
An Embedding Framework for Consistent Polyhedral Surrogates

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We formalize and study the natural approach of designing convex surrogate loss
2 functions via embeddings for problems such as classification or ranking. In this
3 approach, one embeds each of the finitely many predictions (e.g. classes) as a point
4 in \mathbb{R}^d , assigns the original loss values to these points, and convexifies the loss in
5 between to obtain a surrogate. We prove that this approach is equivalent, in a strong
6 sense, to working with polyhedral (piecewise linear convex) losses. Moreover,
7 given any polyhedral loss L , we give a construction of a link function through
8 which L is a consistent surrogate for the loss it embeds. We go on to illustrate
9 the power of this embedding framework with succinct proofs of consistency or
10 inconsistency of various polyhedral surrogates in the literature.

11 1 Introduction

12 Convex surrogate losses are a central building block in machine learning for classification and
13 classification-like problems. A growing body of work seeks to design and analyze convex surrogates
14 for given loss functions, and more broadly, understand when such surrogates can and cannot be found.
15 For example, recent work has developed tools to bound the required number of dimensions of the
16 surrogate’s hypothesis space [12, 24]. Yet in some cases these bounds are far from tight, such as
17 for *abstain loss* (classification with an abstain option) [4, 24, 25, 33, 34]. Furthermore, the kinds of
18 strategies available for constructing surrogates, and their relative power, are not well-understood.

19 We augment this literature by studying a particularly natural approach for finding convex surrogates,
20 wherein one “embeds” a discrete loss. Specifically, we say a convex surrogate L embeds a discrete
21 loss ℓ if there is an injective embedding from the discrete reports (predictions) to a vector space such
22 that (i) the original loss values are recovered, and (ii) a report is ℓ -optimal if and only if the embedded
23 report is L -optimal. If this embedding can be extended to a calibrated link function, which maps
24 approximately L -optimal reports to ℓ -optimal reports, consistency follows [2]. Common examples
25 which follow this general construction include hinge loss as a surrogate for 0-1 loss and the abstain
26 surrogate mentioned above.

27 Using tools from property elicitation, we show a tight relationship between such embeddings and
28 the class of polyhedral (piecewise-linear convex) loss functions. In particular, by focusing on Bayes
29 risks, we show that every discrete loss is embedded by some polyhedral loss, and every polyhedral
30 loss function embeds some discrete loss. Moreover, we show that any polyhedral loss gives rise to
31 a calibrated link function to the loss it embeds, thus giving a very general framework to construct
32 consistent convex surrogates for arbitrary losses.

33 **Related works.** The literature on convex surrogates focuses mainly on smooth surrogate losses [4,
34 5, 7, 8, 26, 30]. Nevertheless, nonsmooth losses, such as the polyhedral losses we consider, have
35 been proposed and studied for a variety of classification-like problems [18, 31, 32]. A notable

addition to this literature is Ramaswamy et al. [25], who argue that nonsmooth losses may enable dimension reduction of the prediction space (range of the surrogate hypothesis) relative to smooth losses, illustrating this conjecture with a surrogate for *abstain loss* needing only $\log n$ dimensions for n labels, whereas the best known smooth loss needs $n - 1$. Their surrogate is a natural example of an embedding (cf. Section 5.1), and serves as inspiration for our work.

While property elicitation has by now an extensive literature [9, 11, 14, 16, 17, 21, 28, 29], these works are mostly concerned with point estimation problems. Literature directly connecting property elicitation to consistency is sparse, with the main reference being Agarwal and Agarwal [2]; note however that they consider single-valued properties, whereas properties elicited by general convex losses are necessarily set-valued.

2 Setting

For discrete prediction problems like classification, due to hardness of directly optimizing a given discrete loss, many machine learning algorithms can be thought of as minimizing a surrogate loss function with better optimization qualities, e.g., convexity. Of course, to show that this surrogate loss successfully addresses the original problem, one needs to establish consistency, which depends crucially on the choice of link function that maps surrogate reports (predictions) to original reports. After introducing notation, and terminology from property elicitation, we thus give a sufficient condition for consistency (Def. 4) which depends solely on the conditional distribution over \mathcal{Y} .

2.1 Notation and Losses

Let \mathcal{Y} be a finite outcome (label) space, and throughout let $n = |\mathcal{Y}|$. The set of probability distributions on \mathcal{Y} is denoted $\Delta_{\mathcal{Y}} \subseteq \mathbb{R}^{\mathcal{Y}}$, represented as vectors of probabilities. We write p_y for the probability of outcome $y \in \mathcal{Y}$ drawn from $p \in \Delta_{\mathcal{Y}}$.

We assume that a given discrete prediction problem, such as classification, is given in the form of a *discrete loss* $\ell : \mathcal{R} \rightarrow \mathbb{R}_+^{\mathcal{Y}}$, which maps a report (prediction) r from a finite set \mathcal{R} to the vector of loss values $\ell(r) = (\ell(r)_y)_{y \in \mathcal{Y}}$ for each possible outcome $y \in \mathcal{Y}$. We will assume throughout that the given discrete loss is *non-redundant*, meaning every report is uniquely optimal (minimizes expected loss) for some distribution $p \in \Delta_{\mathcal{Y}}$. Similarly, surrogate losses will be written $L : \mathbb{R}^d \rightarrow \mathbb{R}_+^{\mathcal{Y}}$, typically with reports written $u \in \mathbb{R}^d$. We write the corresponding expected loss when $Y \sim p$ as $\langle p, \ell(r) \rangle$ and $\langle p, L(u) \rangle$. The *Bayes risk* of a loss $L : \mathbb{R}^d \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ is the function $\underline{L} : \Delta_{\mathcal{Y}} \rightarrow \mathbb{R}_+$ given by $\underline{L}(p) := \inf_{u \in \mathbb{R}^d} \langle p, L(u) \rangle$; naturally for discrete losses we write $\underline{\ell}$ (and the infimum is over \mathcal{R}).

For example, 0-1 loss is a discrete loss with $\mathcal{R} = \mathcal{Y} = \{-1, 1\}$ given by $\ell_{0-1}(r)_y = \mathbb{1}\{r \neq y\}$, with Bayes risk $\underline{\ell}_{0-1}(p) = 1 - \max_{y \in \mathcal{Y}} p_y$. Two important surrogates for ℓ_{0-1} are hinge loss $L_{\text{hinge}}(u)_y = (1 - yu)_+$, where $(x)_+ = \max(x, 0)$, and logistic loss $L(u)_y = \log(1 + \exp(-yu))$ for $u \in \mathbb{R}$.

Most of the surrogates L we consider will be *polyhedral*, meaning piecewise linear and convex; we therefore briefly recall the relevant definitions. In \mathbb{R}^d , a *polyhedral set* or *polyhedron* is the intersection of a finite number of closed halfspaces. A *polytope* is a bounded polyhedral set. A convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *polyhedral* if its epigraph is polyhedral, or equivalently, if it can be written as a pointwise maximum of a finite set of affine functions [27].

Definition 1 (Polyhedral loss). *A loss $L : \mathbb{R}^d \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ is polyhedral if $L(u)_y$ is a polyhedral (convex) function of u for each $y \in \mathcal{Y}$.*

For example, hinge loss is polyhedral, whereas logistic loss is not.

2.2 Property Elicitation

To make headway, we will appeal to concepts and results from the property elicitation literature, which elevates the *property*, or map from distributions to optimal reports, as a central object to study in its own right. In our case, this map will often be multivalued, meaning a single distribution could yield multiple optimal reports. (For example, when $p = (1/2, 1/2)$, both $y = 1$ and $y = -1$ optimize 0-1 loss.) To this end, we will use double arrow notation to mean a mapping to all nonempty subsets, so that $\gamma : \Delta_{\mathcal{Y}} \rightrightarrows \mathcal{R}$ is shorthand for $\Gamma : \Delta_{\mathcal{Y}} \rightarrow 2^{\mathcal{R}} \setminus \emptyset$. See the discussion following Definition 3 for conventions regarding $\mathcal{R}, \Gamma, \gamma, L, \ell$, etc.

85 **Definition 2** (Property, level set). A property is a function $\Gamma : \Delta_{\mathcal{Y}} \rightrightarrows \mathcal{R}$. The level set of Γ for report
86 r is the set $\Gamma_r := \{p : r \in \Gamma(p)\}$.

87 Intuitively, $\Gamma(p)$ is the set of reports which should be optimal for a given distribution p , and Γ_r is the
88 set of distributions for which the report r should be optimal. For example, the *mode* is the property
89 $\text{mode}(p) = \arg \max_{y \in \mathcal{Y}} p_y$, and captures the set of optimal reports for 0-1 loss: for each distribution
90 over the labels, one should report the most likely label. In this case we say 0-1 loss *elicits* the mode,
91 as we formalize below.

92 **Definition 3** (Elicits). A loss $L : \mathcal{R} \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$, elicits a property $\Gamma : \Delta_{\mathcal{Y}} \rightrightarrows \mathcal{R}$ if

$$\forall p \in \Delta_{\mathcal{Y}}, \quad \Gamma(p) = \arg \min_{r \in \mathcal{R}} \langle p, L(r) \rangle. \quad (1)$$

93 As Γ is uniquely defined by L , we write $\text{prop}[L]$ to refer to the property elicited by a loss L .

94 For finite properties (those with $|\mathcal{R}| < \infty$) and discrete losses, we will use lowercase notation γ and
95 ℓ , respectively, with reports $r \in \mathcal{R}$; for surrogate properties and losses we use Γ and L , with reports
96 $u \in \mathbb{R}^d$. For general properties and losses, we will also use Γ and L , as above.

97 2.3 Links and Embeddings

98 To assess whether a surrogate and link function align with the original loss, we turn to the common
99 condition of *calibration*. Roughly, a surrogate and link are calibrated if the best possible expected
100 loss achieved by linking to an incorrect report is strictly suboptimal.

101 **Definition 4.** Let original loss $\ell : \mathcal{R} \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$, proposed surrogate $L : \mathbb{R}^d \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$, and link function
102 $\psi : \mathbb{R}^d \rightarrow \mathcal{R}$ be given. We say (L, ψ) is calibrated with respect to ℓ if for all $p \in \Delta_{\mathcal{Y}}$,

$$\inf_{u \in \mathbb{R}^d : \psi(u) \notin \gamma(p)} \langle p, L(u) \rangle > \inf_{u \in \mathbb{R}^d} \langle p, L(u) \rangle. \quad (2)$$

103 It is well-known that calibration implies *consistency*, in the following sense (cf. [2]). Given feature
104 space \mathcal{X} , fix a fixed distribution $D \in \Delta(\mathcal{X} \times \mathcal{Y})$. Let L^* be the best possible expected L -loss achieved
105 by any hypothesis $H : \mathcal{X} \rightarrow \mathbb{R}^d$, and ℓ^* the best expected ℓ -loss for any hypothesis $h : \mathcal{X} \rightarrow \mathcal{R}$,
106 respectively. Then (L, ψ) is consistent if a sequence of surrogate hypotheses H_1, H_2, \dots whose
107 L -loss limits to L^* , then the ℓ -loss of $\psi \circ H_1, \psi \circ H_2, \dots$ limits to ℓ^* . As Definition 4 does not
108 involve the feature space \mathcal{X} , we will drop it for the remainder of the paper.

109 Several consistent convex surrogates in the literature can be thought of as “embeddings”, wherein one
110 maps the discrete reports to a vector space, and finds a convex loss which agrees with the original
111 loss. A key condition is that the original reports should be optimal exactly when the corresponding
112 embedded points are optimal. We formalize this notion as follows.

113 **Definition 5.** A loss $L : \mathbb{R}^d \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$ embeds a loss $\ell : \mathcal{R} \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$ if there exists some injective
114 embedding $\varphi : \mathcal{R} \rightarrow \mathbb{R}^d$ such that (i) for all $r \in \mathcal{R}$ we have $L(\varphi(r)) = \ell(r)$, and (ii) for all
115 $p \in \Delta_{\mathcal{Y}}, r \in \mathcal{R}$ we have

$$r \in \text{prop}[\ell](p) \iff \varphi(r) \in \text{prop}[L](p). \quad (3)$$

116 Note that it is not clear if embeddings give rise to calibrated links; indeed, apart from mapping the
117 embedded points back to their original reports via $\psi(\varphi(r)) = r$, how to map the remaining values is
118 far from clear. We address the question of when embeddings lead to calibrated links in Section 4.

119 To illustrate the idea of embedding, let us examine hinge loss in detail as a surrogate for 0-1 loss
120 for binary classification. Recall that we have $\mathcal{R} = \mathcal{Y} = \{-1, +1\}$, with $L_{\text{hinge}}(u)_y = (1 - uy)_+$
121 and $\ell_{0-1}(r)_y := \mathbb{1}\{r \neq y\}$, typically with link function $\psi(u) = \text{sgn}(u)$. We will see that hinge
122 loss embeds (2 times) 0-1 loss, via the embedding $\varphi(r) = r$. For condition (i), it is straightforward
123 to check that $L_{\text{hinge}}(r)_y = 2\ell_{0-1}(r)_y$ for all $r, y \in \{-1, 1\}$. For condition (ii), let us compute the
124 property each loss elicits, i.e., the set of optimal reports for each p :

$$\text{prop}[\ell_{0-1}](p) = \begin{cases} 1 & p_1 > 1/2 \\ \{-1, 1\} & p_1 = 1/2 \\ -1 & p_1 < 1/2 \end{cases} \quad \text{prop}[L_{\text{hinge}}](p) = \begin{cases} [1, \infty) & p_1 = 1 \\ 1 & p_1 \in (1/2, 1) \\ [-1, 1] & p_1 = 1/2 \\ -1 & p_1 \in (0, 1/2) \\ (-\infty, -1] & p_1 = 0 \end{cases}.$$

125 In particular, we see that $-1 \in \text{prop}[\ell_{0-1}](p) \iff p_1 \in [0, 1/2] \iff -1 \in \text{prop}[L_{\text{hinge}}](p)$,
 126 and $1 \in \text{prop}[\ell_{0-1}](p) \iff p_1 \in [1/2, 1] \iff 1 \in \text{prop}[L_{\text{hinge}}](p)$. With both conditions of
 127 Definition 5 satisfied, we conclude that L_{hinge} embeds $2\ell_{0-1}$. In this particular case, it is known
 128 (L_{hinge}, ψ) is calibrated for $\psi(u) = \text{sgn}(u)$; we address in Section 4 the interesting question of
 129 whether embeddings lead to calibration in general.

130 3 Embeddings and Polyhedral Losses

131 In this section, we establish a tight relationship between the technique of embedding and the use of
 132 polyhedral (piecewise-linear convex) surrogate losses. We defer to the following section the question
 133 of when such surrogates are consistent.

134 To begin, we observe that, somewhat surprisingly, our embedding condition in Definition 5 is
 135 equivalent to merely matching Bayes risks. This useful fact will drive many of our results.

136 **Proposition 1.** *A loss L embeds discrete loss ℓ if and only if $\underline{L} = \underline{\ell}$.*

137 *Proof.* Throughout we have $L : \mathbb{R}^d \rightarrow \mathbb{R}_+^{\mathcal{Y}}$, $\ell : \mathcal{R} \rightarrow \mathbb{R}_+^{\mathcal{Y}}$, and define $\Gamma = \text{prop}[L]$ and $\gamma = \text{prop}[\ell]$.
 138 Suppose L embeds ℓ via the embedding φ . Letting $\mathcal{U} := \varphi(\mathcal{R})$, define $\gamma' : \Delta_{\mathcal{Y}} \rightrightarrows \mathcal{U}$ by $\gamma' : p \mapsto$
 139 $\Gamma(p) \cap \mathcal{U}$. To see that $\gamma'(p) \neq \emptyset$ for all $p \in \Delta_{\mathcal{Y}}$, note that by the definition of γ as the property elicited
 140 by ℓ we have some $r \in \gamma(p)$, and by the embedding condition (3), $\varphi(r) \in \Gamma(p)$. By Lemma 3, we
 141 see that $L|_{\mathcal{U}}$ (the loss L with reports restricted to \mathcal{U}) elicits γ' and $\underline{L} = \underline{L|_{\mathcal{U}}}$. As $L(\varphi(\cdot)) = \ell(\cdot)$ by
 142 the embedding, we have

$$\underline{\ell}(p) = \min_{r \in \mathcal{R}} \langle p, \ell(r) \rangle = \min_{r \in \mathcal{R}} \langle p, L(\varphi(r)) \rangle = \min_{u \in \mathcal{U}} \langle p, L(u) \rangle = \underline{L|_{\mathcal{U}}},$$

143 for all $p \in \Delta_{\mathcal{Y}}$. Combining with the above, we now have $\underline{L} = \underline{\ell}$.

144 For the reverse implication, assume that $\underline{L} = \underline{\ell}$. In what follows, we implicitly work in the affine
 145 hull of $\Delta_{\mathcal{Y}}$, so that interiors are well-defined, and $\underline{\ell}$ may be differentiable on the (relative) interior of
 146 $\Delta_{\mathcal{Y}}$. Since ℓ is discrete, $-\underline{\ell}$ is polyhedral as the pointwise maximum of a finite set of linear functions.
 147 The projection of its epigraph E_{ℓ} onto $\Delta_{\mathcal{Y}}$ forms a power diagram by Theorem 4, whose cells are
 148 full-dimensional and correspond to the level sets γ_r of $\gamma = \text{prop}[\ell]$.

149 For each $r \in \mathcal{R}$, let p_r be a distribution in the interior of γ_r , and let $u_r \in \Gamma(p)$. Observe that,
 150 by definition of the Bayes risk and Γ , for all $u \in \mathbb{R}^d$ the hyperplane $v \mapsto \langle v, -L(u_r) \rangle$ supports
 151 the epigraph E_L of $-\underline{L}$ at the point $(p, -\langle p, L(u_r) \rangle)$ if and only if $u \in \Gamma(p)$. Thus, the hyperplane
 152 $v \mapsto \langle v, -L(u_r) \rangle$ supports $E_L = E_{\ell}$ at the point $(p_r, -\langle p_r, L(u_r) \rangle)$, and thus does so at the entire
 153 facet $\{(p, -\langle p, L(u_r) \rangle) : p \in \gamma_r\}$; by the above, $u_r \in \Gamma(p)$ for all such distributions as well. We
 154 conclude that $u_r \in \Gamma(p) \iff p \in \gamma_r \iff r \in \gamma(p)$, satisfying condition (3) for $\varphi : r \mapsto u_r$. To
 155 see that the loss values match, we merely note that the supporting hyperplanes to the facets of E_L
 156 and E_{ℓ} are the same, and the loss values are uniquely determined by the supporting hyperplane. (In
 157 particular, if h supports the facet corresponding to γ_r , we have $\ell(r)_y = L(u_r)_y = h(\delta_y)$, where δ_y is
 158 the point distribution on outcome y .) \square

159 From this more succinct embedding condition, we can in turn simplify the condition that a loss
 160 embeds *some* discrete loss: it does if and only if its Bayes risk is polyhedral. (We say a concave
 161 function is polyhedral if its negation is a polyhedral convex function.) Note that Bayes risk, a function
 162 from distributions over \mathcal{Y} to the reals, may be polyhedral even if the loss itself is not.

163 **Proposition 2.** *A loss L embeds a discrete loss if and only if \underline{L} is polyhedral.*

164 *Proof.* If L embeds ℓ , Proposition 1 gives us $\underline{L} = \underline{\ell}$, and its proof already argued that $\underline{\ell}$ is polyhedral.
 165 For the converse, let \underline{L} be polyhedral; we again examine the proof of Proposition 1. The projection
 166 of \underline{L} onto $\Delta_{\mathcal{Y}}$ forms a power diagram by Theorem 4 with finitely many cells C_1, \dots, C_k , which we
 167 can index by $\mathcal{R} := \{1, \dots, k\}$. Defining the property $\gamma : \Delta_{\mathcal{Y}} \rightrightarrows \mathcal{R}$ by $\gamma_r = C_r$ for $r \in \mathcal{R}$, we see
 168 that the same construction gives us points $u_r \in \mathbb{R}^d$ such that $u_r \in \Gamma(p) \iff r \in \gamma(p)$. Defining
 169 $\ell : \mathcal{R} \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ by $\ell(r) = L(u_r)$, the same proof shows that L embeds ℓ . \square

170 Combining Proposition 2 with the observation that polyhedral losses have polyhedral Bayes risks
 171 (Lemma 5), we obtain the first direction of our equivalence between polyhedral losses and embedding.

172 **Theorem 1.** *Every polyhedral loss L embeds a discrete loss.*

173 We now turn to the reverse direction: which discrete losses are embedded by some polyhedral loss?
 174 Perhaps surprisingly, we show that *every* discrete loss is embeddable, using a construction via convex
 175 conjugate duality which has appeared several times in the literature (e.g. [1, 8, 10]). Note however
 176 that the number of dimensions d required could be as large as $|\mathcal{Y}|$.

177 **Theorem 2.** *Every discrete loss ℓ is embedded by a polyhedral loss.*

178 *Proof.* Let $n = |\mathcal{Y}|$, and let $C : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $(-\underline{\ell})^*$, the convex conjugate of $-\underline{\ell}$. From
 179 standard results in convex analysis, C is polyhedral as $-\underline{\ell}$ is, and C is finite on all of \mathbb{R}^n as the
 180 domain of $-\underline{\ell}$ is bounded [27, Corollary 13.3.1]. Note that $-\underline{\ell}$ is a closed convex function, as the
 181 infimum of affine functions, and thus $(-\underline{\ell})^{**} = -\underline{\ell}$. Define $L : \mathbb{R}^n \rightarrow \mathbb{R}^{\mathcal{Y}}$ by $L(u) = C(u)\mathbb{1} - u$,
 182 where $\mathbb{1} \in \mathbb{R}^{\mathcal{Y}}$ is the all-ones vector. We first show that L embeds ℓ , and then establish that the range
 183 of L is in fact $\mathbb{R}_+^{\mathcal{Y}}$, as desired.

184 We compute Bayes risks and apply Proposition 1 to see that L embeds ℓ . For any $p \in \Delta_{\mathcal{Y}}$, we have

$$\begin{aligned} \underline{L}(p) &= \inf_{u \in \mathbb{R}^n} \langle p, C(u)\mathbb{1} - u \rangle \\ &= \inf_{u \in \mathbb{R}^n} C(u) - \langle p, u \rangle \\ &= - \sup_{u \in \mathbb{R}^n} \langle p, u \rangle - C(u) \\ &= -C^*(p) = -(-\underline{\ell}(p))^{**} = \underline{\ell}(p). \end{aligned}$$

185 It remains to show $L(u)_y \geq 0$ for all $u \in \mathbb{R}^n, y \in \mathcal{Y}$. Letting $\delta_y \in \Delta_{\mathcal{Y}}$ be the point distribution on
 186 outcome $y \in \mathcal{Y}$, we have for all $u \in \mathbb{R}^n$, $L(u)_y \geq \inf_{u' \in \mathbb{R}^n} L(u')_y = \underline{L}(\delta_y) = \underline{\ell}(\delta_y) \geq 0$, where
 187 the final inequality follows from the nonnegativity of ℓ . \square

188 4 Consistency via Calibrated Links

189 We have now seen the tight relationship between polyhedral losses and embeddings; in particular,
 190 every polyhedral loss embeds some discrete loss. The embedding itself tells us how to link the
 191 embedded points back to the discrete reports (map $\varphi(r)$ to r), but it is not clear when this link can be
 192 extended to the remaining reports, and whether such a link can lead to consistency. In this section,
 193 we give a construction to generate calibrated links for any polyhedral loss.

194 Appendix D contains the full proof; this section provides a sketch along with the main construction
 195 and result. The first step is to give a link ψ such that exactly minimizing expected surrogate loss L ,
 196 followed by applying ψ , always exactly minimizes expected original loss ℓ . The existence of such a
 197 link is somewhat subtle, because in general some point u that is far from any embedding point can
 198 minimize expected loss for two very different distributions p, p' , making it unclear whether there
 199 exists a link choice $\psi(u)$ that is simultaneously optimal for both. We show that as we vary p over $\Delta_{\mathcal{Y}}$,
 200 there are only finitely many sets of the form $U = \arg \min_{u \in \mathbb{R}^d} \langle p, L(u) \rangle$ (Lemma 4). Associating
 201 each U with $R_U \subseteq \mathcal{R}$, the set of reports whose embedding points are in U , we enforce that all points
 202 in U link to some report in R_U . (As a special case, embedding points must link to their corresponding
 203 reports.) Proving this is well-defined uses a chain of arguments involving the Bayes risk, ultimately
 204 showing that if u lies in multiple U , the corresponding report sets R_U all intersect at some $r =: \psi(u)$.

205 Intuitively, to create separation, we just need to “thicken” this construction by mapping all
 206 approximately-optimal points u to optimal reports r . Let \mathcal{U} contain all optimal report sets U
 207 of the form above. A key step in the following definition will be to narrow down a “link envelope” Ψ
 208 where $\Psi(u)$ denotes the legal or valid choices for $\psi(u)$.

209 **Definition 6.** *Given a polyhedral L that embeds some ℓ , an $\epsilon > 0$, and a norm $\|\cdot\|$, the ϵ -thickened
 210 link ψ is constructed as follows. First, initialize $\Psi : \mathbb{R}^d \rightrightarrows \mathcal{R}$ by setting $\Psi(u) = \mathcal{R}$ for all u . Then for
 211 each $U \in \mathcal{U}$, for all points u such that $\inf_{u^* \in U} \|u^* - u\| < \epsilon$, update $\Psi(u) = \Psi(u) \cap R_U$. Finally,
 212 define $\psi(u) \in \Psi(u)$, breaking ties arbitrarily. If $\Psi(u)$ became empty, then leave $\psi(u)$ undefined.*

213 **Theorem 3.** *Let L be polyhedral, and ℓ the discrete loss it embeds from Theorem 1. Then for small
 214 enough $\epsilon > 0$, the ϵ -thickened link ψ is well-defined and, furthermore, is a calibrated link from L to ℓ .*

215 *Sketch. Well-defined:* For the initial construction above, we argued that if some collection such as
 216 U, U', U'' overlap at a u , then their report sets $R_U, R_{U'}, R_{U''}$ also overlap, so there is a valid choice
 217 $r = \psi(u)$. Now, we thicken all sets $U \in \mathcal{U}$ by a small enough ϵ ; it can be shown that if the *thickened*
 218 sets overlap at u , then U, U', U'' themselves overlap, so again $R_U, R_{U'}, R_{U''}$ overlap and there is a
 219 valid choice $r = \psi(u)$.

220 *Calibrated:* By construction of the thickened link, if u maps to an incorrect report, i.e. $\psi(u) \notin \gamma(p)$,
 221 then u must have at least distance ϵ to the optimal set U . We then show that the minimal gradient
 222 of the expected loss along any direction away from U is lower-bounded, giving a constant excess
 223 expected loss at u . \square

224 5 Application to Specific Surrogates

225 Our results give a framework to construct consistent surrogates and link functions for any discrete
 226 loss, but they also provide a way to verify the consistency or inconsistency of given surrogates. Below,
 227 we illustrate the power of this framework with specific examples from the literature, as well as new
 228 examples. In some cases we simplify existing proofs, while in others we give new results, such as a
 229 new calibrated link for abstain loss, and the inconsistency of the recently proposed Lovász hinge.

230 5.1 Consistency of abstain surrogate and link construction

231 In classification settings with a large number of labels, several authors consider a variant of classifica-
 232 tion, with the addition of a “reject” or *abstain* option. For example, Ramaswamy et al. [25] study the
 233 loss $\ell_\alpha : [n] \cup \{\perp\} \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ defined by $\ell_\alpha(r)_y = 0$ if $r = y$, α if $r = \perp$, and 1 otherwise. Here, the
 234 report \perp corresponds to “abstaining” if no label is sufficiently likely, specifically, if no $y \in \mathcal{Y}$ has
 235 $p_y \geq 1 - \alpha$. Ramaswamy et al. [25] provide a polyhedral surrogate for ℓ_α , which we present here for
 236 $\alpha = 1/2$. Letting $d = \lceil \log_2(n) \rceil$ their surrogate is $L_{1/2} : \mathbb{R}^d \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ given by

$$L_{1/2}(u)_y = (\max_{j \in [d]} B(y)_j u_j + 1)_+, \quad (4)$$

237 where $B : [n] \rightarrow \{-1, 1\}^d$ is an arbitrary injection; let us assume $n = 2^d$ so that we have a bijection.
 238 Consistency is proven for the following link function,

$$\psi(u) = \begin{cases} \perp & \min_{i \in [d]} |u_i| \leq 1/2 \\ B^{-1}(\text{sgn}(-u)) & \text{otherwise} \end{cases}. \quad (5)$$

239 In light of our framework, we can see that $L_{1/2}$ is an excellent example of an embedding, where
 240 $\varphi(y) = B(y)$ and $\varphi(\perp) = 0 \in \mathbb{R}^d$. Moreover, the link function ψ can be recovered from Theorem 3
 241 with norm $\|\cdot\|_\infty$ and $\epsilon = 1/2$; see Figure 1(L). Hence, our framework would have simplified the
 242 process of finding such a link, and the corresponding proof of consistency. To illustrate this point
 243 further, we give an alternate link ψ_1 corresponding to $\|\cdot\|_1$ and $\epsilon = 1$, shown in Figure 1(R):

$$\psi_1(u) = \begin{cases} \perp & \|u\|_1 \leq 1 \\ B^{-1}(\text{sgn}(-u)) & \text{otherwise} \end{cases}. \quad (6)$$

244 Theorem 3 immediately gives calibration of $(L_{1/2}, \psi_1)$ with respect to $\ell_{1/2}$. Aside from its simplicity,
 245 one possible advantage of ψ_1 is that it appears to yield the same constant in generalization bounds as
 246 ψ , yet assigns \perp to much less of the surrogate space \mathbb{R}^d . It would be interesting to compare the two
 247 links in practice.

248 5.2 Inconsistency of top- k losses

249 In certain classification problems, for example in information retrieval, it is common to predict
 250 a set of possible labels. As one instance, for $k < n$ the top- k classification problem has reports
 251 $\mathcal{R} := \{r \subseteq [n] : |r| = k\}$, with label $y \in [n]$. The natural discrete loss $\ell_{\text{top-}k} : \mathcal{R} \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ is given by

$$\ell_{\text{top-}k}(r, y) = \mathbb{1}\{y \notin r\}, \quad (7)$$

252 which simply gives a penalty if the label was not in the reported set.

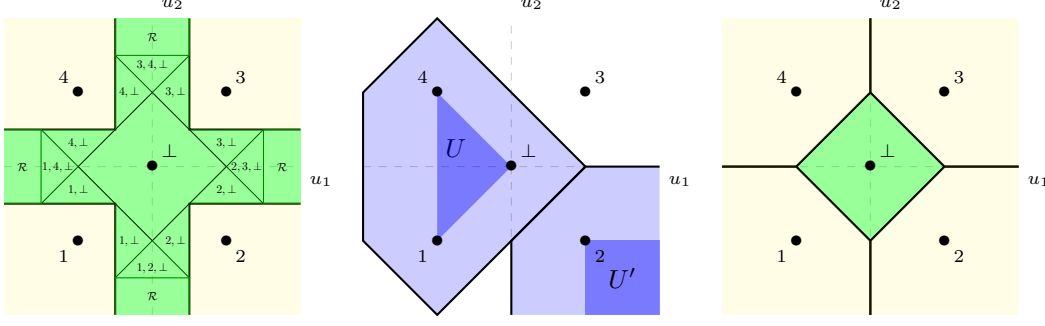


Figure 1: Constructing links for the abstain surrogate $L_{1/2}$ with $d = 2$. The embedding is shown in bold labeled by the corresponding reports. (L) The link envelope Ψ resulting from Theorem 3 using $\| \cdot \|_\infty$ and $\epsilon = 1/2$, and a possible link ψ which matches eq. (5) from [25]. (M) An illustration of the thickened sets from Definition 6 for two sets $U \in \mathcal{U}$, using $\| \cdot \|_1$ and $\epsilon = 1$. (R) The Ψ and ψ from Theorem 3 using $\| \cdot \|_1$ and $\epsilon = 1$.

253 Surrogates for this problem commonly take reports $u \in \mathbb{R}^n$, with the link $\psi(u) = \{u_{[1]}, \dots, u_{[k]}\}$,
 254 where $x_{[u]}$ is the i^{th} largest component of u , with ties broken arbitrarily. Lapin et al. [18, 19, 20]
 255 provide the following convex surrogate loss for this problem, which Yang and Koyejo [31] show to
 256 be inconsistent:¹

$$L'(u)_y := \left(\frac{1}{k} \sum_{i=1}^k (u + \mathbb{1} - e_y)_{[i]} - u_y \right)_+, \quad (8)$$

257 where e_y is 1 in component y and 0 elsewhere.

258 With our framework, we can say more. Specifically, while (L', ψ) is not consistent for $\ell_{\text{top-}k}$, since L'
 259 is polyhedral (Lemma 11), we know from Theorem 1 that it embeds *some* discrete loss ℓ' , and from
 260 Theorem 3 there is a link ψ' such that (L', ψ') is calibrated (and consistent) for ℓ' . We therefore turn
 261 to deriving this discrete loss ℓ' .

262 In Lemma 12, we show that the set $\mathcal{U} = \{u \in \{0, 1\}^n : \|u\|_0 \leq k\}$ is always represented among
 263 the optimizers of L' , meaning for all $p \in \Delta_{\mathcal{Y}}$ we have $\text{prop}[L'](p) \cap \mathcal{U} \neq \emptyset$. From Lemma 3,
 264 then, L' must embed $L'|_{\mathcal{U}}$, which gives us $\ell' : \mathcal{R}' \rightarrow \mathbb{R}_{+}^{\mathcal{Y}}$ via the natural embedding from sets
 265 $\mathcal{R}' = \{r \subseteq [n] : |r| \leq k\}$ to their indicator vectors:

$$\ell'(r)_y = \mathbb{1}\{y \notin r\} + \frac{1}{k}|r \setminus \{y\}| = (1 + \frac{1}{k})\mathbb{1}\{y \notin r\} + \frac{1}{k}(|r| - 1). \quad (9)$$

266 We now see that ℓ' is essentially $\ell_{\text{top-}k}$ (extended to sets smaller than k) plus an additional cardinality
 267 term which rewards smaller sets.

268 Knowledge of what loss L' actually embeds greatly simplifies the task of proving inconsistency
 269 with $\ell_{\text{top-}k}$. Specifically, we see that ℓ' allows sets of cardinality strictly less than k , so we can look
 270 for a distribution making one of these smaller sets optimal. Writing the expected loss, we have
 271 $\langle p, \ell'(r) \rangle = (1 + \frac{1}{k})(1 - p(r)) + \frac{1}{k}(|r| - 1) = 1 + \frac{1}{k}|r| - (1 + \frac{1}{k})p(r)$, where $p(r) = \sum_{i \in r} p_i$. So let
 272 us ask when it is better to drop an element i from r : we have $\langle p, \ell'(r) - \ell'(r \setminus \{i\}) \rangle = \frac{1}{k} - (1 + \frac{1}{k})p_i$,
 273 meaning we will drop elements from r until they all have weight at least $\frac{1}{k+1}$. In particular, for
 274 distributions close to uniform, $r = \emptyset$ is optimal, already giving us inconsistency. More generally, as
 275 $k < n$, the set $P = \{p \in \Delta_{\mathcal{Y}} : \max_{r \in \mathcal{R}'} p(r) < k/(k+1)\}$ is full-dimensional and guarantees that
 276 at least one of the top k labels has probability strictly less than $\frac{1}{k+1}$.

277 5.3 Inconsistency of Lovász hinge

278 Many structured prediction settings can be thought of as making multiple predictions at once, with
 279 a loss function that jointly measures error based on the relationship between these predictions [13,
 280 15, 22]. In the case of k binary predictions, these settings are typically formalized by taking the
 281 predictions and outcomes to be ± 1 vectors, so $\mathcal{R} = \mathcal{Y} = \{-1, 1\}^k$. One then defines a joint
 282 loss function, which is often merely a function of the set of mispredictions, meaning $\ell^g(r)_y =$

¹Yang and Koyejo also introduce a consistent surrogate, but it is non-convex.

283 $g(\{i \in [k] : r_i \neq y_i\})$ for some function $g : 2^{[k]} \rightarrow \mathbb{R}$. For example, Hamming loss is simply
 284 given by $g(S) = |S|$. In an effort to provide a general convex surrogate for these settings when g
 285 is a submodular function, Yu and Blaschko [32] introduce the *Lovász hinge*, which leverages the
 286 well-known convex Lovász extension of submodular functions. While the authors provide theoretical
 287 justification and experiments, consistency of the Lovász hinge is left open, which we resolve.

288 Rather than formally define the Lovász hinge, we defer the complete analysis to the full version of
 289 the paper,² and focus here on the $k = 2$ case. For brevity, we write $g_\emptyset := g(\emptyset)$, $g_{1,2} := g(\{1, 2\})$, etc.
 290 Assuming g is normalized and increasing (meaning $g_{1,2} \geq \{g_1, g_2\} \geq g_\emptyset = 0$), the Lovász hinge
 291 $L : \mathbb{R}^k \rightarrow \mathbb{R}_+^{\mathcal{Y}}$ is given by

$$L^g(u)_y = \max \left\{ (1 - u_1 y_1)_+ g_1 + (1 - u_2 y_2)_+ (g_{1,2} - g_1), \right. \\ \left. (1 - u_2 y_2)_+ g_2 + (1 - u_1 y_1)_+ (g_{1,2} - g_2) \right\}, \quad (10)$$

292 where $(x)_+ = \max\{x, 0\}$. We will explore the range of values of g for which L^g is consistent, where
 293 the link function $\psi : \mathbb{R}^2 \rightarrow \{-1, 1\}^2$ is fixed as $\psi(u)_i = \text{sgn}(u_i)$, with ties broken arbitrarily.

294 Let us consider $g_\emptyset = 0$, $g_1 = g_2 = g_{1,2} = 1$, for which ℓ^g is merely 0-1 loss on \mathcal{Y} . For consistency,
 295 then, for any distribution $p \in \Delta_{\mathcal{Y}}$, we must have that whenever $u \in \arg \min_{u' \in \mathbb{R}^2} p \cdot L^g(u')$, then
 296 $\psi(u)$ must be the most likely outcome, in $\arg \max_{y \in \mathcal{Y}} p(y)$. Simplifying eq. (10), however, we have

$$L^g(u)_y = \max \left\{ (1 - u_1 y_1)_+, (1 - u_2 y_2)_+ \right\} = \max \left\{ 1 - u_1 y_1, 1 - u_2 y_2, 0 \right\}, \quad (11)$$

297 which is exactly the abstain surrogate (4) for $d = 2$. We immediately conclude that L^g cannot be
 298 consistent with ℓ^g , as the origin will be the unique optimal report for L^g under distributions with
 299 $p_y < 0.5$ for all y , and one can simply take a distribution which disagrees with the way ties are broken
 300 in ψ . For example, if we take $\text{sgn}(0) = 1$, then under $p((1, 1)) = p((1, -1)) = p((-1, 1)) = 0.2$
 301 and $p((-1, -1)) = 0.4$, we have $\{0\} = \arg \min_{u \in \mathbb{R}^2} p \cdot L^g(u)$, yet $\psi(0) = (1, 1) \notin \{(-1, -1)\} =$
 302 $\arg \min_{r \in \mathcal{R}} p \cdot \ell^g(r)$.

303 In fact, this example is typical: using our embedding framework, and characterizing when $0 \in \mathbb{R}^2$ is
 304 an embedded point, one can show that L^g is consistent if and only if $g_{1,2} = g_1 + g_2$. Moreover, in the
 305 linear case, which corresponds to g being *modular*, the Lovász hinge reduces to weighted Hamming
 306 loss, which is trivially consistent from the consistency of hinge loss for 0-1 loss. In the full version of
 307 the paper, we generalize this observation for all k : L^g is consistent if and only if g is modular. In
 308 other words, even for $k > 2$, the only consistent Lovász hinge is weighted Hamming loss. These
 309 results cast doubt on the effectiveness of the Lovász hinge in practice.

310 6 Conclusion and Future Directions

311 This paper formalizes an intuitive way to design convex surrogate losses for classification-like
 312 problems—by embedding the reports into \mathbb{R}^d . We establish a close relationship between embedding
 313 and polyhedral surrogates, showing both that every polyhedral loss embeds a discrete loss (Theorem 1)
 314 and that every discrete loss is embedded by some polyhedral loss (Theorem 2). We then construct a
 315 calibrated link function from any polyhedral loss to the discrete loss it embeds, giving consistency
 316 for all such losses. We conclude with examples of how the embedding framework presented can be
 317 applied to understand existing surrogates in the literature, including those for the abstain loss, top- k
 318 loss, and Lovász hinge.

319 One open question of particular interest involves the dimension of the input to a surrogate; given
 320 a discrete loss, can we construct the surrogate that embeds it *of minimal dimension*? If we naïvely
 321 embed the reports into an n -dimensional space, the dimensionality of the problem scales with the
 322 number of possible labels n . As the dimension of the optimization problem is a function of this
 323 *embedding dimension* d , a promising direction is to leverage tools from elicitation complexity [12, 17]
 324 and convex calibration dimension [24] to understand when we can take $d \ll n$.

²Citation withheld for anonymity.

References

- [1] Jacob Abernethy, Yiling Chen, and Jennifer Wortman Vaughan. Efficient market making via convex optimization, and a connection to online learning. *ACM Transactions on Economics and Computation*, 1(2):12, 2013. URL <http://dl.acm.org/citation.cfm?id=2465777>.
- [2] Arpit Agarwal and Shivani Agarwal. On consistent surrogate risk minimization and property elicitation. In *JMLR Workshop and Conference Proceedings*, volume 40, pages 1–19, 2015. URL <http://www.jmlr.org/proceedings/papers/v40/Agarwal15.pdf>.
- [3] Franz Aurenhammer. Power diagrams: properties, algorithms and applications. *SIAM Journal on Computing*, 16(1):78–96, 1987. URL <http://epubs.siam.org/doi/pdf/10.1137/0216006>.
- [4] Peter L Bartlett and Marten H Wegkamp. Classification with a reject option using a hinge loss. *Journal of Machine Learning Research*, 9(Aug):1823–1840, 2008.
- [5] Peter L. Bartlett, Michael I. Jordan, and Jon D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006. URL <http://amstat.tandfonline.com/doi/abs/10.1198/016214505000000907>.
- [6] S.P. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [7] Koby Crammer and Yoram Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of machine learning research*, 2(Dec):265–292, 2001.
- [8] John Duchi, Khashayar Khosravi, Feng Ruan, et al. Multiclass classification, information, divergence and surrogate risk. *The Annals of Statistics*, 46(6B):3246–3275, 2018.
- [9] Tobias Fissler, Johanna F Ziegel, and others. Higher order elicibility and Osband’s principle. *The Annals of Statistics*, 44(4):1680–1707, 2016.
- [10] Rafael Frongillo and Ian Kash. General truthfulness characterizations via convex analysis. In *Web and Internet Economics*, pages 354–370. Springer, 2014.
- [11] Rafael Frongillo and Ian Kash. Vector-Valued Property Elicitation. In *Proceedings of the 28th Conference on Learning Theory*, pages 1–18, 2015.
- [12] Rafael Frongillo and Ian A. Kash. On Elicitation Complexity. In *Advances in Neural Information Processing Systems 29*, 2015.
- [13] Wei Gao and Zhi-Hua Zhou. On the consistency of multi-label learning. In *Proceedings of the 24th annual conference on learning theory*, pages 341–358, 2011.
- [14] T. Gneiting. Making and Evaluating Point Forecasts. *Journal of the American Statistical Association*, 106(494):746–762, 2011.
- [15] Tamir Hazan, Joseph Keshet, and David A McAllester. Direct loss minimization for structured prediction. In *Advances in Neural Information Processing Systems*, pages 1594–1602, 2010.
- [16] Nicolas S. Lambert. Elicitation and evaluation of statistical forecasts. 2018. URL <https://web.stanford.edu/~nlambert/papers/elicitability.pdf>.
- [17] Nicolas S. Lambert, David M. Pennock, and Yoav Shoham. Eliciting properties of probability distributions. In *Proceedings of the 9th ACM Conference on Electronic Commerce*, pages 129–138, 2008.
- [18] Maksim Lapin, Matthias Hein, and Bernt Schiele. Top-k multiclass svm. In *Advances in Neural Information Processing Systems*, pages 325–333, 2015.

- [19] Maksim Lapin, Matthias Hein, and Bernt Schiele. Loss functions for top-k error: Analysis and insights. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1468–1477, 2016.
- [20] Maksim Lapin, Matthias Hein, and Bernt Schiele. Analysis and optimization of loss functions for multiclass, top-k, and multilabel classification. *IEEE transactions on pattern analysis and machine intelligence*, 40(7):1533–1554, 2018.
- [21] Kent Osband and Stefan Reichelstein. Information-eliciting compensation schemes. *Journal of Public Economics*, 27(1):107–115, June 1985. ISSN 0047-2727. doi: 10.1016/0047-2727(85)90031-3. URL <http://www.sciencedirect.com/science/article/pii/0047272785900313>.
- [22] Anton Osokin, Francis Bach, and Simon Lacoste-Julien. On structured prediction theory with calibrated convex surrogate losses. In *Advances in Neural Information Processing Systems*, pages 302–313, 2017.
- [23] Bernardo Avila Pires, Csaba Szepesvari, and Mohammad Ghavamzadeh. Cost-sensitive multiclass classification risk bounds. In *International Conference on Machine Learning*, pages 1391–1399, 2013.
- [24] Harish G Ramaswamy and Shivani Agarwal. Convex calibration dimension for multiclass loss matrices. *The Journal of Machine Learning Research*, 17(1):397–441, 2016.
- [25] Harish G Ramaswamy, Ambuj Tewari, Shivani Agarwal, et al. Consistent algorithms for multiclass classification with an abstain option. *Electronic Journal of Statistics*, 12(1):530–554, 2018.
- [26] M.D. Reid and R.C. Williamson. Composite binary losses. *The Journal of Machine Learning Research*, 9999:2387–2422, 2010.
- [27] R.T. Rockafellar. *Convex analysis*, volume 28 of *Princeton Mathematics Series*. Princeton University Press, 1997.
- [28] L.J. Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, pages 783–801, 1971.
- [29] Ingo Steinwart, Chloé Pasin, Robert Williamson, and Siyu Zhang. Elicitation and Identification of Properties. In *Proceedings of The 27th Conference on Learning Theory*, pages 482–526, 2014.
- [30] Robert C Williamson, Elodie Vernet, and Mark D Reid. Composite multiclass losses. *Journal of Machine Learning Research*, 17(223):1–52, 2016.
- [31] Forest Yang and Sanmi Koyejo. On the consistency of top-k surrogate losses, 01 2018.
- [32] Jiaqian Yu and Matthew B Blaschko. The lovász hinge: A novel convex surrogate for submodular losses. *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [33] Ming Yuan and Marten Wegkamp. Classification methods with reject option based on convex risk minimization. *Journal of Machine Learning Research*, 11(Jan):111–130, 2010.
- [34] Chong Zhang, Wenbo Wang, and Xingye Qiao. On reject and refine options in multicategory classification. *Journal of the American Statistical Association*, 113(522):730–745, 2018. doi: 10.1080/01621459.2017.1282372. URL <https://doi.org/10.1080/01621459.2017.1282372>.