

Learning Classifiers with Fenchel-Young Losses: Generalized Entropies, Margins, and Algorithms

Mathieu Blondel¹, André F. T. Martins², and Vlad Niculae³

¹ *NTT Communication Science Laboratories, Kyoto, Japan*
mathieu@mb Blondel.org

² *Unbabel & Instituto de Telecomunicações, Lisbon, Portugal*
andre.martins@unbabel.com

³ *Cornell University, Ithaca, NY*
vlad@vene.ro

Abstract

We study in this paper Fenchel-Young losses, a generic way to *construct* convex loss functions from a convex regularizer. We provide an in-depth study of their properties in a broad setting and show that they unify many well-known loss functions. When constructed from a generalized entropy, which includes well-known entropies such as Shannon and Tsallis entropies, we show that Fenchel-Young losses induce a predictive probability distribution and develop an efficient algorithm to compute that distribution for separable entropies. We derive conditions for generalized entropies to yield a distribution with sparse support and losses with a separation margin. Finally, we present both primal and dual algorithms to learn predictive models with generic Fenchel-Young losses.

1 Introduction

Loss functions are a cornerstone of statistics and machine learning: They measure the difference, or “loss,” between a ground-truth label and a prediction. For this reason, much work has been devoted to designing loss functions [17, 26, 37, 23] or to studying their theoretical properties [53, 30, 57]. Some loss functions, such as the hinge loss of support vector machines, are intimately connected to the notion of separation margin—a prevalent concept in statistical learning theory, which has been used to prove the famous perceptron mistake bound [47] and many other generalization bounds [55, 49]. For probabilistic classification, proper scoring rules [29, 28] provide a generic way to construct loss functions; however, they expect predictions to lie in the probability simplex and hence typically need to be composed with an invertible link function [57], which can render them non-convex [14]. Loss functions are often tightly connected to the statistical model they learn. For instance, the framework of generalized linear models [40] provides a wealth of loss functions for learning conditional probabilistic models in the exponential family [3]. A well-known instance of this framework is the logistic loss, and its associated “softmax distribution.” The logistic loss has many appealing properties, notably, its well-known relationship with the principle of maximum entropy. Yet, the logistic loss lacks a separation margin. Martins & Astudillo [37] were the first to explicitly seek a **sparse distribution** and to derive an associated loss. This “sparsemax” loss is differentiable and, simultaneously, has a natural separation margin. Despite much progress, a general understanding of how to construct differentiable convex losses (potentially with a margin) is still lacking.

Towards this goal, this paper studies and extends Fenchel-Young losses, recently proposed for structured prediction [43]. We show that Fenchel-Young losses provide a generic and principled way to **construct** a loss function with an associated predictive probability distribution. We further show that there is a tight and fundamental relation between generalized entropies, margins, and sparse probability distributions. In summary, we make the following contributions.

- We provide an in-depth study of Fenchel-Young losses, deriving their properties and showing that they unify many well-known loss functions, including the logistic and sparsemax losses (§2).
- When based on a generalized entropy, which includes well-known instances such as the Shannon, Tsallis, and Rényi entropies, we show that Fenchel-Young losses induce a predictive probability distribution. We extend this view to a mean parametrization, useful for structured prediction (§3).

- We derive conditions for generalized entropies to yield distributions with sparse support and losses with a separation margin. For separable entropies, we provide a simple formula that relates the margin with analytic properties of the entropy function (§4).
- Finally, we show that, despite their generality, Fenchel-Young losses can be used to learn (generalized) linear models very simply, using both primal and dual algorithms. For separable entropies, we show that computing their associated probability distribution reduces to unidimensional root finding, thus allowing efficient computation of a large class of probability distributions (§5).

Notation. We denote the $(d - 1)$ -dimensional probability simplex by $\Delta^d := \{\mathbf{p} \in \mathbb{R}_+^d : \|\mathbf{p}\|_1 = 1\}$. We denote the convex hull of a set $\mathcal{Y} \subseteq \mathbb{R}^d$ by $\text{conv}(\mathcal{Y})$. We denote the domain of a convex function $\Omega : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ by $\text{dom}(\Omega) := \{\boldsymbol{\mu} \in \mathbb{R}^d : \Omega(\boldsymbol{\mu}) < \infty\}$. The (Fenchel) conjugate of Ω is defined by $\Omega^*(\boldsymbol{\theta}) := \sup_{\boldsymbol{\mu} \in \text{dom}(\Omega)} \langle \boldsymbol{\mu}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{\mu})$.

2 Fenchel-Young loss functions

The Fenchel-Young family of loss functions was recently proposed for structured prediction [43]. We extend it to a broader setting: a general learning problem with input variables $\mathbf{x} \in \mathcal{X}$, ground-truth $\mathbf{y} \in \mathcal{Y} \subseteq \mathbb{R}^d$, and parametrized model $\mathbf{f}_W : \mathcal{X} \rightarrow \mathbb{R}^d$ which produces a score vector $\boldsymbol{\theta} := \mathbf{f}_W(\mathbf{x})$. Let Ω be a lower-semicontinuous (l.s.c.) proper convex regularizer, with $\text{dom}(\Omega) \subseteq \mathbb{R}^d$ a closed set. We consider **regularized prediction functions** [42, 39] of the kind

$$\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) \in \partial\Omega^*(\boldsymbol{\theta}) = \underset{\boldsymbol{\mu} \in \text{dom}(\Omega)}{\text{argmax}} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega(\boldsymbol{\mu}). \quad (1)$$

A typical choice is $\text{dom}(\Omega) = \text{conv}(\mathcal{Y})$, implying that $\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta})$ is the mean prediction under some underlying distribution, as described in §3. For this choice of $\text{dom}(\Omega)$, $\hat{\mathbf{y}}_0(\boldsymbol{\theta}) \in \underset{\mathbf{y} \in \mathcal{Y}}{\text{argmax}} \langle \boldsymbol{\theta}, \mathbf{y} \rangle$, recovering the class with highest score. In this broad scenario, Fenchel-Young losses arise as a natural choice to learn the model parameters W (cf. §5 for primal and dual formulations).

Definition 1 Let Ω be a l.s.c. proper convex function, with $\text{dom}(\Omega) \subseteq \mathbb{R}^d$ a closed set. Let $\mathbf{y} \in \mathcal{Y} \subseteq \text{dom}(\Omega)$ be a ground-truth label and $\boldsymbol{\theta} \in \text{dom}(\Omega^*) = \mathbb{R}^d$ be a vector of prediction scores. The **Fenchel-Young loss** $L_\Omega : \text{dom}(\Omega^*) \times \mathcal{Y} \rightarrow \mathbb{R}_+$ generated by Ω is

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{y}) - \langle \boldsymbol{\theta}, \mathbf{y} \rangle.$$

As their name indicates, these losses are grounded in the **Fenchel-Young inequality** [9, Proposition 3.3.4]

$$\Omega^*(\boldsymbol{\theta}) + \Omega(\boldsymbol{\mu}) \geq \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle \quad \forall \boldsymbol{\theta} \in \text{dom}(\Omega^*), \boldsymbol{\mu} \in \text{dom}(\Omega). \quad (2)$$

Under the assumption on Ω in Definition 1, Danskin's theorem [18] implies that (2) becomes an equality (i.e., the duality gap is zero) if and only if $\boldsymbol{\theta} \in \partial\Omega(\boldsymbol{\mu})$. This implies the following properties.

Proposition 1 *Properties of Fenchel-Young losses*

1. **Non-negativity.** $L_\Omega(\boldsymbol{\theta}; \mathbf{y}) \geq 0$ for any $\boldsymbol{\theta} \in \text{dom}(\Omega^*)$ and $\mathbf{y} \in \mathcal{Y} \subseteq \text{dom}(\Omega)$.
2. **Zero loss.** The loss is zero iff $\mathbf{y} \in \partial\Omega^*(\boldsymbol{\theta})$, i.e., iff $\mathbf{y} \in \underset{\boldsymbol{\mu} \in \text{dom}(\Omega)}{\text{argmax}} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega(\boldsymbol{\mu})$.
3. **Convexity & gradient.** L_Ω is convex and the residual vectors are its subgradients: $\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) - \mathbf{y} \in \partial L_\Omega(\boldsymbol{\theta}; \mathbf{y})$. If Ω is strictly convex, then L_Ω is differentiable and $\nabla L_\Omega(\boldsymbol{\theta}; \mathbf{y}) = \hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) - \mathbf{y}$.
4. **Constant invariance.** For any constant $c \in \mathbb{R}$, $L_{\Omega+c}(\boldsymbol{\theta}; \mathbf{y}) = L_\Omega(\boldsymbol{\theta}; \mathbf{y})$.
5. **Temperature scaling.** For any $t > 0$, $L_{t\Omega}(\boldsymbol{\theta}; \mathbf{y}) = tL_\Omega(\frac{\boldsymbol{\theta}}{t}; \mathbf{y})$ and $\hat{\mathbf{y}}_{t\Omega}(\boldsymbol{\theta}) \in \partial\Omega^*(\frac{\boldsymbol{\theta}}{t})$.

The properties above suggest that the minimization of Fenchel-Young losses attempt to adjust the model to produce predictions $\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta})$ that are close to the target \mathbf{y} , reducing the duality gap.

Table 2 shows examples of well-known loss functions that fall into this category: simple choices of Ω yield the squared loss, perceptron loss, logistic loss, sparsemax loss, as well as their extensions to structured classification; see §3.2 for details. More losses can also be recovered, as discussed next.

Table 1: Examples of Fenchel-Young losses. For multi-class classification, we assume $\mathcal{Y} = \{e_i\}_{i=1}^d$ and the ground-truth is $\mathbf{y} = e_k$, where e_i denotes a standard basis (“one-hot”) vector. For structured classification, we assume that elements of \mathcal{Y} are d -dimensional binary vectors with $d \ll |\mathcal{Y}|$, and we denote by $\text{conv}(\mathcal{Y}) = \{\mathbb{E}_{\mathbf{p}}[Y] : \mathbf{p} \in \Delta^{|\mathcal{Y}|}\}$ the corresponding marginal polytope [56]. We denote by $H^s(\mathbf{p}) := -\sum_i p_i \log p_i$ the Shannon entropy of a distribution $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$.

	$\text{dom}(\Omega)$	$\Omega(\boldsymbol{\mu})$	$\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta})$	$L_\Omega(\boldsymbol{\theta}; \mathbf{y})$
Squared loss	$\mathbb{R}^{ \mathcal{Y} }$	$\frac{1}{2} \ \boldsymbol{\mu}\ ^2$	$\boldsymbol{\theta}$	$\frac{1}{2} \ \mathbf{y} - \boldsymbol{\theta}\ ^2$
Perceptron loss [47]	$\Delta^{ \mathcal{Y} }$	0	$\text{argmax}(\boldsymbol{\theta})$	$\max_i \theta_i - \theta_k$
Logistic loss	$\Delta^{ \mathcal{Y} }$	$-H^s(\boldsymbol{\mu})$	$\text{softmax}(\boldsymbol{\theta})$	$\log \sum_i \exp \theta_i - \theta_k$
Sparsemax loss [37]	$\Delta^{ \mathcal{Y} }$	$\frac{1}{2} \ \boldsymbol{\mu}\ ^2$	$\text{sparsemax}(\boldsymbol{\theta})$	$\frac{1}{2} \ \mathbf{y} - \boldsymbol{\theta}\ ^2 - \frac{1}{2} \ \hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) - \boldsymbol{\theta}\ ^2$
Structured perceptron [15]	$\text{conv}(\mathcal{Y})$	0	$\text{MAP}(\boldsymbol{\theta})$	$\max_{\mathbf{y}'} \langle \boldsymbol{\theta}, \mathbf{y}' \rangle - \langle \boldsymbol{\theta}, \mathbf{y} \rangle$
CRF [35]	$\text{conv}(\mathcal{Y})$	$\min_{\mathbb{E}_{\mathbf{p}}[Y]=\boldsymbol{\mu}} -H^s(\mathbf{p})$	$\text{marginals}(\boldsymbol{\theta})$	$\log \sum_{\mathbf{y}'} \exp \langle \boldsymbol{\theta}, \mathbf{y}' \rangle - \langle \boldsymbol{\theta}, \mathbf{y} \rangle$
SparseMAP [43]	$\text{conv}(\mathcal{Y})$	$\frac{1}{2} \ \boldsymbol{\mu}\ ^2$	$\text{sparseMAP}(\boldsymbol{\theta})$	$\frac{1}{2} \ \mathbf{y} - \boldsymbol{\theta}\ ^2 - \frac{1}{2} \ \hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) - \boldsymbol{\theta}\ ^2$

Cost-sensitive losses. Fenchel-Young losses also include the hinge loss of support vector machines. Indeed, from any classification loss L_Ω , we can construct a cost-sensitive version of it as follows. Define $\Psi(\boldsymbol{\mu}; \mathbf{y}) := \Omega(\boldsymbol{\mu}) - \langle \mathbf{c}_{\mathbf{y}}, \boldsymbol{\mu} \rangle$, where $\mathbf{c}_{\mathbf{y}} \in \mathbb{R}_+^{|\mathcal{Y}|}$ is a fixed cost vector that depends on the groundtruth \mathbf{y} ; for example $\mathbf{c}_{\mathbf{y}} = \mathbf{1} - \mathbf{y}$ corresponds to the 0/1 cost. Then, L_Ψ is a cost-sensitive version of L_Ω , which can be written as $L_{\Psi(\cdot; \mathbf{y})}(\boldsymbol{\theta}; \mathbf{y}) = L_\Omega(\boldsymbol{\theta} + \mathbf{c}_{\mathbf{y}}; \mathbf{y}) = \Omega^*(\boldsymbol{\theta} + \mathbf{c}_{\mathbf{y}}) + \Omega(\mathbf{y}) - \langle \boldsymbol{\theta} + \mathbf{c}_{\mathbf{y}}, \mathbf{y} \rangle$. This construction recovers the multi-class hinge loss ([17]; $\Omega = 0$), the softmax-margin loss ([26]; $\Omega = -H^s$), and the cost-augmented sparsemax ([50, Eq. (13)], [43]; $\Omega = \frac{1}{2} \|\cdot\|^2$).

Relation to Bregman divergences. All the losses in Table 2 can be extended to work over $\mathcal{Y} = \text{dom}(\Omega)$. For example, in the case of the logistic loss, where $\Omega(\mathbf{y}) = -H^s(\mathbf{y})$, allowing $\mathbf{y} \in \Delta^{|\mathcal{Y}|}$ instead of $\mathbf{y} \in \{e_i\}_{i=1}^d$ yields the cross-entropy loss, $L_\Omega(\boldsymbol{\theta}; \mathbf{y}) = \text{KL}(\mathbf{y} \parallel \text{softmax}(\boldsymbol{\theta}))$. More generally, there is a relation between Fenchel-Young losses and Bregman divergences. Recall that the Bregman divergence [11] generated by a strictly convex and differentiable Ω is

$$B_\Omega(\mathbf{y} \parallel \boldsymbol{\mu}) := \Omega(\mathbf{y}) - \Omega(\boldsymbol{\mu}) - \langle \nabla \Omega(\boldsymbol{\mu}), \mathbf{y} - \boldsymbol{\mu} \rangle. \quad (3)$$

In other words, this is the difference at \mathbf{y} between Ω and its linearization around $\boldsymbol{\mu}$. Letting $\boldsymbol{\theta} = \nabla \Omega(\boldsymbol{\mu})$ (i.e., $(\boldsymbol{\theta}, \boldsymbol{\mu})$ is a **dual pair**), we have $\Omega^*(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega(\boldsymbol{\mu})$. Substituting in (3), we get $B_\Omega(\mathbf{y} \parallel \boldsymbol{\mu}) = L_\Omega(\boldsymbol{\theta}; \mathbf{y})$. In other words, Fenchel-Young losses can be viewed as a “mixed-form Bregman divergence” [1, Theorem 1.1] where the argument $\boldsymbol{\mu}$ in (3) is **replaced by its dual point $\boldsymbol{\theta}$** .

3 Categorical distribution induced by Fenchel-Young losses

In the previous section, we presented Fenchel-Young losses in a broad setting. We now restrict to classification and show that Fenchel-Young losses induce a probability distribution over classes.

3.1 Generalized entropies and probability spaces

We first restrict to the case $\mathcal{Y} = \{e_i\}_{i=1}^d$. In this case, we have $(\boldsymbol{\theta}, \boldsymbol{\mu}) = (\mathbf{s}, \mathbf{p})$, where $\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|}$ is a vector of prediction scores and $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$ a vector of probabilities (this is not the case in structured prediction, as will become clear in §3.2). We further assume that $\Omega(\mathbf{p}) = -H(\mathbf{p})$, where H is a **generalized entropy** [29]. We say that H is a generalized entropy if $\text{dom}(H) = \Delta^{|\mathcal{Y}|}$ and H is concave. Using a general concave function $H(\mathbf{p})$ to measure the “uncertainty” in a distribution $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$ dates back to at least [21]. The corresponding Fenchel-Young loss is then $L_{-H}(\mathbf{s}; \mathbf{y})$.

Assumptions: We will occasionally make the following assumptions about H .

A.1. Zero entropy: $H(\mathbf{p}) = 0$ if \mathbf{p} is a delta distribution, i.e., $\mathbf{p} \in \mathcal{Y}$.

A.2. Strict concavity: this implies that $\nabla(-H)^*(s)$ exists.

A.3. Symmetry: $H(p) = H(Pp)$ for any permutation matrix P .

Assumptions A.2 and A.3 imply that H is Schur-concave [4], a common requirement in generalized entropies. This in turn implies assumption A.1, up to a constant. Together these assumptions imply that H can be used as a sensible uncertainty measure, as shown next (cf. Appendix B.1 for a proof).

Proposition 2 *Assumptions A.1–A.3 imply the following:*

1. H is non-negative and uniquely maximized by the uniform distribution $p = \mathbf{1}/|\mathcal{Y}|$;
2. The conjugate $(-H)^*$ is convex, symmetric, and differentiable;
3. $\nabla(-H)^*(Ps) = P\nabla(-H)^*(s)$ for any permutation matrix P ;
4. $\nabla(-H)^*(s)$ is “order-preserving”: If $p = \nabla(-H)^*(s)$, then the coordinates of p and s are sorted the same way, i.e., $s_i > s_j \Rightarrow p_i \geq p_j$ and $p_i > p_j \Rightarrow s_i > s_j$.

A particular case of generalized entropies satisfying assumptions A.1–A.3 are separable functions of the form $H(p) = \sum_{j=1}^{|\mathcal{Y}|} h(p_j)$, where $h : [0, 1] \rightarrow \mathbb{R}_+$ is a non-negative strictly concave function such that $h(0) = h(1) = 0$. However, our framework does not require to restrict to this form.

The prediction function $\hat{y}_\Omega(s) = \nabla(-H)^*(s)$ produces a distribution over classes, which can be seen as a special case of generalized exponential family [29, 24]. Because, in the unstructured setting, $\mathbb{E}_p[Y] = p$, the prediction function $\hat{y}_\Omega(s)$ can also be seen as the mean under a distribution p . We now give examples of generalized entropies, some of them enjoying a closed form for $\nabla(-H)^*(s)$.

Shannon entropy [51]. This is the foundation of information theory, defined as $H^s(p) := -\sum_{j=1}^{|\mathcal{Y}|} p_j \log p_j$. As seen in Table 2, the resulting Fenchel-Young loss L_{-H^s} corresponds to the logistic loss. The associated distribution is the classical softmax (cf. [10, Ex. 3.25] for a derivation):

$$\nabla(-H^s)^*(s) = \text{softmax}(s) := \frac{\exp(s)}{\sum_{j=1}^d \exp(s_j)} > 0.$$

Tsallis α -entropies [54]. These are defined as $H_\alpha^T(p) := k(\alpha - 1)^{-1} \left(1 - \sum_{j=1}^{|\mathcal{Y}|} p_j^\alpha\right)$, where $\alpha \geq 0$ and k is an arbitrary positive constant. These entropies arise as a generalization of the Shannon-Khinchin axioms to non-extensive systems [52] and have numerous scientific applications [25, 38]. For convenience, we set $k = \alpha^{-1}$ for the rest of this paper. Tsallis entropies satisfy assumptions A.1–A.3 and can also be written in separable form:

$$H_\alpha^T(p) := \sum_{j=1}^{|\mathcal{Y}|} h_\alpha(p_j) \quad \text{with} \quad h_\alpha(t) := \frac{t - t^\alpha}{\alpha(\alpha - 1)}.$$

The limit case $\alpha \rightarrow 1$ corresponds to the Shannon entropy. When $\alpha = 2$, we recover the Gini index [27], a popular “impurity measure” for decision trees:

$$H_2^T(p) = \frac{1}{2} \sum_{j=1}^{|\mathcal{Y}|} p_j(1 - p_j) = \frac{1}{2}(1 - \|p\|_2^2). \quad (4)$$

It can be checked that $L_{-H_2^T}$ equals the sparsemax loss [37] (cf. Table 2). The gradient mapping is

$$\nabla(-H_2^T)^*(s) = \text{sparsemax}(s) := \underset{p \in \Delta^{|\mathcal{Y}|}}{\operatorname{argmax}} \langle p, s \rangle + H_2^T(p) = \underset{p \in \Delta^{|\mathcal{Y}|}}{\operatorname{argmin}} \|p - s\|^2,$$

which is the Euclidean projection onto the probability simplex. This recovers the sparsemax transformation introduced in [37]. The resulting distribution has **sparse** support: it may assign exactly zero probability to some classes. The projection can be computed **exactly** in $O(|\mathcal{Y}|)$ time [31, 13, 22, 16].

Another interesting case is $\alpha \rightarrow +\infty$, which gives $H_\infty^T(p) = 0$, hence $L_{-H_\infty^T}$ is the perceptron loss in Table 2. The resulting “argmax” distribution puts all probability mass on the top-scoring classes. In sum, the prediction functions for $\alpha = 1, 2, \infty$ are respectively softmax, sparsemax, and argmax. Tsallis entropies can be seen as a **continuous parametric family** subsuming these important cases.

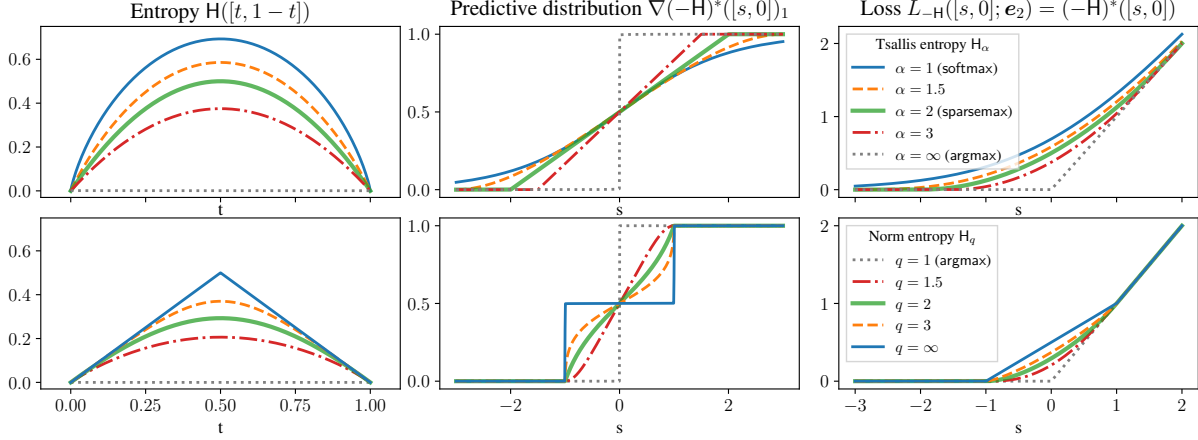


Figure 1: Tsallis and norm entropies (left) along with their prediction functions (middle) and Fenchel-Young losses (right) for the binary case, where $\mathbf{p} = (t, 1 - t) \in \Delta^2$ and $\mathbf{s} = (s, 0) \in \mathbb{R}^2$. Except for softmax, which never exactly reaches 0, all mappings shown on the center lead to **sparse outputs**.

Other entropy families. An interesting class of non-separable entropies are entropies generated by a q -norm, defined as $H_q^N(\mathbf{p}) := 1 - \|\mathbf{p}\|_q$; we call them **norm entropies**. Since norms are convex, these entropies satisfy assumptions A.1–A.3 for $q \geq 1$. They differ from Tsallis entropies in that the norm is not raised to q : a subtle but important difference. We illustrate H , $\nabla(-H)^*$, and L_{-H} for Tsallis and norm entropies in Figure 1. The limit case $q \rightarrow \infty$ is particularly interesting: in this case, we obtain $H_\infty^N = 1 - \|\mathbf{p}\|_\infty$, recovering the Berger-Parker dominance index [7], widely used in ecology to measure species diversity. Other interesting entropies include the **squared-norm entropies** [42] and **Rényi entropies** [45]; cf. Appendix A for more details.

3.2 Structured prediction and mean spaces

We now extend this probabilistic perspective to the structured prediction setting, where \mathcal{Y} is a set of structured objects. In this setting, probability distributions and means no longer coincide, i.e., $\mathbb{E}_{\mathbf{p}}[Y] \neq \mathbf{p}$, and they live in different spaces. We assume that Ω can be written in the form

$$H(\mathbf{p}) = -\Omega(\mathbb{E}_{\mathbf{p}}[Y]) \quad \text{for all } \mathbf{p} \in \Delta^{|\mathcal{Y}|}. \quad (5)$$

That is, $\text{dom}(\Omega) = \text{conv}(\mathcal{Y}) \subseteq \mathbb{R}^d$. In particular, this is the case when using a generalized maximum entropy principle [29, 24]: for all $\boldsymbol{\mu} \in \text{conv}(\mathcal{Y})$ we define $-\Omega(\boldsymbol{\mu}) = \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} H(\mathbf{p})$ s.t. $\mathbb{E}_{\mathbf{p}}[Y] = \boldsymbol{\mu}$.

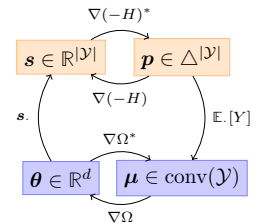
As with exponential families [3], we can alternatively characterize a distribution induced by L_Ω with its mean. Let $\mathbf{s}_\theta \in \mathbb{R}^{|\mathcal{Y}|}$ be a vector of predictions scores, with elements defined by $s_\theta(\mathbf{y}) := \langle \mathbf{y}, \boldsymbol{\theta} \rangle \forall \mathbf{y} \in \mathcal{Y}$. Assuming that Ω is of the form (5), we have the simple but far-reaching identity

$$(-H)^*(\mathbf{s}_\theta) = \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \langle \mathbf{p}, \mathbf{s}_\theta \rangle + H(\mathbf{p}) = \sup_{\boldsymbol{\mu} \in \text{conv}(\mathcal{Y})} \langle \boldsymbol{\mu}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{\mu}) = \Omega^*(\boldsymbol{\theta}),$$

where we used $\langle \mathbf{p}, \mathbf{s}_\theta \rangle = \langle \mathbb{E}_{\mathbf{p}}[Y], \boldsymbol{\theta} \rangle$. This identity connects the conjugates $(-H)^*$ in probability space $\Delta^{|\mathcal{Y}|}$ and Ω^* in mean space $\text{conv}(\mathcal{Y}) \subseteq \mathbb{R}^d$. This is useful in structured prediction, since $\Omega^*(\boldsymbol{\theta})$ just involves a d -dimensional optimization problem instead of a $|\mathcal{Y}|$ -dimensional one for $(-H)^*$. The optimal distribution \mathbf{p}^* is related to $\boldsymbol{\mu}^*$ by $\boldsymbol{\mu}^* = \sum_{\mathbf{y} \in \mathcal{Y}} p^*(\mathbf{y})\mathbf{y}$ and, from Carathéodory's theorem, the support of \mathbf{p}^* contains at most $d \ll |\mathcal{Y}|$ elements. If Ω^* is twice-differentiable, its gradient and Hessian equal the first and second moments under \mathbf{p}^* :

$$\mathbb{E}_{\mathbf{p}^*}[Y] = \nabla \Omega^*(\boldsymbol{\theta}) = \hat{\mathbf{y}}_\Omega(\boldsymbol{\theta}) \in \text{conv}(\mathcal{Y}) \quad \text{and} \quad \text{cov}_{\mathbf{p}^*}[Y] = \nabla^2 \Omega^*(\boldsymbol{\theta}).$$

As indicated in Table 2, $\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta})$ recovers MAP, marginals and sparseMAP for suitable choices of Ω . In the sparseMAP case [43], the mean $\hat{\mathbf{y}}_\Omega(\boldsymbol{\theta})$ along with the distribution \mathbf{p}^* can be computed using conditional gradient or active set algorithms even when $\text{conv}(\mathcal{Y})$ does not enjoy an exact expression. Each iteration of these algorithms only requires a call to a maximum a-posteriori (MAP) oracle, a.k.a. linear minimization oracle (LMO) in the conditional gradient literature [32].



4 Separation margin of Fenchel-Young losses

In this section, we are going to see that the simple assumptions A.1–A.3 about the generalized entropy H are enough to obtain results about the separation margin associated with L_{-H} . The notion of separation margin is well-known in machine learning, lying at the heart of support vector machines and leading to generalization error bounds [55, 49, 30]. We provide a definition and will see next that many other Fenchel-Young losses also have a “margin,” for suitable conditions on H .

Definition 2 Let $L(s; \mathbf{y})$ be a loss function over $\mathbb{R}^{|\mathcal{Y}|} \times \{\mathbf{e}_i\}_{i=1}^{|\mathcal{Y}|}$, where $\mathbf{y} = \mathbf{e}_k$ is a ground truth label. We say that L has a **separation margin property** if there exists $m > 0$ such that:

$$s_k \geq m + \max_{j \neq k} s_j \quad \Rightarrow \quad L(s; \mathbf{y}) = 0. \quad (6)$$

The smallest possible m that satisfies (6) is called the **margin of L** , denoted $\text{margin}(L)$.

The multi-class hinge loss, $L(s; \mathbf{e}_k) = \max\{0, \max_{j \neq k} 1 + s_j - s_k\}$, which we saw in §2 to be an instance of a Fenchel-Young loss, is a famous example with margin 1.

To see that other entropies also induce a margin, we need first to characterize the gradient mappings $\partial(-H)$ and $\nabla(-H)^*$ associated with these entropies. (Note that the former mapping is never single-valued: if \mathbf{s} is in $\partial(-H)(\mathbf{p})$, then so is $\mathbf{s} + c\mathbf{1}$, for any constant $c \in \mathbb{R}$.) We show in Appendix B.3 (Lemma 1) that $\mathbf{s} \in \partial(-H)(\mathbf{e}_k)$ iff $s_k = (-H)^*(\mathbf{s})$. This confirms that the Shannon entropy H^S has $\partial(-H^S)(\mathbf{e}_k) = \emptyset$: since $(-H^S)^*$ is the log-partition function, which upper bounds the maximal score, there is no vector \mathbf{s} satisfying the required condition and hence L_{-H^S} **does not have a margin**.

The next result, proved in Appendix B.2, characterizes more precisely the images of $\nabla(-H)^*$.

Proposition 3 Let H satisfy assumptions A.1–A.3. If $\partial(-H)(\mathbf{p}) \neq \emptyset$ for any $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$, then:

1. The mapping $\nabla(-H)^*$ covers the full simplex, i.e., $\nabla(-H)^*(\mathbb{R}^{|\mathcal{Y}|}) = \Delta^{|\mathcal{Y}|}$;
2. Suppose $\nabla(-H)^*(\mathbf{s}) = \mathbf{p}$ where \mathbf{p} is a relative boundary point, with support $\mathcal{I} := \{1 \leq i \leq |\mathcal{Y}| : p_i > 0\}$. Then, we have $\nabla(-H)^*(\mathbf{s}') = \mathbf{p}$ for any \mathbf{s}' such that $s'_i = s_i$ for $i \in \mathcal{I}$ and $s'_i \leq s_i$ otherwise.

If H is separable with $H(\mathbf{p}) = \sum_i h(p_i)$, the conditions required by Proposition 3 are met if $\lim_{t \rightarrow 0^+} h'(t) < \infty$ and $\lim_{t \rightarrow 1^-} h'(t) > -\infty$. This is the case with Tsallis entropies for $\alpha > 1$, but not Shannon entropy, for which $h'(t) = -1 - \log t$. Functions whose gradient “explode” in the boundary of their domain (hence failing to meet the condition in Proposition 3) are called “essentially smooth” [46]. For those functions, $\nabla(-H)^*$ maps only to the relative interior of $\Delta^{|\mathcal{Y}|}$, never attaining boundary points [56]. This prevents these functions from generating a sparse $\nabla(-H)^*$ or a loss L_{-H} with a margin. Hence, there is a **close link between sparse distributions and losses with a margin**.

We are now ready to state the main result of this section, proved in Appendix B.3: an explicit condition for a Fenchel-Young loss to have a margin, along with an expression to compute it.

Proposition 4 The loss L_{-H} has the separation margin property iff there is a $m > 0$ such that $m\mathbf{e}_k \in \partial(-H)(\mathbf{e}_k)$. The margin of L_{-H} is given by the smallest such m , or equivalently by:

$$\text{margin}(L_{-H}) = \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \frac{H(\mathbf{p})}{1 - \|\mathbf{p}\|_\infty}. \quad (7)$$

Interestingly, the denominator of (7) is the generalized entropy $H_\infty^N(\mathbf{p})$ introduced in §3. As Figure 1 suggests, this entropy provides an upper bound for convex losses with unit margin. The next proposition, which we prove in full detail in Appendix B.4, takes a step even further and provides a remarkably simple formula for twice-differentiable generalized entropies.

Proposition 5 Assume H satisfies the conditions in Proposition 3 and is twice-differentiable on the simplex. Then, for arbitrary $j \neq k$:

$$\text{margin}(L_{-H}) = \nabla_j H(\mathbf{e}_k) - \nabla_k H(\mathbf{e}_k).$$

In particular, if H is separable, i.e., $H(\mathbf{p}) = \sum_{i=1}^{|\mathcal{Y}|} h(p_i)$, where $h : [0, 1] \rightarrow \mathbb{R}_+$ is concave, twice differentiable, and satisfies $h(0) = h(1) = 0$, we have

$$\text{margin}(L_{-H}) = h'(0) - h'(1) = - \int_0^1 h''(t) dt. \quad (8)$$

This provides a geometric characterization of separable entropies and their margins: (8) tells us that only the slopes of h at the two extremities of $[0, 1]$ are relevant to determine its margin. It is also a **constructive result**, since it allows to design an entropy that leads to a loss with a prescribed margin.

Margins of Tsallis and norm entropies. We now apply the results above to obtain margin values of Tsallis and norm entropies (which can be confirmed visually in Figure 1).

Proposition 6 1. Let H_α^\top be the Tsallis α -entropy. Then, for $\alpha \in [0, 1]$, $L_{-H_\alpha^\top}$ does not have the separation margin property. For $\alpha > 1$, we have $\text{margin}(L_{-H_\alpha^\top}) = 1/(\alpha - 1)$.

2. Let H_q^N be the q -norm entropy. Then, for any $q \geq 1$, we have $\text{margin}(L_{-H_q^N}) = 1$.

Proof: As seen in §3, Tsallis entropies are separable with $h(t) = (t - t^\alpha)/(\alpha(\alpha - 1))$. For $\alpha > 1$, the derivative is $h'(t) = (1 - \alpha t^{\alpha-1})/(\alpha(\alpha - 1))$, so we have $h'(0) = 1/(\alpha(\alpha - 1))$ and $h'(1) = -1/\alpha$. Proposition 5 then yields $\text{margin}(L_{-H_\alpha^\top}) = h'(0) - h'(1) = 1/(\alpha - 1)$. Norm entropies are not separable, but their gradient can be expressed as $\nabla H_q^N(\mathbf{p}) = -(\mathbf{p}/\|\mathbf{p}\|_q)^{q-1}$. We thus have $\nabla H_q^N(\mathbf{e}_k) = -\mathbf{e}_k$, so $\text{margin}(H_q^N) = \nabla_j H_q^N(\mathbf{e}_k) - \nabla_k H_q^N(\mathbf{e}_k) = 1$. ■

5 Learning with Fenchel-Young losses

We now turn to learning predictive models with a Fenchel-Young loss L_Ω , for arbitrary $\text{dom}(\Omega)$.

Primal gradient-based algorithms. Let $\mathbf{f}_W : \mathcal{X} \rightarrow \mathbb{R}^d$ be a model parametrized by W . To learn W from training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathcal{Y}$, we minimize the regularized empirical risk

$$\min_W \sum_{i=1}^n L_\Omega(\mathbf{f}_W(\mathbf{x}_i); \mathbf{y}_i) + G(W), \quad (9)$$

where G is a regularizer. This objective is very broad and allows to learn, e.g., neural networks using Fenchel-Young losses. If \mathbf{f}_W is a linear model and G is convex, then (9) is convex. Assuming further that Ω is strongly convex, then L_Ω is smooth and (9) can be solved globally using proximal gradient algorithms [58, 5, 20]. From Proposition 1 (item 3), the gradient of the loss term with respect to $\mathbf{f}_W(\mathbf{x})$ involves the regularized prediction function (1), $\hat{\mathbf{y}}_\Omega = \nabla \Omega^*$. The two key computational ingredients to solve (9) are therefore $\nabla \Omega^*$ and the proximity operator of G , if G is non-differentiable.

Computing $\nabla \Omega^*(\boldsymbol{\theta}) = \arg\max_{\mathbf{y} \in \text{dom}(\Omega)} \langle \boldsymbol{\theta}, \mathbf{y} \rangle - \Omega(\mathbf{y})$ is a concave optimization problem, which enjoys a closed form when $\Omega = -H$ is the negative Shannon entropy or Gini index (§3). More generally, for separable entropies defined over the simplex, we next show that computing $\nabla(-H)^*$ **reduces to unidimensional root finding** (cf. §B.5 for a proof).

Proposition 7 Let $H(\mathbf{p}) = \sum_{j=1}^{|\mathcal{Y}|} h(p_j)$, where h is strictly concave, with derivative h' . Then we have $\nabla(-H)^*(\mathbf{s}) = \hat{\mathbf{p}}(\tau^*)$, where τ^* solves the equation $f(\tau^*) = 0$, and where we defined

$$\hat{\mathbf{p}}(\tau) := (-h')^{-1}(\max\{\mathbf{s} - \tau, -h'(0)\}), \quad f(\tau) := \hat{\mathbf{p}}(\tau)^\top \mathbf{1} - 1.$$

The search interval is $[\tau_{\min}, \tau_{\max}]$, where $\tau_{\min} := \max(\mathbf{s}) + h'(1)$, $\tau_{\max} := \max(\mathbf{s}) + h'(1/|\mathcal{Y}|)$.

An approximate τ such that $|f(\tau)| \leq \epsilon$ can be found in $O(1/\log \epsilon)$ time by, e.g., bisection. Root finding has been employed for related projections [33] but our result applies more broadly (cf. Appendix B.5).

Dual algorithms. When $f_W(\mathbf{x}) = W\mathbf{x}$, we may instead solve the corresponding dual problem

$$\max_{\beta} -D(\beta) \text{ s.t. } \beta_i \in \text{dom}(\Omega) \forall i \in [n], \quad \text{where} \quad D(\beta) := \sum_i \Omega(\beta_i) - \Omega(\mathbf{y}_i) + G^*(V(\beta)) \quad (10)$$

and where we defined $V(\beta) := \sum_{i=1}^n (\mathbf{y}_i - \beta_i) \mathbf{x}_i^\top$. Assuming G is a λ -strongly convex regularizer, given an optimal dual solution β^* , we retrieve the optimal primal solution by $W^* = \nabla G^*(V(\beta^*))$. We can solve (10) using stochastic block coordinate ascent. At every iteration, we pick $i \in [n]$ and follow [50, Option I] to solve an approximate sub-problem w.r.t. β_i . By using the specific form of L_Ω losses, we are able to formulate the update w.r.t. β_i as a proximity operator. Namely, we perform $\beta_i \leftarrow \text{prox}_{\frac{1}{\sigma_i}\Omega}(\mathbf{v}_i/\sigma_i)$, where $\mathbf{v}_i := \nabla G^*(V(\beta))\mathbf{x}_i + \sigma_i\beta_i$, $\sigma_i := \frac{\|\mathbf{x}_i\|^2}{\lambda}$ and where we defined

$$\text{prox}_{\tau\Omega}(\mathbf{x}) := \underset{\mathbf{y} \in \text{dom}(\Omega)}{\text{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 + \tau\Omega(\mathbf{y}). \quad (11)$$

Note that if $\text{prox}_{\tau\Omega}$ is sparse, **then so are the dual variables** $\{\beta_i\}_{i=1}^n$. Sparsity in the dual variables is useful when kernelizing models, as it makes computing predictions more efficient. In many cases, (11) can be solved in closed form; see §B.7 for examples and for a detailed derivation. However, computing $\text{prox}_{\tau\Omega}$ can sometimes be more challenging than $\nabla\Omega^*$. For instance, when Ω is the Shannon negative entropy $-\mathcal{H}^s$, $\nabla\Omega^*$ enjoys a closed form but not $\text{prox}_{\tau\Omega}$. For such cases, assuming $\text{dom}(\Omega)$ is a compact set, another option to solve (10) is the block Frank-Wolfe algorithm [34].

6 Related work

Proper scoring rules (a.k.a. proper losses) [29, 28, 57] are a key tool in the context of probability forecasting [19] to express the goodness-of-fit (lower is better) between a ground truth $\mathbf{y} \in \mathcal{Y}$ and a vector of predictions $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$ (note that \mathbf{p} must be in the simplex). The fact that any scoring rule induces a generalized entropy is well-known [21, 29]. The reverse mapping has also been studied. Let \mathcal{H} be a generalized entropy. Then we can construct a proper scoring rule $S_{-\mathcal{H}}: \Delta^{|\mathcal{Y}|} \times \mathcal{Y} \rightarrow \mathbb{R}$ as follows [48, 28, 44]

$$S_\Omega(\mathbf{p}; \mathbf{y}) := \langle \nabla\Omega(\mathbf{p}), \mathbf{p} \rangle - \Omega(\mathbf{p}) - \nabla\Omega(\mathbf{p})(\mathbf{y}) = \Omega^*(\nabla\Omega(\mathbf{p})) - \nabla\Omega(\mathbf{p})(\mathbf{y}).$$

As an example, choosing the Gini index $\mathcal{H}(\mathbf{p}) = 1 - \|\mathbf{p}\|^2$ generates the **Brier score** [12] $S_{-\mathcal{H}}(\mathbf{p}; \mathbf{y}) = \sum_{\mathbf{y}' \in \mathcal{Y}} ([\mathbf{y} = \mathbf{y}'] - p(\mathbf{y}))^2$. Thus sparsemax and the Brier score share the same generating function. A crucial difference between S_Ω and L_Ω , however, is that S_Ω is **not necessarily convex** in its first argument ([57, Proposition 17] shows that it is in fact quasi-convex) while L_Ω **always is**. In addition, the first argument is **constrained** to $\Delta^{|\mathcal{Y}|}$ for S_Ω , while it is **unconstrained** for L_Ω . For this reason, S_Ω is usually composed with an **invertible** link function [14, 57], while this is not necessary with L_Ω . Finally, a classification loss construction was recently proposed in [23, Proposition 3]. Fenchel-Young losses are broader, supporting e.g. regression, and in fact include [23, Proposition 3] as a special case.

Smoothing techniques. Another way to view our work is through the prism of smoothing techniques [41, 6, 42, 39]. Indeed, if $\text{dom}(\Omega) = \Delta^{|\mathcal{Y}|}$, Ω^* can be seen as a smoothed max operator (see Figure 1, right), with $(t\Omega)^*(\mathbf{s}) \rightarrow \max(\mathbf{s})$ when $t \rightarrow 0$. Smoothing techniques were used extensively in [50] to create smoothed losses. However, these techniques were applied on a per-loss basis, while we propose a generic loss construction, with clear links between smoothing / regularization and the probability distribution over classes induced by $\nabla\Omega^*$.

7 Conclusion

Fenchel-Young losses not only unify many existing loss functions, they offer a principled way to construct new ones from any convex regularizer Ω . When the regularizer is a generalized entropy, we showed that Fenchel-Young losses induce a probability distribution over classes and under precise conditions we derived, enjoy a separation margin: a useful property leading to generalization error bounds [55, 49], and as we showed, closely related to sparse gradient mappings. Importantly, Fenchel-Young losses support regularizers Ω defined over arbitrary domains $\text{dom}(\Omega)$. Leveraging this fact allows to construct loss functions for a large variety of applications.

Acknowledgements

We are grateful to Tim Vieira for pointing us to [24], which kicked off this project. We thank Mário Figueiredo and Nicolas Keriven for reading a draft of this paper and for their insightful comments. AM was partially supported by the European Research Council (ERC StG DeepSPIN 758969) and by the Fundação para a Ciência e Tecnologia through contracts UID/EEA/50008/2013, PTDC/EEI-SII/7092/2014 (LearnBig), and CMUPERI/TIC/0046/2014 (GoLocal).

References

- [1] S. Amari. *Information Geometry and Its Applications*. Springer, 2016.
- [2] K. Ball, E. A. Carlen, and E. H. Lieb. Sharp uniform convexity and smoothness inequalities for trace norms. *Inventiones Mathematicae*, 115(1):463–482, 1994.
- [3] O. Barndorff-Nielsen. *Information and Exponential Families: In Statistical Theory*. John Wiley & Sons, 1978.
- [4] H. H. Bauschke, P. L. Combettes, et al. *Convex analysis and monotone operator theory in Hilbert spaces*, volume 2011. Springer, 2017.
- [5] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [6] A. Beck and M. Teboulle. Smoothing and first order methods: A unified framework. *SIAM Journal on Optimization*, 22(2):557–580, 2012.
- [7] W. H. Berger and F. L. Parker. Diversity of planktonic foraminifera in deep-sea sediments. *Science*, 168(3937):1345–1347, 1970.
- [8] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific Belmont, 1999.
- [9] J. Borwein and A. S. Lewis. *Convex analysis and nonlinear optimization: theory and examples*. Springer Science & Business Media, 2010.
- [10] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [11] L. M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- [12] G. W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- [13] P. Brucker. An $O(n)$ algorithm for quadratic knapsack problems. *Operations Research Letters*, 3(3):163–166, 1984.
- [14] A. Buja, W. Stuetzle, and Y. Shen. Loss functions for binary class probability estimation and classification: Structure and applications. Technical report, University of Pennsylvania, 2005.
- [15] M. Collins. Discriminative training methods for Hidden Markov Models: Theory and experiments with perceptron algorithms. In *Proc. of EMNLP*, 2002.
- [16] L. Condat. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, 158(1-2):575–585, 2016.
- [17] K. Crammer and Y. Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:265–292, 2001.
- [18] J. M. Danskin. The theory of max-min, with applications. *SIAM Journal on Applied Mathematics*, 14(4):641–664, 1966.
- [19] A. P. Dawid. Probability forecasting. *Encyclopedia of Statistical Sciences*, 1986.
- [20] A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Proc. of NIPS*, 2014.
- [21] M. H. DeGroot. Uncertainty, information, and sequential experiments. *The Annals of Mathematical Statistics*, pages 404–419, 1962.
- [22] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In *Proc. of ICML*, 2008.

- [23] J. C. Duchi, K. Khosravi, and F. Ruan. Multiclass classification, information, divergence, and surrogate risk. *Annals of Statistics*, 2018.
- [24] R. Frongillo and M. D. Reid. Convex foundations for generalized MaxEnt models. In *Proc. of AIP*, 2014.
- [25] M. Gell-Mann and C. Tsallis. *Nonextensive Entropy: Interdisciplinary Applications*. Oxford University Press, 2004.
- [26] K. Gimpel and N. A. Smith. Softmax-margin CRFs: Training log-linear models with cost functions. In *Proc. of NAACL*, 2010.
- [27] C. Gini. Variabilità e mutabilità. *Reprinted in Memorie di metodologica statistica (Ed. Pizetti E, Salvemini, T). Rome: Libreria Eredi Virgilio Veschi*, 1912.
- [28] T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [29] P. D. Grünwald and A. P. Dawid. Game theory, maximum entropy, minimum discrepancy and robust Bayesian decision theory. *Annals of Statistics*, pages 1367–1433, 2004.
- [30] Y. Guermeur. Vc theory of large margin multi-category classifiers. *Journal of Machine Learning Research*, 8:2551–2594, 2007.
- [31] M. Held, P. Wolfe, and H. P. Crowder. Validation of subgradient optimization. *Mathematical Programming*, 6(1):62–88, 1974.
- [32] M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proc. of ICML*, 2013.
- [33] W. Krichene, S. Krichene, and A. Bayen. Efficient Bregman projections onto the simplex. In *Proc. of CDC*. IEEE, 2015.
- [34] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *Proc. of ICML*, 2012.
- [35] J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. of ICML*, 2001.
- [36] O. Mangasarian. Pseudo-convex functions. *Journal of the Society for Industrial and Applied Mathematics, Series A: Control*, 3(2):281–290, 1965.
- [37] A. F. Martins and R. F. Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *Proc. of ICML*, 2016.
- [38] A. F. T. Martins, M. A. T. Figueiredo, P. M. Q. Aguiar, N. A. Smith, and E. P. Xing. Nonextensive information theoretic kernels on measures. *Journal of Machine Learning Research*, 10:935–975, 2009.
- [39] A. Mensch and M. Blondel. Differentiable Dynamic Programming for Structured Prediction and Attention. *arXiv preprint arXiv:1802.03676*, 2018.
- [40] J. A. Nelder and R. J. Baker. *Generalized Linear Models*. Wiley Online Library, 1972.
- [41] Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, 2005.
- [42] V. Niculae and M. Blondel. A regularized framework for sparse and structured neural attention. In *Proc. of NIPS*, 2017.
- [43] V. Niculae, A. F. Martins, M. Blondel, and C. Cardie. SparseMAP: Differentiable sparse structured inference. In *Proc. of ICML*, 2018.
- [44] M. D. Reid, R. M. Frongillo, R. C. Williamson, and N. Mehta. Generalized mixability via entropic duality. In *Proc. of COLT*, 2015.
- [45] A. Rényi. On measures of entropy and information. In *Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics, and Probability*, volume 1, pages 547–561, Berkeley, 1961. University of California Press.
- [46] R. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [47] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386, 1958.
- [48] L. J. Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801, 1971.
- [49] B. Schölkopf and A. J. Smola. *Learning With Kernels*. The MIT Press, Cambridge, MA, 2002.

- [50] S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Mathematical Programming*, 155(1):105–145, 2016.
- [51] C. Shannon and W. Weaver. *The Mathematical Theory of Communication*. University of Illinois Press, Urbana, Ill., 1949.
- [52] H. Suyari. Generalization of Shannon-Khinchin axioms to nonextensive systems and the uniqueness theorem for the nonextensive entropy. *IEEE Trans. Information Theory*, 50(8):1783–1787, 2004.
- [53] A. Tewari and P. L. Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8(May):1007–1025, 2007.
- [54] C. Tsallis. Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52:479–487, 1988.
- [55] V. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [56] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [57] R. C. Williamson, E. Vernet, and M. D. Reid. Composite multiclass losses. *Journal of Machine Learning Research*, 2016.
- [58] S. Wright, R. Nowak, and M. A. T. Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7):2479–2493, 2009.
- [59] H. Zou, J. Zhu, and T. Hastie. New multicategory boosting algorithms based on multicategory fisher-consistent losses. *The Annals of Applied Statistics*, 2(4):1290, 2008.

Appendix

A More examples of generalized entropies and illustrations

We now discuss further examples of generalized entropy H . In all these examples, $\nabla(-H)^*$ can typically not be computed in closed form and we must resort to projected gradient solvers such as FISTA [5] to solve it. Figure 2 shows plots of these entropies, as well as their corresponding prediction functions and loss functions (we also repeat the Tsallis and norm entropies defined in §3).

A.1 Rényi β -entropies.

Rényi entropies [45] are defined for any $\beta \geq 0$ as:

$$H_\beta^R(\mathbf{p}) := \frac{1}{1-\beta} \log \sum_{j=1}^{|\mathcal{Y}|} p_j^\beta.$$

Unlike Shannon and Tsallis entropies, Rényi entropies are not separable, with the exception of $\beta \rightarrow 1$, which also recovers Shannon entropy as a limit case. The case $\beta \rightarrow +\infty$ gives $H_\beta^R(\mathbf{p}) = -\log \|\mathbf{p}\|_\infty$. For $\beta \in [0, 1]$, Rényi entropies satisfy assumptions A.1–A.3; for $\beta > 1$, Rényi entropies fail to be concave. They are however pseudo-concave [36], meaning that, for all $\mathbf{p}, \mathbf{q} \in \Delta^{|\mathcal{Y}|}$, $\langle \nabla H_\beta^R(\mathbf{p}), \mathbf{q} - \mathbf{p} \rangle \leq 0$ implies $H_\beta^R(\mathbf{q}) \leq H_\beta^R(\mathbf{p})$. This implies, among other things, that points $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$ with zero gradient are maximizers of $\langle \mathbf{p}, \boldsymbol{\theta} \rangle + H_\beta^R(\mathbf{p})$, which allows to compute the predictive distribution $\nabla(-H_\beta^R)^*$ with gradient-based methods. This predictive distribution and the corresponding loss function are shown in Figure 2.

A.2 Squared q -norm entropies

Inspired by [42], as a simple extension of the Gini index (4), we consider the following generalized entropy based on squared q -norms:

$$H_q^{\text{sq}}(\mathbf{p}) := \frac{1}{2}(1 - \|\mathbf{p}\|_q^2) = \frac{1}{2} - \frac{1}{2} \left(\sum_{j=1}^{|\mathcal{Y}|} p_j^q \right)^{\frac{2}{q}} \geq 0 \quad \forall \mathbf{p} \in \Delta^{|\mathcal{Y}|}.$$

Contrary to [42], we include the constant term $\frac{1}{2}$ to ensure positivity.

For $q \in (1, 2]$, it is known that the squared q -norm is strongly convex w.r.t. $\|\cdot\|_q$ [2], implying that $(-H_q^{\text{sq}})^*$ is smooth. Although $\nabla(-H_q^{\text{sq}})^*(\mathbf{s}) = \operatorname{argmax}_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \langle \mathbf{p}, \mathbf{s} \rangle + H_q^{\text{sq}}(\mathbf{p})$ can no longer be solved in closed form for $q \in (1, 2)$, it can be solved efficiently using projected gradient descent.

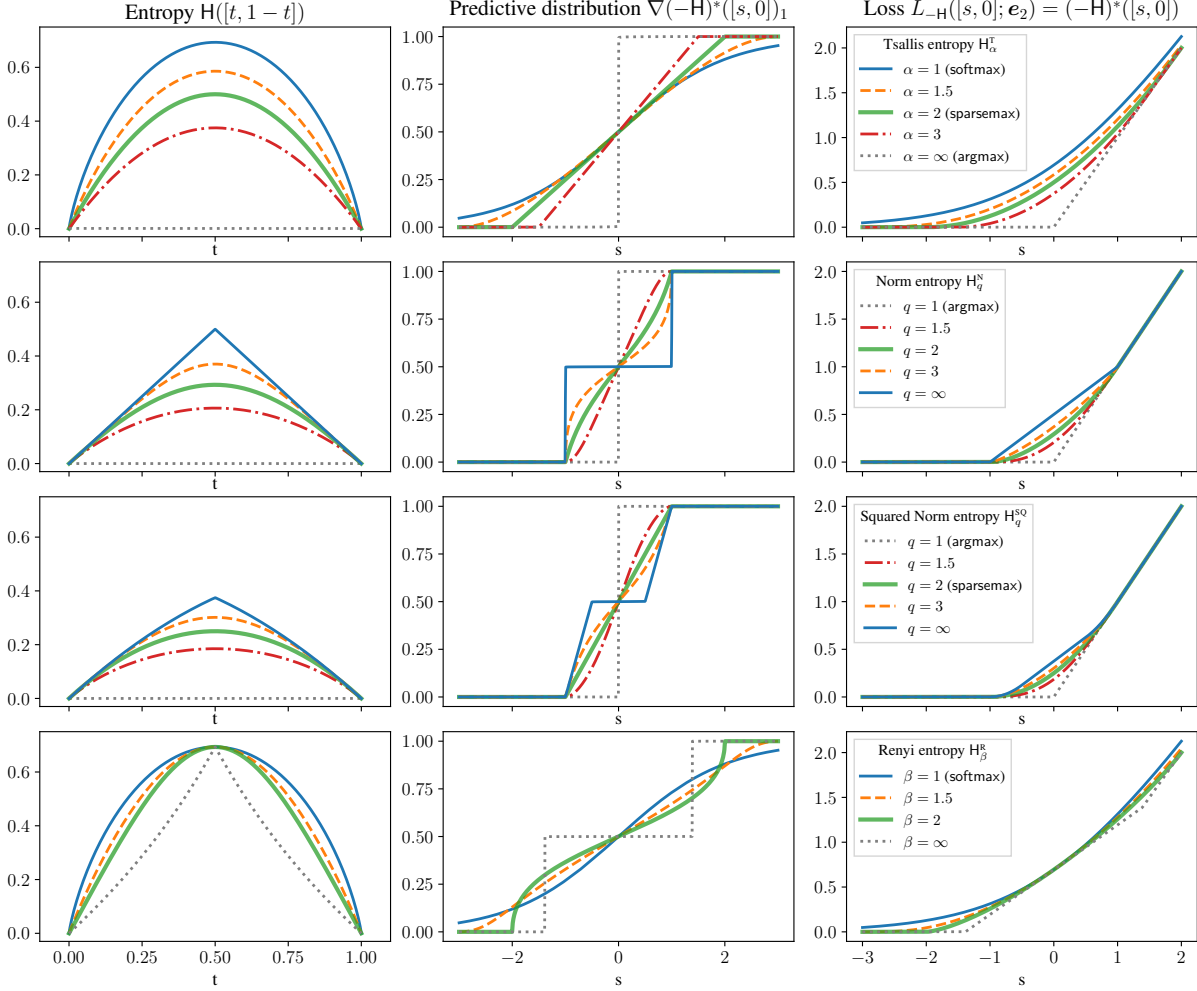


Figure 2: More examples of entropies (left) along with their prediction functions (middle) and Fenchel-Young losses (right) for the binary case, where $\mathbf{p} = (t, 1 - t) \in \Delta^2$ and $\mathbf{s} = (s, 0) \in \mathbb{R}^2$.

A.3 Pairwise hinge entropies

We now turn to the following optimization problem, dating back to [21] (see also [29, 23]), which provides a *reverse* construction from a loss $L: \mathbb{R}^{|\mathcal{Y}|} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ to an entropy $H_L: \Delta^{|\mathcal{Y}|} \rightarrow \mathbb{R}_+$:

$$H_L(\mathbf{p}) := \inf_{\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|}} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}) L(\mathbf{s}; \mathbf{y}). \quad (12)$$

This is the infimum of a linear and thus concave function of \mathbf{p} . Hence, by Danskin's theorem [18] [8, Proposition B.25], $H_L(\mathbf{p})$ is concave (note that this is true even if L is not convex).

The following proposition shows how to compute the entropy generated by pairwise hinge losses.

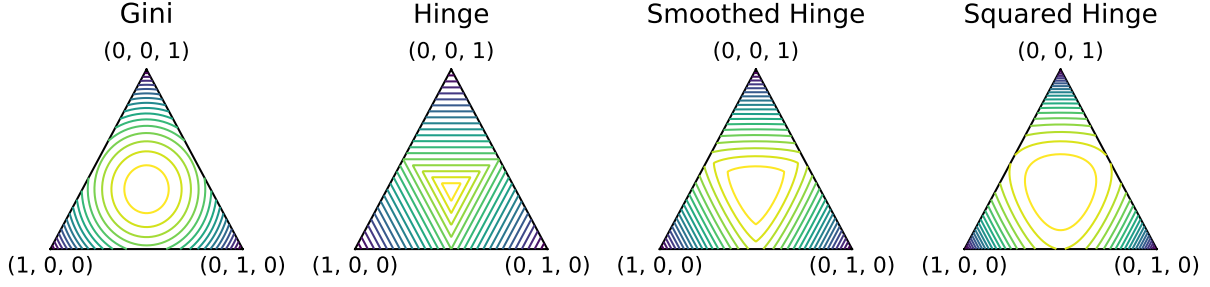


Figure 3: Contour comparison of Gini index H_2^T and pairwise hinge entropies H_L ; cf. Proposition 8.

Proposition 8 *Entropy generated by pairwise hinge losses*

Let $L(\mathbf{s}; \mathbf{e}_k) := \sum_{j \neq k} \phi(s_j - s_k)$ and $\tau(\mathbf{p}) := \min_j 1 - p_j$. Then, for all $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$:

1. If $\phi(t)$ is the hinge function $[1 + t]_+$ then $H_L(\mathbf{p}) = \tau(\mathbf{p})|\mathcal{Y}|$.
2. If $\phi(t)$ is the smoothed hinge function (20), then $H_L(\mathbf{p}) = \frac{-\tau(\mathbf{p})^2}{2} \sum_{j=1}^{|\mathcal{Y}|} \frac{1}{1-p_j} + \tau(\mathbf{p})|\mathcal{Y}|$.
3. If $\phi(t)$ is the squared hinge function $\frac{1}{2}[1 + t]_+^2$, then $H_L(\mathbf{p}) = \frac{\frac{1}{2}m^2}{\sum_{j=1}^{|\mathcal{Y}|} 1/(1-p_j)}$.

See Appendix B.6 for a proof. The first one recovers [23, Example 5] with a simpler proof. The last two are new. These entropies are illustrated in Figure 3. For the last two choices, H_L is strictly concave over $\Delta^{|\mathcal{Y}|}$ and the associated gradient mapping $\nabla(-H_L)^*$ typically sparse, although it cannot be computed in closed form.

Our proof also readily supports loss functions of the form $L(\mathbf{s}; \mathbf{e}_k) = \phi(s_k)$ with $\text{dom}(L) = \{\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|} : \mathbf{s}^\top \mathbf{1} = 0\}$ [53, 59]. It can be checked that we then simply need to replace all occurrences of $1 - p_j$ with p_j in Proposition 8.

B Proofs and derivations

B.1 Proof of Proposition 2

The two facts stated in (i) (H is always non-negative and maximized by the uniform distribution) follow directly from Jensen's inequality. Indeed, for any $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$, we have:

- $H(\mathbf{p}) \geq \sum_{j=1}^{|\mathcal{Y}|} p_j H(\mathbf{e}_j) = 0$;
- $H(\mathbf{1}/|\mathcal{Y}|) = H\left(\sum_{\mathbf{P} \in \mathcal{P}} \frac{1}{|\mathcal{Y}|!} \mathbf{P}\mathbf{p}\right) \geq \sum_{\mathbf{P} \in \mathcal{P}} \frac{1}{|\mathcal{Y}|!} H(\mathbf{P}\mathbf{p}) = H(\mathbf{p})$,

where \mathcal{P} is the set of permutation matrices. In the second case, the maximizer is unique due to strict concavity.

Let $\Omega(\mathbf{p}) := -H(\mathbf{p})$. To see (ii), note that the convexity and differentiability of Ω^* follow from basic Fenchel duality and the fact that Ω is *strictly* convex. It remains to prove that Ω^* is symmetric. In fact, we have $\Omega^*(\mathbf{P}\mathbf{s}) = \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} (\mathbf{P}\mathbf{s})^\top \mathbf{p} - \Omega(\mathbf{p}) = \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \mathbf{s}^\top \mathbf{P}^\top \mathbf{p} - \Omega(\mathbf{P}^\top \mathbf{p}) = \Omega^*(\mathbf{s})$. The last equality was obtained by a change of variable $\mathbf{p}' = \mathbf{P}^\top \mathbf{p}$, from which \mathbf{p} is recovered as $\mathbf{p} = \mathbf{P}\mathbf{p}'$, which proves (iii). Let us turn to (iv). Since Ω is convex, so is Ω^* , hence the gradient operator $\nabla \Omega^*$ is monotone, i.e., we have $(\mathbf{s}' - \mathbf{s})^\top (\mathbf{p}' - \mathbf{p}) \geq 0$ for any $\mathbf{s}, \mathbf{s}' \in \mathbb{R}^{|\mathcal{Y}|}$, $\mathbf{p} = \nabla \Omega^*(\mathbf{s})$ and $\mathbf{p}' = \nabla \Omega^*(\mathbf{s}')$. Let \mathbf{s}' be obtained from \mathbf{s} by swapping two coordinates, i.e., $s'_j = s_i$, $s'_i = s_j$, and $s'_k = s_k$ for any $k \notin \{i, j\}$. Then, since Ω is symmetric, we obtain:

$$2(s_j - s_i)(p_j - p_i) \geq 0,$$

which implies $s_i > s_j \Rightarrow p_i \geq p_j$ and $p_i > p_j \Rightarrow s_i \geq s_j$. To fully prove (iv), we need to show that the last inequality is strict: to do this, we simply invoke (iii) with a permutation matrix that permutes i and j , from which we must have $s_i = s_j \Rightarrow p_i = p_j$.

B.2 Proof of Proposition 3

Let $\Omega = -H$, and assume that Ω has non-empty subdifferential everywhere in $\Delta^{|\mathcal{Y}|}$. We need to show that for any $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$, there is some $\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|}$ such that $\mathbf{p} \in \operatorname{argmin}_{\mathbf{p}' \in \Delta^{|\mathcal{Y}|}} \Omega(\mathbf{p}') - \langle \mathbf{s}, \mathbf{p}' \rangle$. The Lagrangian associated with this minimization problem is:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\mu}, \lambda) = \Omega(\mathbf{p}) - \langle \mathbf{s} + \boldsymbol{\mu}, \mathbf{p} \rangle + \lambda(\mathbf{1}^\top \mathbf{p} - 1).$$

The KKT conditions are:

$$\begin{cases} 0 \in \partial_{\mathbf{p}} \mathcal{L}(\mathbf{p}, \boldsymbol{\mu}, \lambda) = \partial \Omega(\mathbf{p}) - \mathbf{s} - \boldsymbol{\mu} + \lambda \mathbf{1} \\ \langle \mathbf{p}, \boldsymbol{\mu} \rangle = 0 \\ \mathbf{p} \in \Delta^{|\mathcal{Y}|}, \boldsymbol{\mu} \geq 0. \end{cases}$$

For a given $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$, we seek \mathbf{s} such that $(\mathbf{p}, \boldsymbol{\mu}, \lambda)$ are a solution to the KKT conditions for some $\boldsymbol{\mu} \geq 0$ and $\lambda \in \mathbb{R}$.

We will show that such \mathbf{s} exists by simply choosing $\boldsymbol{\mu} = \mathbf{0}$ and $\lambda = 0$. Those choices are dual feasible and guarantee that the slackness complementary condition is satisfied. In this case, we have from the first condition that $\mathbf{s} \in \partial \Omega(\mathbf{p})$. From the assumption that Ω has non-empty subdifferential in all the simplex, we have that for any $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$ we can find a $\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|}$ such that (\mathbf{p}, \mathbf{s}) are a dual pair, i.e., $\mathbf{p} = \nabla \Omega^*(\mathbf{s})$, which proves that $\nabla \Omega^*(\mathbb{R}^{|\mathcal{Y}|}) = \Delta^{|\mathcal{Y}|}$.

To prove item 2, let us now fix \mathbf{p} and pick $\mathbf{s} \in \partial \Omega(\mathbf{p})$. Naturally, $\mathbf{s}' = \mathbf{s} + \lambda \mathbf{1}$ also maps to \mathbf{p} . Moreover, let $\mathcal{I} = \{1 \leq i \leq |\mathcal{Y}| : p_i = 0\}$. Then, any $\mu_i \geq 0$ for $i \in \mathcal{I}$ also leads to a solution, which shows any \mathbf{s}' with $s'_i = s_i$ for $i \notin \mathcal{I}$ and $s'_i \leq s_i$ for $i \in \mathcal{I}$ also maps to \mathbf{p} .

B.3 Proof of Proposition 4

We start by proving the following lemma.

Lemma 1 *Let H satisfy assumptions A.1–A.3. Then:*

1. *We have $\mathbf{s} \in \partial(-H)(\mathbf{e}_k)$ iff $s_k = (-H)^*(\mathbf{s})$. That is:*

$$\partial(-H)(\mathbf{e}_k) = \{\mathbf{s} \in \mathbb{R}^{|\mathcal{Y}|} : s_k \geq \langle \mathbf{s}, \mathbf{p} \rangle + H(\mathbf{p}), \forall \mathbf{p} \in \Delta^{|\mathcal{Y}|}\}.$$

2. *If $\mathbf{s} \in \partial(-H)(\mathbf{e}_k)$, then, we also have $\mathbf{s}' \in \partial(-H)(\mathbf{e}_k)$ for any \mathbf{s}' such that $s'_k = s_k$ and $s'_i \leq s_i$, for all $i \neq k$.*

Proof: (of the lemma.) Let $\Omega = -H$. From Proposition 2, we can consider $\partial \Omega(\mathbf{e}_1)$ without loss of generality, in which case any $\mathbf{s} \in \partial \Omega(\mathbf{e}_1)$ satisfies $s_1 = \max_j s_j$. We have $\mathbf{s} \in \partial \Omega(\mathbf{e}_1)$ iff $\Omega(\mathbf{e}_1) = \langle \mathbf{s}, \mathbf{e}_1 \rangle - \Omega^*(\mathbf{s}) = s_1 - \Omega^*(\mathbf{s})$. Since $\Omega(\mathbf{e}_1) = 0$, we must have $s_1 = \Omega^*(\mathbf{s}) \geq \sup_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \langle \mathbf{s}, \mathbf{p} \rangle - \Omega(\mathbf{p})$, which proves part 1. To see 2, note that we have $s'_k = s_k \geq \langle \mathbf{s}, \mathbf{p} \rangle - \Omega(\mathbf{p}) \geq \langle \mathbf{s}', \mathbf{p} \rangle - \Omega(\mathbf{p})$, for all $\mathbf{p} \in \Delta^{|\mathcal{Y}|}$, from which the result follows. ■

We now proceed to the proof of Proposition 4. Let $\Omega = -H$, and suppose that L_Ω has the separation margin property. Then, $\mathbf{s} = m\mathbf{e}_1$ satisfies the margin condition $s_1 \geq m + \max_{j \neq 1} s_j$, hence $L_\Omega(m\mathbf{e}_1, \mathbf{e}_1) = 0$. From the first part of Proposition 1, this implies $m\mathbf{e}_1 \in \partial \Omega(\mathbf{e}_1)$.

Conversely, let us assume that $m\mathbf{e}_1 \in \partial \Omega(\mathbf{e}_1)$. From the second part of Proposition 1, this implies that $\mathbf{s} \in \partial \Omega(\mathbf{e}_1)$ for any \mathbf{s} such that $s_1 = m$ and $s_i \leq 0$ for all $i \geq 2$; and more generally we have $\mathbf{s} + c\mathbf{1} \in \partial \Omega(\mathbf{e}_1)$. That is, any \mathbf{s} with $s_1 \geq m + \max_{i \neq 1} s_i$ satisfies $\mathbf{s} \in \partial \Omega(\mathbf{e}_1)$. From Proposition 1, this is equivalent to $L_\Omega(\mathbf{s}; \mathbf{e}_1) = 0$.

Let us now determine the margin of L_Ω , i.e., the smallest m such that $m\mathbf{e}_1 \in \partial \Omega(\mathbf{e}_1)$. From Proposition 1, this is equivalent to $m \geq m\mathbf{p}_1 - \Omega(\mathbf{p})$ for any $\mathbf{p} \in \Delta^d$, i.e., $-\Omega(\mathbf{p})(1 - p_1) \leq m$. Note that by Proposition 2 the “most competitive” \mathbf{p} ’s are sorted as \mathbf{e}_1 , so we may write $p_1 = \|\mathbf{p}\|_\infty$ without loss of generality. The margin of L_Ω is the smallest possible such margin, given by (7).

B.4 Proof of Proposition 5

Define $\Omega = -H$. Let us start by writing the margin expression (7) as a unidimensional optimization problem. This is done by noticing that the max-generalized entropy problem constrained to $\max(\mathbf{p}) = 1 - t$ gives

$\mathbf{p} = \left[1 - t, \frac{t}{d-1}, \dots, \frac{t}{d-1}\right]$, for $t \in [0, 1 - \frac{1}{d}]$ by a similar argument as the one used in Proposition 2, item 1. We obtain:

$$\text{margin}(L_\Omega) = \sup_{t \in [0, 1 - \frac{1}{d}]} \frac{-\Omega\left(\left[1 - t, \frac{t}{d-1}, \dots, \frac{t}{d-1}\right]\right)}{t}.$$

We write the argument above as $A(t) = \frac{-\Omega(\mathbf{e}_1 + t\mathbf{v})}{t}$, where $\mathbf{v} := [-1, \frac{1}{d-1}, \dots, \frac{1}{d-1}]$. We will first prove that A is decreasing in $[0, 1 - \frac{1}{d}]$, which implies that the supremum (and the margin) equals $A(0)$. Note that we have the following expression for the derivative of any function $f(\mathbf{e}_1 + t\mathbf{v})$:

$$(f(\mathbf{e}_1 + t\mathbf{v}))' = \mathbf{v}^\top \nabla f(\mathbf{e}_1 + t\mathbf{v}).$$

Using this fact, we can write the derivative $A'(t)$ as:

$$A'(t) = \frac{-t\mathbf{v}^\top \nabla \Omega(\mathbf{e}_1 + t\mathbf{v}) + \Omega(\mathbf{e}_1 + t\mathbf{v})}{t^2} := \frac{B(t)}{t^2}.$$

In turn, the derivative $B'(t)$ is:

$$\begin{aligned} B'(t) &= -\mathbf{v}^\top \nabla \Omega(\mathbf{e}_1 + t\mathbf{v}) - t(\mathbf{v}^\top \nabla \Omega(\mathbf{e}_1 + t\mathbf{v}))' + \mathbf{v}^\top \nabla \Omega(\mathbf{e}_1 + t\mathbf{v}) \\ &= -t(\mathbf{v}^\top \nabla \Omega(\mathbf{e}_1 + t\mathbf{v}))' \\ &= -t\mathbf{v}^\top \nabla \nabla \Omega(\mathbf{e}_1 + t\mathbf{v})\mathbf{v} \\ &\leq 0, \end{aligned}$$

where we denote by $\nabla \nabla \Omega$ the Hessian of Ω , and used the fact that it is positive semi-definite, due to the convexity of Ω . This implies that B is decreasing, hence for any $t \in [0, 1]$, $B(t) \leq B(0) = \Omega(\mathbf{e}_1) = 0$, where we used the fact $\|\nabla \Omega(\mathbf{e}_1)\| < \infty$, assumed as a condition of Proposition 3. Therefore, we must also have $A'(t) = \frac{B(t)}{t^2} \leq 0$ for any $t \in [0, 1]$, hence A is decreasing, and $\sup_{t \in [0, 1 - 1/d]} A(t) = \lim_{t \rightarrow 0+} A(t)$. By L'Hôpital's rule:

$$\begin{aligned} \lim_{t \rightarrow 0+} A(t) &= \lim_{t \rightarrow 0+} (-\Omega(\mathbf{e}_1 + t\mathbf{v}))' \\ &= -\mathbf{v}^\top \nabla \Omega(\mathbf{e}_1) \\ &= \nabla_1 \Omega(\mathbf{e}_1) - \frac{1}{d-1} \sum_{j \geq 2} \nabla_j \Omega(\mathbf{e}_1) \\ &= \nabla_1 \Omega(\mathbf{e}_1) - \nabla_2 \Omega(\mathbf{e}_1), \end{aligned}$$

which proves the first part.

If Ω is separable, then $\nabla \Omega(\mathbf{p}) = -\sum_{i=1}^d h(p_i)$, from which we obtain $\nabla_1 \Omega(\mathbf{e}_1) = -h'(1)$ and $\nabla_2 \Omega(\mathbf{e}_1) = -h'(0)$, yielding $\text{margin}(L_\Omega) = h'(0) - h'(1)$. Since h is twice differentiable, this equals $-\int_0^1 h''(t)dt$, completing the proof.

B.5 Proof of Proposition 7

Let H be a separable entropy on the simplex, i.e., $H(\mathbf{p}) = \sum_{j=1}^{|\mathcal{Y}|} h(p_j)$, where $h : [0, 1] \rightarrow \mathbb{R}_+$ is a non-negative, strictly concave, differentiable function such that $h(0) = h(1) = 0$, and $\text{dom}(H) = \Delta^{|\mathcal{Y}|}$; implying, importantly, that h' is invertible. As discussed in §3, such a H satisfies assumptions A.1–A.3. We show how computing $\nabla(-H)^*$ reduces to finding the root of a monotonic scalar function, for which efficient algorithms exist.

Recall the optimization problem we are trying to solve

$$\nabla(-H)^* := \operatorname{argmax}_{\mathbf{p} \in \Delta^{|\mathcal{Y}|}} \langle \mathbf{p}, \mathbf{s} \rangle + H(\mathbf{p}).$$

The Lagrangian is

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\mu}, \tau) := -H(\mathbf{p}) - \langle \mathbf{s} + \boldsymbol{\mu}, \mathbf{p} \rangle + \tau(\mathbf{1}^\top \mathbf{p} - 1).$$

The unique solution $(\mathbf{p}^*, \boldsymbol{\mu}^*, \tau^*)$ must satisfy the KKT conditions, which, in the separable case, are

$$\begin{aligned} h'(p_i) + s_i - \tau + \mu_i &= 0 \quad \forall 1 \leq i \leq |\mathcal{Y}| \\ \langle \mathbf{p}, \boldsymbol{\mu} \rangle &= 0 \\ \mathbf{p} &\in \Delta^{|\mathcal{Y}|}, \quad \boldsymbol{\mu} \geq 0. \end{aligned}$$

Let us define

$$\tau_{\min} := \max(\mathbf{s}) + h'(1) \quad \text{and} \quad \tau_{\max} := \max(\mathbf{s}) + h' \left(\frac{1}{|\mathcal{Y}|} \right).$$

Since h is strictly concave, h' is decreasing and so $\tau_{\min} < \tau_{\max}$. For any $\tau \in [\tau_{\min}, \tau_{\max}]$, we construct $\boldsymbol{\mu}$ as

$$\mu_i := \begin{cases} 0, & s_i - \tau \geq -h'(0) \\ -h'(0) - s_i + \tau, & s_i - \tau < -h'(0) \end{cases}$$

By construction, $\mu_i \geq 0$, satisfying dual feasibility. Injecting μ_i in $-h'(p_i) = s_i + \mu_i - \tau$ and combining the two cases, we obtain

$$-h'(p_i) = \max\{s_i - \tau, -h'(0)\}. \quad (17)$$

We show that i) the stationarity conditions have a unique solution given τ , and ii) $[\tau_{\min}, \tau_{\max}]$ forms a bracketing interval containing the primal-feasible solution, which can then be found by bisection. Since any such point verifies the KKT conditions, it must be the global solution.

Solving the stationarity conditions. Since $-h$ is strictly convex, its derivative $-h'$ is continuous and strictly increasing, and is thus a one-to-one mapping between $[0, 1]$ and $[-h'(0), -h'(1)]$. Denote by $(-h')^{-1} : [-h'(0), -h'(1)] \rightarrow [0, 1]$ its inverse. If $s_i - \tau \geq -h'(0)$, we have that

$$\begin{aligned} -h'(0) &\leq -h'(p_i) = s_i - \tau \leq \max(\mathbf{s}) - \tau_{\min} \\ &= \max(\mathbf{s}) - \max(\mathbf{s}) - h'(1) \\ &= -h'(1). \end{aligned}$$

Otherwise, $-h'(p_i) = -h'(0)$. This verifies that the r.h.s. of Eq. (17) is always within the domain of $(-h')^{-1}$. We can thus apply the inverse to both sides to solve for p_i , obtaining

$$p_i(\tau) = (-h')^{-1}(\max\{s_i - \tau, -h'(0)\}). \quad (18)$$

Validating the bracketing interval. Consider the primal feasibility function $f(\tau) = \mathbf{1}^\top \mathbf{p}(\tau) - 1$; $\mathbf{p}(\tau)$ is primal feasible iff τ is a root of f . We show that f is decreasing on $[\tau_{\min}, \tau_{\max}]$, and that it has opposite signs at the two extremities. From the intermediate value theorem, the unique root τ^* must satisfy $\tau^* \in [\tau_{\min}, \tau_{\max}]$.

As the inverse of an increasing function, $(-h')^{-1}$ is increasing. Therefore, for all i , $p_i(\tau)$ is decreasing, and so is the sum $f(\tau) = \sum_i p_i(\tau) - 1$. It remains to check the signs at the boundaries.

$$\begin{aligned} \sum_i p_i(\tau_{\max}) &= \sum_i (-h')^{-1}(\max\{s_i - \max(\mathbf{s}) - h'(1/|\mathcal{Y}|), -h'(0)\}) \\ &\leq |\mathcal{Y}|(-h')^{-1}(\max\{-h'(1/|\mathcal{Y}|), -h'(0)\}) \\ &= |\mathcal{Y}|(-h')^{-1}(-h'(1/|\mathcal{Y}|)) = 1 \end{aligned}$$

upper-bounding each term of the sum by the largest one. At the other end,

$$\begin{aligned} \sum_i p_i(\tau_{\min}) &= \sum_i (-h')^{-1}(\max\{s_i - \max(\mathbf{s}) - h'(1), -h'(0)\}) \\ &\geq (-h')^{-1}(\max\{-h'(1), -h'(0)\}) \\ &= (-h')^{-1}(-h'(1)) = 1 \end{aligned}$$

using that a sum of non-negative terms is no less than its largest term. Therefore, $f(\tau_{\min}) \geq 0$ and $f(\tau_{\max}) \leq 0$. This implies that there must exist τ^* in $[\tau_{\min}, \tau_{\max}]$ satisfying $f(\tau^*) = 0$. A corresponding triplet $(\mathbf{p}(\tau^*), \boldsymbol{\mu}(\tau^*), \tau^*)$ thus satisfies all of the KKT conditions, confirming that it is the globally optimal solution.

Algorithm 1 is an example of a bisection algorithm for finding the optimal solution; more advanced root finding methods can also be used. We note that the resulting algorithm resembles the method provided in [33], with a non-trivial difference being the order of the thresholding and $(-h)^{-1}$ in Eq. (18). Our result thus applies for important cases like the Tsallis entropy with $\alpha = 1.5$, where [33] is not applicable.

Algorithm 1 Compute $\nabla(-H)^*(s)$

Input: $s \in \mathbb{R}^{|\mathcal{Y}|}$, $H(p) = \sum_i h(p_i)$
 $\hat{p}(\tau) := (-h')^{-1}(\max\{s - \tau, -h'(0)\})$
 $f(\tau) := \hat{p}(\tau)^\top \mathbf{1} - 1$
 $\tau_{\min} \leftarrow \max(s) + h'(1)$;
 $\tau_{\max} \leftarrow \max(s) + h'(1/|\mathcal{Y}|)$
 $\tau \leftarrow (\tau_{\min} + \tau_{\max})/2$
while $|f(\tau)| > \epsilon$
 $\tau \leftarrow (\tau_{\min} + \tau_{\max})/2$
 if $f(\tau) < 0$ $\tau_{\max} \leftarrow \tau$
 else $\tau_{\min} \leftarrow \tau$
Output: $\nabla(-H)^*(s) \approx \hat{p}(\tau)$

B.6 Proof of Proposition 8

We first need the following lemma.

Lemma 2 Let $L(s; e_k)$ be defined as

$$L(s; e_k) := \begin{cases} \sum_j c_{k,j} \phi(s_j) & \text{if } s^\top \mathbf{1} = 0 \\ \infty & \text{o.w.} \end{cases},$$

where $\phi: \mathbb{R} \rightarrow \mathbb{R}_+$ is convex. Then, H_L defined in (12) equals

$$-H_L(p) = \min_{\tau \in \mathbb{R}} \sum_j (p^\top c_j) \phi^* \left(\frac{-\tau}{p^\top c_j} \right).$$

Proof: We want to solve

$$H_L(p) = \min_{s^\top \mathbf{1}=0} \sum_j p_j L(s; e_j) = \min_{s^\top \mathbf{1}=0} \sum_j p_j \sum_i c_{j,i} \phi(s_i) = \min_{s^\top \mathbf{1}=0} \sum_i (p^\top c_i) \phi(s_i),$$

where c_i is a vector gathering $c_{j,i}$ for all j . Introducing a Lagrange multiplier we get

$$H_L(p) = \min_{s \in \mathbb{R}^{|\mathcal{Y}|}} \max_{\tau \in \mathbb{R}} \sum_i (p^\top c_i) \phi(s_i) + \tau \sum_i s_i.$$

Strong duality holds and we can swap the order of the min and max. Routine calculations then give

$$-H_L(p) = \min_{\tau \in \mathbb{R}} \sum_i (p^\top c_i) \phi^* \left(\frac{-\tau}{p^\top c_i} \right). \quad (19)$$

■

We now prove Proposition 8. First, we rewrite $L(s; e_k) = \sum_{j \neq k} \phi(s_j - s_k)$ as $L(s; e_k) = \sum_j c_{k,j} \phi(s_j - s_k)$, where we choose $c_{k,j} = 0$ if $k = j$, 1 otherwise, leading to $p^\top c_j = 1 - p_j$. Because $\phi(s_j - s_k)$ is shift invariant w.r.t. s , without loss of generality, we can further rewrite the loss as $L(s; e_k) = \sum_j c_{k,j} \phi(s_j)$ with $\text{dom}(L) = \{s \in \mathbb{R}^{|\mathcal{Y}|} : s^\top \mathbf{1} = 0\}$. Hence, Lemma 2 applies. We now derive closed form for specific choices of ϕ (we let $d := |\mathcal{Y}|$ for convenience).

Hinge loss. When $\phi(t) = [1 + t]_+$, the conjugate is

$$\phi^*(u) = \begin{cases} -u & \text{if } u \in [0, 1] \\ \infty & \text{o.w.} \end{cases}.$$

The constraint set for τ is therefore $\mathcal{C} := \bigcap_{j \in [d]} [-\mathbf{p}^\top \mathbf{c}_j, 0] = \left[-\min_{j \in [d]} \mathbf{p}^\top \mathbf{c}_j, 0 \right]$. Hence

$$-\mathbf{H}_L(\mathbf{p}) = \min_{\tau \in \mathcal{C}} d\tau = -d \min_{j \in [d]} \mathbf{p}^\top \mathbf{c}_j.$$

This recovers [23, Example 5, §A.6] with a simpler proof. We next turn to the following new results.

Smoothed hinge loss. We add quadratic regularization to the conjugate [50]:

$$\phi^*(u) = \begin{cases} -u + \frac{1}{2}u^2 & \text{if } u \in [0, 1] \\ \infty & \text{o.w.} \end{cases}.$$

Going back to ϕ , we obtain:

$$\phi(t) = \begin{cases} 0 & \text{if } t \leq -1 \\ 1 + t - \frac{\gamma}{2} & \text{if } t \geq \gamma - 1 \\ \frac{1}{2\gamma}(1+t)^2 & \text{o.w.} \end{cases} \quad (20)$$

The constraint set for τ is the same as before, $\mathcal{C} = \left[-\min_{j \in [d]} \mathbf{p}^\top \mathbf{c}_j, 0 \right]$.

Plugging ϕ^* into (19), we obtain

$$-\mathbf{H}_L(\mathbf{p}) = \min_{\tau \in \mathcal{C}} \frac{\tau^2}{2} \sum_{j=1}^d \frac{1}{\mathbf{p}^\top \mathbf{c}_j} + d\tau. \quad (21)$$

Since the problem is unidimensional, let us solve for τ unconstrained first:

$$\tau = -d / \left(\sum_{j=1}^d 1/(\mathbf{p}^\top \mathbf{c}_j) \right). \quad (22)$$

We notice that $\tau \leq -\min_j \mathbf{p}^\top \mathbf{c}_j$ for $\mathbf{c}_j \geq \mathbf{0}$ since $\sum_j \frac{\min_i \mathbf{p}^\top \mathbf{c}_i}{\mathbf{p}^\top \mathbf{c}_j} \in [0, d]$. This expression of τ is not feasible. Hence the optimal solution is at the boundary and $\tau^* = -\min_j \mathbf{p}^\top \mathbf{c}_j$. Plugging that expression back into (21) gives the claimed expression of \mathbf{H}_L .

Squared hinge loss. When $\phi(t) = \frac{1}{2}[1+t]_+^2$, the conjugate is

$$\phi^*(u) = \begin{cases} -u + \frac{1}{2}u^2 & \text{if } u \geq 0 \\ \infty & \text{o.w.} \end{cases}.$$

The constraint is now $\tau \leq 0$. Hence, the optimal solution of the problem w.r.t. τ in (19) is now (22) for all $\mathbf{p} \in \Delta^d, \mathbf{c}_j \geq \mathbf{0}$. Simplifying, we get

$$\mathbf{H}_L(\mathbf{p}) = \frac{\frac{1}{2}d^2}{\sum_{j=1}^d 1/(\mathbf{p}^\top \mathbf{c}_j)}.$$

B.7 Derivation of the dual coordinate ascent algorithm

We rewrite the primal objective (9) as

$$\min_{W, \{\boldsymbol{\theta}_i\}} \sum_{i=1}^n L_\Omega(\boldsymbol{\theta}_i; \mathbf{y}_i) + G(W) \text{ s.t. } \boldsymbol{\theta}_i = W\mathbf{x}_i \forall i \in [n].$$

After routine calculations, we find that its dual is

$$-\min_{\boldsymbol{\alpha}} \sum_{i=1}^n L_\Omega^*(-\boldsymbol{\alpha}_i; \mathbf{y}_i) + G^*\left(\sum_{i=1}^n \boldsymbol{\alpha}_i \mathbf{x}_i^\top\right).$$

Using $L_\Omega^*(-\alpha; \mathbf{y}) = \Omega(\mathbf{y} - \alpha) - \Omega(\mathbf{y})$ and using the change of variable $\beta_i := \mathbf{y}_i - \alpha_i$, we get

$$-\min_{\beta} \sum_{i=1}^n \Omega(\beta_i) - \Omega(\mathbf{y}_i) + G^*(V(\beta)) \text{ s.t. } \beta_i \in \text{dom}(\Omega) \forall i \in [n],$$

where $V(\beta) := \sum_{i=1}^n (\mathbf{y}_i - \beta_i) \mathbf{x}_i^\top$. The sub-problem associated with $i \in [n]$ is

$$\underset{\beta_i \in \text{dom}(\Omega)}{\text{argmin}} \quad \Omega(\beta_i) + G^*(V(\beta)). \quad (23)$$

Since this problem could be hard to solve, we follow [50, Option I] and consider instead an upper-bound. Since G^* is $\frac{1}{\lambda}$ -smooth w.r.t. the dual norm $\|\cdot\|$, it holds that

$$G^*(V(\beta)) \leq G^*(V(\bar{\beta})) + \langle \nabla G^*(V(\bar{\beta})), V(\beta) - V(\bar{\beta}) \rangle + \frac{1}{2\lambda} \|V(\beta) - V(\bar{\beta})\|^2,$$

where $\bar{\beta}$ denotes the current iterate of β . Using $V(\beta) - V(\bar{\beta}) = \sum_{i=1}^n (\bar{\beta}_i - \beta_i) \mathbf{x}_i^\top$, we get

$$G^*(V(\beta)) \leq G^*(V(\bar{\beta})) + \sum_{i=1}^n \langle \nabla G^*(V(\bar{\beta})) \mathbf{x}_i, \bar{\beta}_i - \beta_i \rangle + \frac{\sigma_i}{2} \|\bar{\beta}_i - \beta_i\|^2,$$

where $\sigma_i := \frac{\|\mathbf{x}_i\|^2}{\lambda}$. Substituting $G^*(V(\beta))$ by the above upper bound into (23) and ignoring constant terms, we get the following approximate sub-problem:

$$\underset{\beta_i \in \text{dom}(\Omega)}{\text{argmin}} \quad \Omega(\beta_i) - \beta_i^\top \mathbf{v}_i + \frac{\sigma_i}{2} \|\beta_i\|^2 = \text{prox}_{\frac{1}{\sigma_i} \Omega} \left(\frac{\mathbf{v}_i}{\sigma_i} \right), \quad (24)$$

where $\mathbf{v}_i := \nabla G^*(V(\bar{\beta})) \mathbf{x}_i + \sigma_i \bar{\beta}_i$ and where we defined the proximity operator

$$\text{prox}_{\tau \Omega}(\mathbf{x}) := \underset{\mathbf{y} \in \text{dom}(\Omega)}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 + \tau \Omega(\mathbf{y}).$$

Note that when G^* is a quadratic function, (24) is an optimal solution of (23).

Summary. A block-wise update is $\beta_i \leftarrow \text{prox}_{\frac{1}{\sigma_i} \Omega}(\mathbf{v}_i / \sigma_i)$, where $\mathbf{v}_i := \nabla G^*(V(\beta)) \mathbf{x}_i + \sigma_i \beta_i$ and $\sigma_i := \frac{\|\mathbf{x}_i\|^2}{\lambda}$. A block coordinate ascent algorithm picks $i \in [n]$, performs this update and repeats, until converging to an optimal dual solution β^* . Since G is a λ -strongly convex regularizer, given β^* , we retrieve the optimal primal solution by $W^* = \nabla G^*(V(\beta^*))$.

Examples of regularization. When $G(W) = \frac{\lambda}{2} \|W\|_F^2$, $G^*(V) = \frac{1}{2\lambda} \|V\|_F^2$ and $\nabla G^*(V) = \frac{1}{\lambda} V$.

When $G(W) = \frac{\lambda}{2} \|W\|_F^2 + \lambda \rho R(W)$, for some regularizer $R(W)$, we have

$$\begin{aligned} \nabla G^*(V) &= \underset{W}{\text{argmax}} \quad \langle W, V \rangle - \frac{\lambda}{2} \|W\|_F^2 - \lambda \rho R(W) \\ &= \underset{W}{\text{argmin}} \quad \frac{1}{2} \left\| W - \frac{V}{\lambda} \right\|_F^2 + \rho R(W) \\ &= \text{prox}_{\rho R}(V/\lambda). \end{aligned}$$

The conjugate is equal to $G^*(V) = \langle \nabla G^*(V), V \rangle - G(\nabla G^*(V))$.

For instance, if we choose $R(W) = \|W\|_1 := \sum_{i,j} |w_{i,j}|$, then $G(W)$ is the **elastic-net** regularization and $\text{prox}_{\rho R}$ is the well-known **soft-thresholding** operator

$$\text{prox}_{\rho R}(V) = \text{sign}(V)[|V| - \rho]_+,$$

where all operations above are performed element-wise.

Examples of losses. For the squared loss, $\text{dom}(\Omega) = \mathbb{R}^d$ and $\Omega(\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|^2$. Hence:

$$\text{prox}_{\tau\Omega}(\mathbf{x}) = \underset{\mathbf{y} \in \mathbb{R}^d}{\text{argmin}} \frac{1}{2}\|\mathbf{y} - \mathbf{x}\|^2 + \frac{\tau}{2}\|\mathbf{y}\|^2 = \frac{\mathbf{x}}{\tau + 1}.$$

For the perceptron loss [47, 15], $\text{dom}(\Omega) = \Delta^{|\mathcal{Y}|}$ and $\Omega = 0$. Hence:

$$\text{prox}_{\tau\Omega}(\mathbf{x}) = \underset{\mathbf{p} \in \Delta^{|\mathcal{Y}|}}{\text{argmin}} \frac{1}{2}\|\mathbf{p} - \mathbf{x}\|^2.$$

For the sparsemax loss [37], $\text{dom}(\Omega) = \Delta^{|\mathcal{Y}|}$ and $\Omega(\mathbf{p}) = \frac{1}{2}(\|\mathbf{p}\|^2 - 1)$. Hence:

$$\text{prox}_{\tau\Omega}(\mathbf{x}) = \underset{\mathbf{p} \in \Delta^{|\mathcal{Y}|}}{\text{argmin}} \frac{1}{2}\|\mathbf{p} - \mathbf{x}\|^2 + \frac{\tau}{2}\|\mathbf{p}\|^2 = \underset{\mathbf{p} \in \Delta^{|\mathcal{Y}|}}{\text{argmin}} \frac{1}{2}\left\|\mathbf{p} - \frac{\mathbf{x}}{\tau + 1}\right\|^2.$$

For the cost-sensitive multiclass hinge loss, $\text{dom}(\Omega) = \Delta^{|\mathcal{Y}|}$ and since $\Omega^*(\boldsymbol{\theta} + \mathbf{c}_{\mathbf{y}}) = \max_{\mathbf{y}' \in \mathcal{Y}} \theta(\mathbf{y}') + c_{\mathbf{y}}(\mathbf{y}')$, we have $\Omega(\mathbf{p}) = -\langle \mathbf{p}, \mathbf{c}_{\mathbf{y}} \rangle$. Hence:

$$\text{prox}_{\tau\Omega}(\mathbf{x}) = \underset{\mathbf{p} \in \Delta^{|\mathcal{Y}|}}{\text{argmin}} \frac{1}{2}\|\mathbf{p} - \mathbf{x}\|^2 - \tau\langle \mathbf{p}, \mathbf{c}_{\mathbf{y}} \rangle = \underset{\mathbf{p} \in \Delta^{|\mathcal{Y}|}}{\text{argmin}} \frac{1}{2}\|\mathbf{p} - (\mathbf{x} + \tau\mathbf{c}_{\mathbf{y}})\|^2.$$

Choosing $\mathbf{c}_{\mathbf{y}} = \mathbf{1} - \mathbf{y}$, where $\mathbf{y} \in \{e_1, \dots, e_{|\mathcal{Y}|}\}$, recovers the usual multiclass hinge loss [17].

For the multiclass logistic loss, $\text{dom}(\Omega) = \Delta^{|\mathcal{Y}|}$ and Ω is the negative Shannon entropy. Unlike $\nabla\Omega^*$, $\text{prox}_{\tau\Omega}(\mathbf{x})$ cannot be computed in closed form but we can compute it approximately using a few iterations of proximal gradient descent.