To bound the rest of the terms, we use the fact that $\nabla \Phi(\mathbf{q}^T) = \nabla \Phi(\mathbf{q}^0) - \eta \sum_{t=1}^T (\ell_{x^t}(A^t) + \mathbf{v}^t)$, and thus by letting $\Phi_\eta := \eta^{-1} \Phi$,

$$\begin{aligned} \eta^{-1} [\Phi^*(\nabla \Phi(\mathbf{q}^0)) - \Phi^*(\nabla \Phi(\mathbf{q}^T) + \eta \mathbf{v}^T)] &= \Phi_\eta^*(\nabla \Phi_\eta(\mathbf{q}^0)) \\ &\quad - \Phi_\eta^* \left(\nabla \Phi_\eta(\mathbf{q}^0) - \sum_{t=1}^T \ell_{x^t}(A^t) - \sum_{t=1}^{T-1} \mathbf{v}^t \right), \\ &= \inf_{\mathbf{q} \in \Delta_k} \left\langle \mathbf{q}, \sum_{t=1}^T \ell_{x^t}(A^t) + \sum_{t=1}^{T-1} \mathbf{v}^t \right\rangle + \frac{D_\Phi(\mathbf{q}, \mathbf{q}^0)}{\eta}, \\ &\leq \sum_{t=1}^T \ell_{x^t}(\mathbf{a}_\theta^t) + \sum_{t=1}^{T-1} v_\theta^t + \frac{D_\Phi(\mathbf{e}_\theta, \mathbf{q}^0)}{\eta}, \forall \theta \in [k]. \end{aligned}$$

Substituting this last inequality and (99) back into (98) yields the desired bound. \square

D Defining the Bayes Risk Using the Superprediction Set

In this section, we argue that when a loss $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ is mixable, in the classical or generalized sense, it does not matter whether we define the Bayes risk \underline{L}_ℓ using the full superprediction set \mathcal{S}_ℓ^∞ or its finite part \mathcal{S}_ℓ . Recall the definition of the Bayes risk;

Definition 2 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

$$\forall \mathbf{u} \in \mathbb{R}^n, \quad \underline{L}_\ell(\mathbf{u}) := \inf_{\mathbf{z} \in \mathcal{S}_\ell} \langle \mathbf{u}, \mathbf{z} \rangle. \quad (100)$$

Note that the right hand side of (100) does not change if we substitute \mathcal{S}_ℓ for its closure $\overline{\mathcal{S}_\ell}$ with respect to $[0, +\infty]^n$. Thus, it suffices to show that $\mathcal{S}_\ell^\infty \subseteq \overline{\mathcal{S}_\ell}$ when the loss ℓ is mixable. We show this in Theorem 26, but first we give a characterization of the (finite part) of the superprediction set for a proper loss.

Proposition 25. Let $\ell: \Delta_n \rightarrow [0, +\infty]^n$ be a proper loss. If \underline{L}_ℓ is differentiable on $]0, +\infty[^n$, then

$$\overline{\mathcal{S}_\ell} \supseteq \mathcal{C}_\ell := \{\mathbf{u} \in [0, +\infty]^n : \forall \mathbf{p} \in \Delta_n, \underline{L}_\ell(\mathbf{p}) \leq \langle \mathbf{p}, \mathbf{u} \rangle\}. \quad (101)$$

Proof. Let $\mathbf{v} \in \mathcal{C}_\ell \cap [0, +\infty]^n$. Let $f: \text{ri } \Delta_n \times [n] \rightarrow \mathbb{R}$ be defined by $f(\mathbf{p}, x) := \ell_x(\mathbf{p}) - v_x$. By the choice of \mathbf{v} , we have $\mathbb{E}_{x \sim \mathbf{p}} f(\mathbf{p}, x) = \langle \mathbf{p}, \ell(\mathbf{p}) \rangle - \langle \mathbf{p}, \mathbf{v} \rangle \leq 0$ for all $\mathbf{p} \in \Delta_n$. Since \underline{L}_ℓ is differentiable on $]0, +\infty[^n$, by assumption, ℓ is continuous on $\text{ri } \Delta_n$, and thus f is continuous in the first argument. Since \mathbf{v} has finite components, the map f satisfies all the conditions of Lemma 5. Therefore, there exists $(\mathbf{p}_m) \subset \text{ri } \Delta_n$ such that

$$\forall m \in \mathbb{N}, \forall x \in [n], \ell_x(\mathbf{p}_m) \leq v_x + \frac{1}{m}. \quad (102)$$

Without loss of generality, we can assume by extracting a subsequence if necessary that $\ell(\mathbf{p}_m)$ converges to $\mathbf{s} \in [0, +\infty]^n$. By definition, we have $\mathbf{s} \in \overline{\mathcal{S}_\ell}$, and from (102) it follows that $\mathbf{s} \leq \mathbf{v}$ coordinate-wise. Thus, \mathbf{v} is in $\overline{\mathcal{S}_\ell}$.

The above argument shows that $\mathcal{C}_\ell \cap [0, +\infty]^n \subseteq \overline{\mathcal{S}_\ell}$, and since $\overline{\mathcal{S}_\ell}$ is closed in $[0, +\infty]^n$ we have $\overline{\mathcal{C}} \subseteq \overline{\mathcal{S}_\ell}$, where $\overline{\mathcal{C}}$ is the closure of $\mathcal{C}_\ell \cap [0, +\infty]^n$ in $[0, +\infty]^n$. Now it suffice to show that $\mathcal{C}_\ell \subseteq \overline{\mathcal{C}}$ to complete the proof.

Let $\mathbf{u} \in \mathcal{C}_\ell$ and $\mathfrak{l} := \{x \in [n] : u_x < +\infty\}$. Define $(\mathbf{u}_m) \subset [0, +\infty]^n$ by $u_{m,x} = u_x$ if $x \in \mathfrak{l}$; and m otherwise. Let $\mathbf{p} \in \Delta_n$. It follows that

$$\begin{aligned} \langle \mathbf{p}, \mathbf{u}_m \rangle &= \sum_{x' \in \mathfrak{l}} p_{x'} u_{m,x'} + \sum_{x \notin \mathfrak{l}} p_x u_{m,x}, \\ &= \sum_{x' \in \mathfrak{l}} p_{x'} u_{m,x'} + \sum_{x \notin \mathfrak{l}} p_x u_{m,x}. \end{aligned} \quad (103)$$

Claim 2. $\forall \epsilon > 0, \exists m_\epsilon \geq 1, \forall \mathbf{p} \in \Delta_k, \underline{L}_\ell(\mathbf{p}) \leq \langle \mathbf{p}, \mathbf{u}_{m_\epsilon} \rangle - \epsilon$.

Suppose that Claim 2 is false. This means that there exists $\delta > 0$ such that

$$\forall m \geq 1, \exists \mathbf{p}_m \in \Delta_n, \langle \mathbf{p}_m, \mathbf{u}_m \rangle - \delta < \underline{L}_\ell(\mathbf{p}_m). \quad (104)$$

We may assume, by extracting a subsequence if necessary (Δ_n is compact), that (\mathbf{p}_m) converges to $\mathbf{p}_* \in \Delta_n$. Taking the limit $m \rightarrow \infty$ in (104) would lead to the contradiction ' $\langle \mathbf{p}_*, \mathbf{u} \rangle < \underline{L}_\ell(\mathbf{p}_*)$ ', since from (103) we have $\lim_{m \rightarrow \infty} \langle \mathbf{p}_m, \mathbf{u}_m \rangle = \langle \mathbf{p}_*, \mathbf{u} \rangle$. Therefore, Claim 2 is true. For $\epsilon = \frac{1}{k}$ let $m_k := m_\epsilon$ be as in Claim (2). The claim then implies that $\liminf_{k \rightarrow \infty} \langle \mathbf{p}, \mathbf{u}_{m_k} \rangle \geq \underline{L}_\ell(\mathbf{p})$ uniformly for $\mathbf{p} \in \Delta_k$. By the claim we also have that $\mathbf{u}_{m_k} \in \mathcal{C}_\ell \cap [0, +\infty]^n$ for all $k \in \mathbb{N}$, and by construction of \mathbf{v}_m , we have $\lim_{k \rightarrow \infty} \mathbf{u}_{m_k} = \mathbf{u}$. This shows that $\mathcal{C}_\ell \subseteq \overline{\mathcal{C}}$, which completes the proof. \square

Theorem 26. Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss. If $\mathcal{S}_\ell^\infty \not\subseteq \overline{\mathcal{S}_\ell}$, then ℓ is not mixable.

Proof. Suppose that ℓ is mixable and let $\underline{\ell}$ be a proper support loss of ℓ . From Proposition 12, \underline{L}_ℓ is differentiable on $]0, +\infty[^n$, and thus Theorem 5 implies that $\overline{\mathcal{S}}_\ell = \mathcal{S}_{\underline{\ell}}$. Therefore, Lemma 25 implies that $\overline{\mathcal{S}}_\ell \supseteq \{\mathbf{u} \in [0, +\infty]^n : \forall \mathbf{p} \in \Delta_n, \underline{L}_\ell(\mathbf{p}) \leq \langle \mathbf{p}, \mathbf{u} \rangle\}$. Thus, if $\mathcal{S}_\ell^\infty \not\subseteq \overline{\mathcal{S}}_\ell$, there exists $\epsilon > 0$, $\mathbf{p}_\epsilon \in \Delta_k$, and $\mathbf{s} \in \mathcal{S}_\ell^\infty \setminus \overline{\mathcal{S}}_\ell$ such that

$$\langle \mathbf{p}_\epsilon, \mathbf{s} \rangle < \underline{L}_\ell(\mathbf{p}_\epsilon) - 2\epsilon. \quad (105)$$

Note that \mathbf{p}_ϵ cannot be in $\text{ri } \Delta_n$; otherwise, (105) would imply that \mathbf{s} has all finite components, and thus would be included in $\overline{\mathcal{S}}_\ell$, which is a contradiction. Assume from now on that $\mathbf{p}_\epsilon \in \text{rbd } \Delta_n$. From the definition of the support loss, there exists a sequence $(\mathbf{p}_m) \subseteq \text{ri } \Delta_n$ such that $\mathbf{p}_m \xrightarrow{m \rightarrow \infty} \mathbf{p}_\epsilon$ and $\underline{\ell}(\mathbf{p}_m) \xrightarrow{m \rightarrow \infty} \underline{\ell}(\mathbf{p}_\epsilon)$. Therefore, Theorem 5 implies that there exists $\mathbf{a}_\epsilon \in \mathcal{A}$ such that

$$\langle \mathbf{p}_\epsilon, \ell(\mathbf{a}_\epsilon) \rangle < \langle \mathbf{p}_\epsilon, \underline{\ell}(\mathbf{p}_\epsilon) \rangle + \epsilon. \quad (106)$$

To see this, note that since $(\mathbf{p}_m) \subset \text{ri } \Delta_n \subseteq \text{dom } \underline{\ell}$, Theorem 5 guarantees the existence of a sequence $(\mathbf{a}_m) \subset \mathcal{A}$ such that $\ell(\mathbf{a}_m) = \underline{\ell}(\mathbf{p}_m)$. On the other hand, for any $x \in [n]$ such that $\ell_x(\mathbf{p}_\epsilon) = +\infty$, we have $p_{\epsilon,x} = 0$ — otherwise, $\underline{L}_\ell(\mathbf{p}_\epsilon)$ would be infinite. It follows, by continuity of the inner product that $\langle \mathbf{p}_\epsilon, \ell(\mathbf{a}_m) \rangle \xrightarrow{m \rightarrow \infty} \langle \mathbf{p}_\epsilon, \underline{\ell}(\mathbf{p}_m) \rangle$, and thus it suffices to pick \mathbf{a}_ϵ equal to \mathbf{a}_m for m large enough.

Now since ℓ is η -mixable, there exists $\eta > 0$ and $\mathbf{a}_* \in \mathcal{A}$ such that

$$\ell(\mathbf{a}_*) \leq -\eta^{-1} \log \left(\frac{1}{2} e^{-\eta \mathbf{s}} + \frac{1}{2} e^{-\eta \ell(\mathbf{a}_\epsilon)} \right),$$

and due to the convexity of $-\log$,

$$\leq \frac{1}{2} \mathbf{s} + \frac{1}{2} \ell(\mathbf{a}_\epsilon).$$

Using (105) and (106) yields

$$\langle \mathbf{p}_\epsilon, \ell(\mathbf{a}_*) \rangle \leq \underline{L}_\ell(\mathbf{p}_\epsilon) - \epsilon/2. \quad (107)$$

On the other hand, by definition of a proper support loss, $\langle \mathbf{p}_\epsilon, \underline{\ell}(\mathbf{p}_\epsilon) \rangle \leq \langle \mathbf{p}_\epsilon, \ell(\mathbf{a}_*) \rangle$. This combined with (107), lead to the contradiction $\langle \mathbf{p}_\epsilon, \underline{\ell}(\mathbf{p}_\epsilon) \rangle < \underline{L}_\ell(\mathbf{p}_\epsilon)$. \square

E The Update Step of the GAA and the Mirror Descent Algorithm

In this section, we demonstrate that the update steps of the GAA and the Mirror Descent Algorithm are essentially the same (at least for finite losses) according to the definition of the MDA given by Beck and Teboulle [2];

Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss and $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ an entropy such that $\tilde{\Phi}$ is differentiable on $\text{int } \tilde{\Delta}_k$. Let \mathbf{q}^t be the update distribution of the GAA at round t and $\tilde{\mathbf{q}}^t = \Pi_k(\mathbf{q}^t)$. It follows from the definition of \mathbf{q}^t (see Algorithm 2) that

$$\begin{aligned} \tilde{\mathbf{q}}^t &= \underset{\tilde{\mathbf{q}} \in \tilde{\Delta}_k}{\text{argmin}} \langle \Pi_k(\tilde{\mathbf{q}}), \ell_{x^t}(A^t) \rangle + \eta^{-1} D_{\tilde{\Phi}}(\tilde{\mathbf{q}}, \tilde{\mathbf{q}}^{t-1}), \\ &= \underset{\tilde{\mathbf{q}} \in \tilde{\Delta}_k}{\text{argmin}} \langle \tilde{\mathbf{q}}, J_k^\top \ell_{x^t}(A^t) \rangle + \eta^{-1} D_{\tilde{\Phi}}(\tilde{\mathbf{q}}, \tilde{\mathbf{q}}^{t-1}), \\ &= \underset{\tilde{\mathbf{q}} \in \tilde{\Delta}_k}{\text{argmin}} \langle \tilde{\mathbf{q}}, \nabla l_t(\tilde{\mathbf{q}}^{t-1}) \rangle + \eta^{-1} D_{\tilde{\Phi}}(\tilde{\mathbf{q}}, \tilde{\mathbf{q}}^{t-1}), \end{aligned} \quad (108)$$

where $l_t(\tilde{\boldsymbol{\mu}}) := \langle \Pi_k(\tilde{\boldsymbol{\mu}}), \ell_{x^t}(A^t) \rangle = \langle \boldsymbol{\mu}, \ell_{x^t}(A^t) \rangle$. Update (108) is, by definition [2], the MDA with the sequence of losses l_t on $\text{int } \tilde{\Delta}_k$, ‘distance’ function $D_{\tilde{\Phi}}(\cdot, \cdot)$, and learning rate η . Therefore, the MDA is exactly the update step of the GAA.

F The Generalized Aggregating Algorithm Using the Shannon Entropy S

The purpose of this appendix is to show that the GAA reduces to the AA when the former uses the Shannon entropy. In this case, generalized and classical mixability are equivalent. In what follows, we make use of the following proposition which is proved in C.5.

Proposition 20 For the Shannon entropy S , it holds that $\tilde{S}^*(\mathbf{v}) = \log(\langle \exp(\mathbf{v}), \mathbf{1}_k \rangle + 1)$, $\forall \mathbf{v} \in \mathbb{R}^{k-1}$, and $S^*(\mathbf{z}) = \log \langle \exp(\mathbf{z}), \mathbf{1}_k \rangle$, $\forall \mathbf{z} \in \mathbb{R}^k$.

Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss and Φ be as in Proposition 14 and suppose that Φ and $\tilde{\Phi}^*$ are differentiable on $\text{ri } \Delta_k$ and \mathbb{R}^{k-1} , respectively. It was shown in [9] that

$$\nabla \Phi^*(\nabla \Phi(\mathbf{q}) - \ell_x(A)) = \underset{\mu \in \Delta_k}{\operatorname{argmin}} \langle \mu, \ell_x(A) \rangle + D_\Phi(\mu, \mathbf{q}), \quad (109)$$

$$\text{Mix}_\Phi(\ell_x(A), \mathbf{q}) = \Phi^*(\nabla \Phi(\mathbf{q})) - \Phi^*(\nabla \Phi(\mathbf{q}) - \ell_x(A)). \quad (110)$$

Let $\mathbf{q} \in \text{ri } \Delta_k$. By definition of S , $\nabla S(\mathbf{q}) = \log \mathbf{q} + \mathbf{1}_k$, and due to Proposition 20, $S^*(\mathbf{z}) = \log \langle \exp \mathbf{z}, \mathbf{1}_k \rangle$, $\mathbf{z} \in \mathbb{R}^k$. Therefore, $\nabla S(\mathbf{q}) - \eta \ell_x(A) = \log(\exp(-\eta \ell_x(A)) \odot \mathbf{q}) + \mathbf{1}_k$ and $\nabla S^*(\mathbf{z}) = \frac{\exp \mathbf{z}}{\langle \exp \mathbf{z}, \mathbf{1}_k \rangle}$, $\forall (\mathbf{z}, A) \in [n] \times (\text{dom } \ell)^k$. Thus,

$$\nabla S^*(\nabla S(\mathbf{q}) - \eta \ell_x(A)) = \frac{\exp(-\eta \ell_x(A)) \odot \mathbf{q}}{\langle \exp(-\eta \ell_x(A)), \mathbf{q} \rangle}. \quad (111)$$

Let $S_\eta := \eta^{-1} S$. Then $\nabla S = \eta \nabla S_\eta$ and $\forall \mathbf{z} \in \mathbb{R}^k$, $\nabla S_\eta^*(\mathbf{z}) = \nabla S^*(\eta \mathbf{z})$ [9].² Then the left hand side of (111) can be written as $\nabla S_\eta^*(\nabla S_\eta(\mathbf{q}) - \ell_x(A))$. Using this fact, (109) and (111) show that the update distribution \mathbf{q}^t of the GAA (Algorithm 2) coincides with that of the AA after substituting \mathbf{q}, x , and A by \mathbf{q}^{t-1}, x^t , and $A^t := [\mathbf{a}_\theta]_{\theta \in [k]}$, respectively.

Now using the fact that $\text{Mix}_S^\eta(\ell_x(A), \mathbf{q}) = \eta^{-1} \text{Mix}_S(\eta \ell_x(A), \mathbf{q})$ [9] and (110), we get

$$\begin{aligned} \text{Mix}_S^\eta(\ell_x(A), \mathbf{q}) &= \eta^{-1} [S^*(\nabla S(\mathbf{q})) - S^*(\nabla S(\mathbf{q}) - \eta \ell_x(A))], \\ &= -\eta^{-1} \log \langle \exp(-\eta \ell_x(A)), \mathbf{q} \rangle. \end{aligned} \quad (112)$$

Equation 112 shows that the η -mixability condition is equivalent to the (η, S) -mixability condition for a finite loss. This remains true for losses taking infinite values — see the proof of Theorem 11 in Appendix C.5.

G Legendre Φ , but no Φ -mixable ℓ

In this appendix, we construct a *Legendre type* entropy [11] for which there are no Φ -mixable losses satisfying a weak condition (see below).

Let $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ be a loss satisfying condition 1. According to Alexandrov's Theorem, a concave function is twice differentiable almost everywhere (see e.g. [4, Thm. 6.7]). Now we give a version of Theorem 14 which does not assume the twice differentiability of the Bayes risk. The proof is almost identical to that of Theorem 14 with only minor modifications.

Theorem 27. Let $\Phi: \mathbb{R}^k \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy such that $\tilde{\Phi}$ is twice differentiable on $\text{int } \tilde{\Delta}_k$, and $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ a loss satisfying Condition 1 and such that $\exists(\tilde{\mathbf{p}}, \mathbf{v}) \in \mathcal{D} \times \mathbb{R}^{\tilde{n}}, \mathbf{H}_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v} \neq \mathbf{0}_{\tilde{n}}$, where $\mathcal{D} \subset \text{int } \tilde{\Delta}_n$ is a set of Lebesgue measure 1 where $\tilde{\ell}$ is twice differentiable, and define

$$\underline{\eta}_\ell^* := \inf_{\tilde{\mathbf{p}} \in \mathcal{D}} (\lambda_{\max}([\mathbf{H}_{\tilde{\ell}}(\tilde{\mathbf{p}})]^{-1} \mathbf{H}_{\tilde{\ell}}(\tilde{\mathbf{p}})))^{-1}. \quad (113)$$

Then ℓ is Φ -mixable only if $\underline{\eta}_\ell^* \Phi - S$ is convex on Δ_k .

The new condition on the Bayes risk is much weaker than requiring $\underline{\ell}$ to be twice differentiable on $]0, +\infty[^n$. In the next example, we will show that there exists a Legendre type entropy for which there are no Φ -mixable losses satisfying the condition of Theorem 27.

Example 1. Let $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy such that

$$\forall q \in]0, 1[, \Phi(q, 1-q) = \tilde{\Phi}(q) = \int_{1/2}^q \log \left(\frac{\log(1-t)}{\log t} \right) dt.$$

²Reid et al. [9] showed the equality $\nabla \Phi_\eta^*(\mathbf{u}) = \nabla \Phi^*(\eta \mathbf{u})$, $\forall \mathbf{u} \in \text{dom } \Phi^*$, for any entropy differentiable on Δ_k - not just for the Shannon Entropy.

$\tilde{\Phi}$ is differentiable and strictly convex on the open set $(0, 1)$. Furthermore, it satisfies (10) which makes it a function of Legendre type [11, Lem. 26.2]. In fact, (10) is satisfied due to

$$\left| \frac{d}{dq} \tilde{\Phi}(q) \right| = \left| \log \left(\frac{\log(1-q)}{\log q} \right) \right| \xrightarrow{q \rightarrow b} +\infty, \text{ where } b \in \{0, 1\},$$

$$\frac{d^2}{dq^2} \tilde{\Phi}(q) = \frac{-1}{q \log q} + \frac{-1}{(1-q) \log(1-q)} > 0, \forall q \in]0, 1[.$$

The Shannon entropy on Δ_2 is defined by $S(q, 1-q) = \tilde{S}(q) = q \log q + (1-q) \log(1-q)$, for $q \in]0, 1[$. Thus, $\frac{d^2}{dq^2} \tilde{S}(q) = \frac{1}{q(1-q)}$.

Suppose now that there exists a Φ -mixable loss $\ell: \mathcal{A} \rightarrow [0, +\infty]^n$ satisfying condition 1 and such that $\exists(\tilde{\mathbf{p}}, \mathbf{v}) \in \mathcal{D} \times \mathbb{R}^{\tilde{n}}, H_{\tilde{\ell}}(\tilde{\mathbf{p}})\mathbf{v} \neq \mathbf{0}_{\tilde{n}}$. Let $\underline{\eta}_\ell^*$ be as in (113). By definition, we have $\underline{\eta}_\ell^* < +\infty$, and thus

$$\underline{\eta}_\ell^* \left[\frac{d^2}{dq^2} \tilde{\Phi}(q) \right] \left[\frac{d^2}{dq^2} \tilde{S}(q) \right]^{-1} = \underline{\eta}_\ell^* \left(\frac{q-1}{\log q} + \frac{-q}{\log(1-q)} \right) \xrightarrow{q \rightarrow b} 0, \quad (114)$$

where $b \in \{0, 1\}$. From Lemma 3, (114) implies that $\underline{\eta}_\ell^* \Phi - S$ is not convex on Δ_k , which is a contradiction according to Theorem 27.

H Loss Surface and Superprediction Set

In this appendix, we derive an expression for the curvature of the image of a proper loss function. We will need the following lemma.

Lemma 28. Let $\sigma: [0, +\infty[^n \rightarrow \mathbb{R}$ be a 1-homogeneous, twice differentiable function on $]0, +\infty[^n$. Then σ is concave on $]0, +\infty[^n$ if and only if $\tilde{\sigma} = \sigma \circ \Pi_n$ is concave on $\text{int } \tilde{\Delta}_n$.

Proof. The forward implication is immediate; if σ is concave on $]0, +\infty[^n$, then $\sigma \circ \Pi_k$ is concave on $\text{int } \tilde{\Delta}_k$, since Π_k is an affine function.

Now assume that $\tilde{\sigma}$ is concave on $\text{int } \tilde{\Delta}_k$. Let $\lambda \in [0, 1]$ and $(\mathbf{p}, \mathbf{q}) \in [0, +\infty[^n \times [0, +\infty[^n$. We need to show that

$$\lambda \sigma(\mathbf{p}) + (1-\lambda) \sigma(\mathbf{q}) \leq \sigma(\lambda \mathbf{p} + (1-\lambda) \mathbf{q}). \quad (115)$$

Note that if $\mathbf{p} = \mathbf{0}$ or $\mathbf{q} = \mathbf{0}$, (115) is trivially with equality due to the 1-homogeneity of σ . Now assume that \mathbf{p} and \mathbf{q} are non-zero and let $c := \lambda \|\mathbf{p}\|_1 + (1-\lambda) \|\mathbf{q}\|_1$. For convenience, we also denote $\mathbf{p}_1 = \mathbf{p} / \|\mathbf{p}\|_1$ and $\mathbf{q}_1 = \mathbf{q} / \|\mathbf{q}\|_1$ which are both in Δ_n . It follows that

$$\begin{aligned} \lambda \sigma(\mathbf{p}) + (1-\lambda) \sigma(\mathbf{q}) &= cM \left(\lambda \frac{\|\mathbf{p}\|_1}{c} \sigma(\mathbf{p}_1) + (1-\lambda) \frac{\|\mathbf{q}\|_1}{c} \sigma(\mathbf{q}_1) \right), \\ &= c \left(\lambda \frac{\|\mathbf{p}\|_1}{c} \tilde{\sigma}(\tilde{\mathbf{p}}_1) + (1-\lambda) \frac{\|\mathbf{q}\|_1}{c} \tilde{\sigma}(\tilde{\mathbf{q}}_1) \right), \\ &\leq c \tilde{\sigma} \left(\lambda \frac{\|\mathbf{p}\|_1}{c} \tilde{\mathbf{p}}_1 + (1-\lambda) \frac{\|\mathbf{q}\|_1}{c} \tilde{\mathbf{q}}_1 \right), \\ &= c \sigma \left(\lambda \frac{\|\mathbf{p}\|_1}{c} \mathbf{p}_1 + (1-\lambda) \frac{\|\mathbf{q}\|_1}{c} \mathbf{q}_1 \right), \\ &= \sigma(\lambda \mathbf{p} + (1-\lambda) \mathbf{q}), \end{aligned}$$

where the first and last equalities are due the 1-homogeneity of σ and the inequality is due to $\tilde{\sigma}$ being concave on the $\text{int } \tilde{\Delta}_n$. \square

H.1 Convexity of the Superprediction Set

In the literature, many theoretical results involving loss functions relied on the fact that the superprediction set of a proper loss is convex [16, 6]. An earlier proof of this result by [16] was incomplete³. In the next theorem we restate this result.

Theorem 29. *If $\ell: \Delta_n \rightarrow [0, +\infty]^n$ is a continuous proper loss, then $\mathcal{S}_\ell = \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})}$. In particular, \mathcal{S}_ℓ is convex.*

$\mathcal{S}_\ell \subseteq \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})}$: Let $\mathbf{v} \in \mathcal{S}_\ell$, $\mathbf{u} \in [0, +\infty]^n$, and $\mathbf{q} \in \Delta_n$ such that $\mathbf{v} = \ell(\mathbf{q}) + \mathbf{u}$. Since ℓ is proper then $\forall \mathbf{p} \in \Delta_n$, $\underline{L}_\ell(\mathbf{p}) = \langle \mathbf{p}, \ell(\mathbf{p}) \rangle \leq \langle \mathbf{p}, \ell(\mathbf{q}) \rangle \leq \langle \mathbf{p}, \ell(\mathbf{q}) + \mathbf{u} \rangle = \langle \mathbf{p}, \mathbf{v} \rangle$. Therefore, $\mathbf{v} \in \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})}$.

$[\bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})} \subseteq \mathcal{S}_\ell]$: Let $\mathbf{v} \in \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})}$. Let $\Omega = [n]$, $\Delta(\Omega) = \Delta_n$, and $Q(\mathbf{p}, x) = \ell_x(\mathbf{p}) - v_x$ for all $(\mathbf{p}, x) \in \Delta_n \times [n]$. Since $\mathbf{v} \in \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{-\mathbf{p}, -\underline{L}_\ell(\mathbf{p})}$, $\mathbb{E}_{x \sim \mathbf{p}} Q(\mathbf{p}, x) = \langle \mathbf{p}, \ell(\mathbf{p}) \rangle - \langle \mathbf{p}, \mathbf{v} \rangle \leq 0$ for all $\mathbf{p} \in \Delta_n$. Lemma 4, implies that there exists $\mathbf{p}_* \in \Delta_n$ such that $Q(\mathbf{p}_*, x) = \ell_x(\mathbf{p}_*) - v_x \leq 0$, for all $x \in [n]$. This shows that $\mathbf{v} \in \mathcal{S}_\ell$. \square

H.2 Curvature of the Loss Surface

The *normal curvature* of a \tilde{n} -manifold \mathcal{S} [13] at a point $\mathbf{r} \in \mathcal{S}$ in the direction of $\mathbf{w} \in T_{\mathbf{r}}\mathcal{S}$, where $T_{\mathbf{r}}\mathcal{S}$ is the *tangent space* of \mathcal{S} at $\mathbf{r} \in \mathcal{S}$, is defined by

$$\kappa(\mathbf{r}, \mathbf{w}) = \frac{\langle \mathbf{w}, \text{DN}^{\mathcal{S}}(\mathbf{r})\mathbf{w} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle}, \quad (116)$$

where $\text{N}^{\mathcal{S}}(\mathbf{r})$ is the normal vector to the surface at \mathbf{r} . The *minimum principal curvature* of \mathcal{S} at \mathbf{r} is expressed as $\underline{\kappa}(\mathbf{r}) := \inf\{\kappa(\mathbf{r}, \mathbf{w}) : \mathbf{w} \in T_{\mathbf{r}}\mathcal{S} \cap \mathcal{B}(\mathbf{r}, 1)\}$.

In the next theorem, we establish a direct link between the curvature of a loss surface and the Hessian of the loss' Bayes risk.

Theorem 30. *Let $\ell: \text{ri } \Delta_n \rightarrow [0, +\infty]^n$ be a loss whose Bayes risk is twice differentiable and strictly concave on $]0, +\infty[^n$. Let $\mathbf{p} \in \text{ri } \Delta_n$, $X_{\mathbf{p}} := I_{\tilde{n}} - \tilde{\mathbf{p}}\mathbf{1}_{\tilde{n}}^T$, and $\mathbf{w} \in T_{\tilde{\ell}(\tilde{\mathbf{p}})}\mathcal{S}_\ell$. Then*

1. $\exists \mathbf{v} \in \mathbb{R}^{n-1}$ such that $\text{D}\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} = \mathbf{w}$.
2. \mathcal{S}_ℓ is a \tilde{n} -manifold.
3. The normal curvature of \mathcal{S}_ℓ at $\ell(\mathbf{p}) = \tilde{\ell}(\tilde{\mathbf{p}})$ in the direction \mathbf{w} is given by

$$\kappa_\ell(\ell(\mathbf{p}), \mathbf{w}) = \left\| \begin{bmatrix} X_{\mathbf{p}} \\ -\tilde{\mathbf{p}}^T \end{bmatrix} (-\text{H}\tilde{L}_\ell(\tilde{\mathbf{p}}))^{\frac{1}{2}} \mathbf{u} \right\|^{-1}, \quad (117)$$

where $\mathbf{u} = (-\text{H}\tilde{L}_\ell(\tilde{\mathbf{p}}))^{\frac{1}{2}} \mathbf{v} / \|(-\text{H}\tilde{L}_\ell(\tilde{\mathbf{p}}))^{\frac{1}{2}} \mathbf{v}\|$.

It becomes clear from (117) that smaller eigenvalues of $-\text{H}\tilde{L}_\ell(\tilde{\mathbf{p}})$ will tend to make the loss surface more curved at $\ell(\mathbf{p})$, and vice versa.

Before proving Theorem 30, we first define parameterizations on manifolds.

Definition 31 (Local and Global Parameterization). *Let $\mathcal{S} \subseteq \mathbb{R}^n$ be a \tilde{n} -manifold and \mathcal{U} an open set in $\mathbb{R}^{\tilde{n}}$. The map $\varphi: \mathcal{U} \rightarrow \mathcal{S}$ is called a *local parameterization* of \mathcal{S} if $\text{D}\varphi(\mathbf{u}): \mathbb{R}^{\tilde{n}} \rightarrow T_{\varphi(\mathbf{u})}\mathcal{S}$ is injective for all $\mathbf{u} \in \mathcal{U}$, where $T_{\varphi(\mathbf{u})}\mathcal{S}$ is the tangent space of \mathcal{S} at $\varphi(\mathbf{u}) \in \mathcal{S}$. φ is called a *global parameterization* of \mathcal{S} if it is, additionally, onto.*

Let φ be a global parameterization of \mathcal{S} and $\text{N}^\varphi := \text{N}^{\mathcal{S}} \circ \varphi$. By a direct application of the chain rule, (116) can be written as

$$\kappa(\varphi(\mathbf{u}), \mathbf{w}) = \frac{\langle \mathbf{w}, \text{DN}^\varphi(\mathbf{u})\mathbf{v} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle}, \quad (118)$$

³It was claimed that if \mathcal{S}_ℓ is non-convex, there exists a point \mathbf{s}_0 on the loss surface \mathcal{S}_ℓ such that no hyperplane supports \mathcal{S}_ℓ at \mathbf{s}_0 . The non-convexity of a set by itself is not sufficient to make such a claim; the continuity of the loss ℓ is required.

where \mathbf{v} is such that $D\varphi(\mathbf{u})\mathbf{v} = \mathbf{w}$. The existence of such a \mathbf{v} is guaranteed by the fact that $D\varphi$ is injective and $\dim \mathbb{R}^{\tilde{n}} = \dim T_{\varphi(\mathbf{u})}\mathcal{S} = \tilde{n}$.

Theorem 30. First we show that \mathcal{S}_ℓ is a \tilde{n} -manifold. Consider the map $\tilde{\ell} : \text{int } \tilde{\Delta}_n \rightarrow \mathcal{S}_\ell$ and note that $\text{int } \tilde{\Delta}_n$ is trivially a \tilde{n} -manifold. Due to the strict concavity of the Bayes risk, $\tilde{\ell}$ is injective [14] and from Lemmas 8 and 28, $D\tilde{\ell}(\tilde{\mathbf{p}}) : \mathbb{R}^{\tilde{n}} \rightarrow T_{\tilde{\ell}(\tilde{\mathbf{p}})}\mathcal{S}_\ell$ is also injective. Therefore, $\tilde{\ell}$ is an *immersion* [10]. $\tilde{\ell}$ is also *proper* in the sense that the preimage of every compact subset of \mathcal{S}_ℓ is compact. Therefore, $\tilde{\ell}$ is a proper injective immersion, and thus it is an embedding from the \tilde{n} -manifold $\text{int } \tilde{\Delta}_n$ to \mathcal{S}_ℓ (ibid.). Hence, \mathcal{S}_ℓ is a manifold.

Now we prove (117). The map $\tilde{\ell}$ is a global parameterization of \mathcal{S}_ℓ . In fact, from Lemma 8, $D\tilde{\ell}(\tilde{\mathbf{p}})$ has rank \tilde{n} , for all $\tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n$, which implies that $D\tilde{\ell}(\tilde{\mathbf{p}})$ is onto from $\mathbb{R}^{\tilde{n}}$ to $T_{\tilde{\ell}(\tilde{\mathbf{p}})}\mathcal{S}_\ell$. Therefore, given $\mathbf{w} \in T_{\tilde{\ell}(\tilde{\mathbf{p}})}\mathcal{S}_\ell$, there exists $\mathbf{v} \in \mathbb{R}^{\tilde{n}}$ such that $\mathbf{w} = D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}$. Furthermore, Lemma 8 implies that $N^{\tilde{\ell}}(\tilde{\mathbf{p}}) = \mathbf{p}$, since $\langle \mathbf{p}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_n^\top$. Substituting $N^{\tilde{\ell}}$ into (118) yields

$$\begin{aligned} \kappa_\ell(\tilde{\ell}(\tilde{\mathbf{p}}), \mathbf{w}) &= \frac{\mathbf{v}^\top (D\tilde{\ell}(\tilde{\mathbf{p}}))^\top \begin{bmatrix} I_{\tilde{n}}, \\ \mathbf{1}_{\tilde{n}} \end{bmatrix} \mathbf{v}}{\langle D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}, D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \rangle}, \\ &= \frac{\mathbf{v}^\top H\tilde{L}_\ell(\tilde{\mathbf{p}}) \begin{bmatrix} X_{\tilde{\mathbf{p}}}^\top, & -\tilde{\mathbf{p}} \end{bmatrix} \begin{bmatrix} I_{\tilde{n}} \\ \mathbf{1}_{\tilde{n}} \end{bmatrix} \mathbf{v}}{\langle D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v}, D\tilde{\ell}(\tilde{\mathbf{p}})\mathbf{v} \rangle}, \\ &= \frac{\mathbf{v}^\top H\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{v}}{\mathbf{v}^\top H\tilde{L}_\ell(\tilde{\mathbf{p}}) \begin{bmatrix} X_{\tilde{\mathbf{p}}}^\top, & -\tilde{\mathbf{p}} \end{bmatrix} \begin{bmatrix} X_{\tilde{\mathbf{p}}} \\ -\tilde{\mathbf{p}}^\top \end{bmatrix} H\tilde{L}_\ell(\tilde{\mathbf{p}})\mathbf{v}}. \end{aligned} \quad (119)$$

Setting $\mathbf{u} = (-H\tilde{L}_\ell(\tilde{\mathbf{p}}))^{-\frac{1}{2}}\mathbf{v} / \|(-H\tilde{L}_\ell(\tilde{\mathbf{p}}))^{-\frac{1}{2}}\mathbf{v}\|$ in (119) gives the desired result. \square

I Classical Mixability Revisited

In this appendix, we provide a more concise proof of the necessary and sufficient conditions for the convexity of the superprediction set [14].

Theorem 32. Let $\ell : \Delta_n \rightarrow [0, +\infty]^n$ be a strictly proper loss whose Bayes risk is twice differentiable on $]0, +\infty[^n$. The following points are equivalent;

- (i) $\forall \tilde{\mathbf{p}} \in \text{int } \tilde{\Delta}_n, \eta H\tilde{L}_\ell(\tilde{\mathbf{p}}) \succeq H\tilde{L}_{\log}(\tilde{\mathbf{p}})$.
- (ii) $e^{-\eta \mathcal{S}_\ell} = \bigcap_{\mathbf{p} \in \Delta_n} \mathcal{H}_{\tau(\mathbf{p}), 1} \cap [0, +\infty[^n$, where $\tau(\mathbf{p}) := \mathbf{p} \odot e^{\eta \ell(\mathbf{p})}$.
- (iii) $e^{-\eta \mathcal{S}_\ell}$ is convex.

Proof. We already showed (i) \implies (ii) \implies (iii) in the proof of Theorem 7.

We now show (iii) \implies (i). Since $e^{-\eta \mathcal{S}_\ell}$ is convex, any point $\mathbf{s} \in \text{bd } e^{-\eta \mathcal{S}_\ell}$ is supported by a hyperplane [7, Lem. A.4.2.1]. Since $\mathbf{u} \mapsto e^{-\eta \mathbf{u}}$ is a homeomorphism, it maps boundaries to boundaries. From this and Lemma 17, $\text{bd } e^{-\eta \mathcal{S}_\ell} = e^{-\eta \mathcal{S}_\ell}$. Thus, for $\mathbf{p} \in \text{ri } \Delta_n$, there exists a unit-norm vector $\mathbf{u} \in \mathbb{R}^n$ such that for all $\mathbf{s} \in \mathcal{S}_\ell$ it either holds that $\langle \mathbf{u}, e^{-\eta \ell(\mathbf{p})} \rangle \leq \langle \mathbf{u}, e^{-\eta \mathbf{s}} \rangle$; or $\langle \mathbf{u}, e^{-\eta \ell(\mathbf{p})} \rangle \geq \langle \mathbf{u}, e^{-\eta \mathbf{s}} \rangle$. It is easy to see that it is the latter case that holds, since we can choose $\mathbf{s} = \ell(\mathbf{r}) + c\mathbf{1} \in \mathcal{S}_\ell$, for $\mathbf{r} \in \Delta_n$, and make $\langle \mathbf{u}, e^{-\eta \mathbf{s}} \rangle$ arbitrarily small by making $c \in \mathbb{R}$ large. Therefore, $\forall \mathbf{r} \in \text{ri } \Delta_n, \langle \mathbf{u}, e^{-\eta \tilde{\ell}(\tilde{\mathbf{p}})} \rangle = \langle \mathbf{u}, e^{-\eta \ell(\mathbf{p})} \rangle \geq \langle \mathbf{u}, e^{-\eta \ell(\mathbf{r})} \rangle = \langle \mathbf{u}, e^{-\eta \tilde{\ell}(\tilde{\mathbf{r}})} \rangle$ and $\tilde{\mathbf{p}}$ is a critical point of the function $f(\tilde{\mathbf{r}}) := \langle \mathbf{u}, e^{-\eta \tilde{\ell}(\tilde{\mathbf{r}})} \rangle$ on $\text{int } \tilde{\Delta}_n$. This implies that $\nabla f(\tilde{\mathbf{p}}) = \mathbf{0}_{\tilde{n}}$; that is, $-\eta \langle \mathbf{u}, \text{diag}(e^{-\eta \tilde{\ell}(\tilde{\mathbf{p}})}) D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = -\eta \langle \text{diag}(e^{-\eta \tilde{\ell}(\tilde{\mathbf{p}})}) \mathbf{u}, D\tilde{\ell}(\tilde{\mathbf{p}}) \rangle = \mathbf{0}_n^\top$. From Lemma 8, there exists $\lambda \in \mathbb{R}$ such that $\text{diag}(e^{-\eta \tilde{\ell}(\tilde{\mathbf{p}})}) \mathbf{u} = \lambda \mathbf{p}$. Therefore, $\mathbf{u} = \lambda \mathbf{p} \odot e^{\eta \tilde{\ell}(\tilde{\mathbf{p}})}$, where $\lambda = \|\mathbf{p} \odot e^{\eta \tilde{\ell}(\tilde{\mathbf{p}})}\|^{-1}$,

since $\|u\| = 1$. For $v \in \mathbb{R}^{n-1}$, let $\tilde{\alpha}^t := \tilde{p} + tv$, where $t \in \{s : \tilde{p} + sv \in \text{int } \tilde{\Delta}_n\}$. Since f is twice differentiable and attains a maximum at \tilde{p} ,

$$\begin{aligned} 0 &\geq \frac{1}{\lambda\eta} \frac{d^2}{dt^2} f \circ \tilde{\alpha}^t \Big|_{t=0} = \frac{1}{\lambda} \frac{d}{dt} \left\langle u, \text{diag } e^{-\eta \tilde{\ell}(\tilde{\alpha}^t)} D\tilde{\ell}(\tilde{\alpha}^t) v \right\rangle \Big|_{t=0}, \\ &= \frac{d}{dt} \left\langle p \odot e^{\eta \tilde{\ell}(\tilde{p})}, \text{diag } e^{-\eta \tilde{\ell}(\tilde{\alpha}^t)} D\tilde{\ell}(\tilde{p}) v \right\rangle \Big|_{t=0} + \frac{d}{dt} \left\langle p, D\tilde{\ell}(\tilde{\alpha}^t) v \right\rangle \Big|_{t=0}, \\ &= \eta v^T H\tilde{L}_\ell(\tilde{p}) (H\tilde{L}_{\log}(\tilde{p}))^{-1} H\tilde{L}_\ell(\tilde{p}) v - v^T H\tilde{L}_\ell(\tilde{p}) v, \end{aligned} \quad (120)$$

where in the second equality we substituted u by $\lambda p \odot e^{\eta \tilde{\ell}(\tilde{p})}$ and in (120) we used (5) and (6) from Lemma 9. Note that by the assumptions on ℓ it follows that the Bayes risk \tilde{L}_ℓ is strictly concave [14, Lemma 6] and $-H\tilde{L}_\ell(\tilde{p})$ is symmetric negative-definite. In particular, $H\tilde{L}_\ell(\tilde{p})$ is invertible. Setting $\hat{v} := H\tilde{L}_\ell(\tilde{p})v$ in (120) yields

$$0 \geq \eta \hat{v} (H\tilde{L}_{\log}(\tilde{p}))^{-1} \hat{v} - \hat{v} (H\tilde{L}_\ell(\tilde{p}))^{-1} \hat{v}.$$

Since $v \in \mathbb{R}^{n-1}$ was chosen arbitrarily, $(H\tilde{L}_\ell(\tilde{p}))^{-1} \succeq \eta (H\tilde{L}_{\log}(\tilde{p}))^{-1}, \forall \tilde{p} \in \text{int } \tilde{\Delta}_n$. This is equivalent to the condition $\forall \tilde{p} \in \text{int } \tilde{\Delta}_n, \eta H\tilde{L}_\ell(\tilde{p}) \succeq H\tilde{L}_{\log}(\tilde{p})$. \square

J An Experiment on Football Prediction Dataset

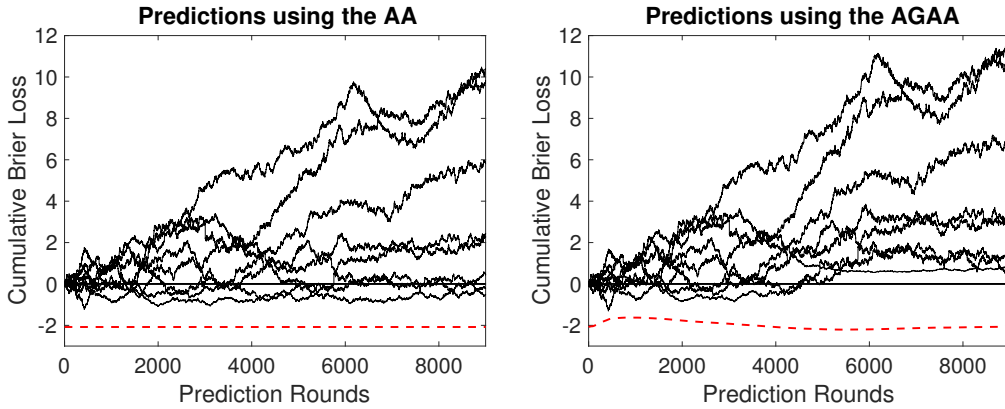


Figure 1: The figure corresponds to the 2005/2006, 2006/2007, 2007/2008, and 2008/2009 seasons. The solid lines represent, at each round t , the difference between the cumulative losses of the experts and that of the learner who uses either the AA (left) or the AGAA (right); that is, $\text{Loss}_{\theta}^{\ell_{\text{Brier}}}(t) - \text{Loss}_{\mathfrak{M}}^{\ell_{\text{Brier}}}(t)$, for $\mathfrak{M} \in \{\text{AA}, \text{AGAA}\}$. The red dashed lines represent the negative of the regret bound in (12) with respect to the best expert θ^* ; that is, $-R_{\ell_{\text{Brier}}}^S - \Delta R_{\theta^*}(t)$ at each round t .

J.1 Testing the AGAA

To test the AGAA empirically, we used prediction data⁴ from the British football leagues, including the Premier Leagues, Championships, Leagues 1-2, and Conferences. The first dataset contains predictions for the 2005/2006, 2006/2007, 2007/2008, and 2008/2009 seasons, matching the dataset used in [15]. The second dataset contains predictions for the 2009/2010, 2010/2011, 2011/2012, and 2012/2013 seasons. For this set, we considered predictions from 9 bookmakers; Bet365, Bet&Win, Blue Square, Gamebookers, Interwetten, Ladbrokes, Stan James, VC Bet, and William Hill.

On each dataset, we compared the performance of the AGAA with that of the AA using the Brier score (the Brier loss is 1-mixable). For the AGAA, we chose β^t according to Theorem 19 with $v^t := -\frac{1}{2^t} \sum_{s=1}^t \ell_{x^s}(A^s)$ and we set $\Phi = S$, i.e. the Shannon entropy. The results in Figure 1 [resp. Figure 2] correspond to the seasons from 2005 to 2009 [resp. 2009 to 2013]. For fair comparison

⁴The data was collected from <http://www.football-data.co.uk/>.

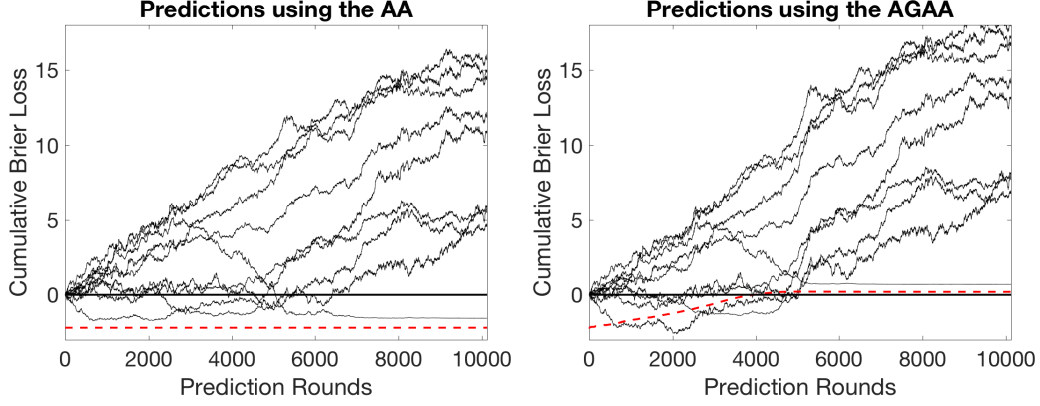


Figure 2: The figure corresponds to the 2009/2010, 2010/2011, 2011/2012, and 2012/2013 seasons. The solid lines represent, at each round t , the difference between the cumulative losses of the experts and that of the learner who uses either the AA (left) or the AGAA (right); that is, $\text{Loss}_\theta^{\ell_{\text{Brier}}}(t) - \text{Loss}_{\mathfrak{M}}^{\ell_{\text{Brier}}}(t)$, for $\mathfrak{M} \in \{\text{AA}, \text{AGAA}\}$. The red dashed lines represent the negative of the regret bound in (12) with respect to the best expert θ^* ; that is, $-R_{\ell_{\text{Brier}}}^S - \Delta R_{\theta^*}(t)$ at each round t .

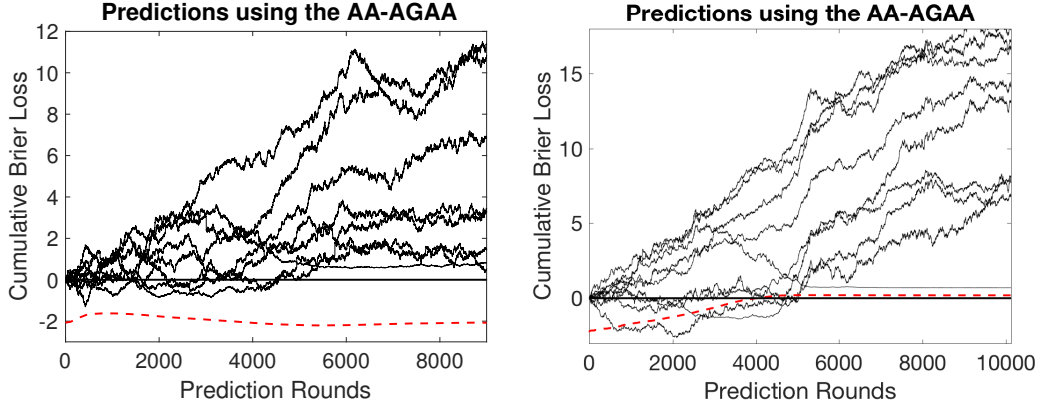


Figure 3: The figure on the left [resp. right] hand side corresponds to the football seasons from 2005 to 2009 [resp. 2009 to 2013]. The solid lines represent, at each round t , the difference between the cumulative losses of the experts and that of the learner using the AA-AGAA meta algorithm (refer to text); that is, $\text{Loss}_\theta^{\ell_{\text{Brier}}}(t) - \text{Loss}_{\text{AA-AGAA}}^{\ell_{\text{Brier}}}(t)$. The red dashed lines represent the negative of the regret bound in (12) with respect to the best expert θ^* ; that is, $-R_{\ell_{\text{Brier}}}^S - \Delta R_{\theta^*}(t)$ at each round t .

with the results of Vovk [15], we 1) used the same substitution function as [15]; 2) used the same method for converting odds to probabilities; and 3) sorted the data first by date then by league and then by name of the host team (For more detail see [15]).

In all figures the solid lines represent, at each round t , the difference between the cumulative losses of the experts and that of the learners; that is, $\text{Loss}_\theta^{\ell_{\text{Brier}}}(t) - \text{Loss}_{\mathfrak{M}}^{\ell_{\text{Brier}}}(t)$, for $\mathfrak{M} = \text{AA}, \text{AGAA}$. The red dashed lines represent the negative of the regret bound in (12) with respect to the best expert θ^* ; that is, $-R_{\ell_{\text{Brier}}}^S - \Delta R_{\theta^*}(t) = -R_{\ell_{\text{Brier}}}^S - \sum_{s=1}^{t-1} (v_\theta^s - \langle v_\theta^s, q^s \rangle)$ at each round t , where (q^s) are the distributions over experts.

From Figures 1 and 2 it can be seen that the learners using the AGAA perform better than the best expert (and better than the AA) at the end of the games.

J.2 Testing a AA-AGAA Meta-Learner

Consider the algorithm (referred to as AA-AGAA) that takes the outputs of the AGAA and the AA as in the previous section and aggregates them using the AA to yield a *meta prediction*. The worst case

regret of this algorithm is guaranteed not to exceed that of the original AA and AGAA by more than $\eta^{-1} \log 2$ for an η -mixable loss. Figure 3 shows the results for this algorithm for the same datasets as the previous section. The AA-AGAA still achieves a negative regret at the end of the game.

References

- [1] Ravi P Agarwal, Maria Meehan, and Donal O'Regan. *Fixed point theory and applications*, volume 141. Cambridge university press, 2001.
- [2] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [3] Dennis S. Bernstein. *Matrix Mathematics: Theory, Facts, and Formulas*. Princeton University Press, second edition, 2011.
- [4] Jonathan M. Borwein, Jon D. Vanderwerff, et al. *Convex functions: constructions, characterizations and counterexamples*, volume 109. Cambridge University Press Cambridge, 2010.
- [5] Alexey Chernov, Yuri Kalnishkan, Fedor Zhdanov, and Vladimir Vovk. Supermartingales in prediction with expert advice. *Theoretical Computer Science*, 411(29-30):2647–2669, 2010.
- [6] A. Philip Dawid. The geometry of proper scoring rules. *Annals of the Institute of Statistical Mathematics*, 59(1):77–93, 2007.
- [7] J-B. Hiriart-Urruty and C. Lemaréchal. *Fundamentals of convex analysis.*, 2001.
- [8] Nelson Merentes and Kazimierz Nikodem. Remarks on strongly convex functions. *Aequationes mathematicae*, 80(1):193–199, 2010.
- [9] 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.
- [10] Joel W. Robbin and Dietmar A. Salamon. Introduction to differential geometry, eth, lecture notes. *ETH, Lecture Notes, preliminary version, January*, 2011.
- [11] R. Tyrrell Rockafellar. *Convex analysis*. Princeton University Press, Princeton, NJ, 1997.
- [12] Walter Rudin et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- [13] John A. Thorpe. *Elementary Topics in Differential Geometry*. Springer Science & Business Media, 1994.
- [14] 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.
- [15] Vladimir Vovk and Fedor Zhdanov. Prediction with expert advice for the brier game. *Journal of Machine Learning Research*, 10(Nov):2445–2471, 2009.
- [16] Robert C. Williamson, Elodie Vernet, and Mark D. Reid. Composite multiclass losses. *Journal of Machine Learning Research*, 17:223:1–223:52, 2016.