



**Problem:** Characterise the Bayes-optimal scorers for the bipartite ranking risk with a surrogate loss  $\ell$ .

**Approach:** Exploit the reduction of bipartite ranking to classification over pairs, and the machinery of proper composite losses.

**Results:** Under a condition on the link function for  $\ell$ , we obtain the Bayes-optimal scorer, and surrogate regret bounds. Bayes-optimal scorers can also be established more generally, including for the  $p$ -norm push risk.

## Bipartite ranking

Class-conditional densities and base rate

**Input** IID samples from  $D = (P, Q, \pi)$  over  $\mathcal{X} \times \{\pm 1\}$

**Output** Scorer  $s : \mathcal{X} \rightarrow \mathbb{R}$

**Performance**  $\text{AUC}^D(s) = \mathbb{E}_{X \sim P, X' \sim Q} \left[ \mathbb{I}[s(X) > s(X')] + \frac{1}{2} \mathbb{I}[s(X) = s(X')] \right]$

AUC maximisation is challenging due to the non-convex indicator function. Typically, for a **surrogate loss**  $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$ , one minimises

$$\mathbb{L}_{\text{Bipart}, \ell}^D(s) := \frac{1}{2} \cdot \mathbb{E}_{X \sim P, X' \sim Q} [\ell_1(s(X) - s(X')) + \ell_{-1}(s(X') - s(X))].$$

This is intuitive, but is it consistent for AUC maximisation?

## Bayes-optimal scorers

For a loss  $\ell$ , the **Bayes-optimal scorers** are minimisers of the bipartite risk:

$$\mathcal{S}_{\text{Bipart}, \ell}^{D,*} := \underset{s : \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{L}_{\text{Bipart}, \ell}^D(s).$$

For consistency, we minimally need  $\mathcal{S}_{\text{Bipart}, \ell}^{D,*} \cap \mathcal{S}_{\text{Bipart}, 01}^{D,*} \neq \emptyset$ . The Neyman-Pearson lemma implies that  $\mathcal{S}_{\text{Bipart}, 01}^{D,*}$  comprises all monotone transformations of  $\eta : x \mapsto \Pr[Y = 1 | X = x]$ . **What is  $\mathcal{S}_{\text{Bipart}, \ell}^{D,*}$ ?**

We answer this for **proper composite**  $\ell$  with invertible link  $\Psi : [0, 1] \rightarrow \mathbb{R}$ . These are the fundamental losses of class-probability estimation, since

$$\mathcal{S}_{\text{Class}, \ell}^{D,*} = \Psi \circ \eta.$$

## Reduction to classification

For the pair-classification distribution  $\text{Bipart}(D) = (P \times Q, Q \times P, 1/2)$ ,

$$\mathbb{L}_{\text{Bipart}, \ell}^D(s) = \mathbb{L}_{\text{Class}, \ell}^{\text{Bipart}(D)}(\text{Diff}(s)) = \mathbb{E}_{((X, X'), Y) \sim \text{Bipart}(D)} [\ell(Y, (\text{Diff}(s))(X, X'))],$$

where  $\text{Diff}(s) : (x, x') \mapsto s(x) - s(x')$ .

It now seems we can simply compute  $\mathcal{S}_{\ell}^{\text{Bipart}(D),*}$ . However, the risk only considers a restricted function class of **decomposable** scorers,  $\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) : s : \mathcal{X} \rightarrow \mathbb{R}\}$ . Thus, computing  $\mathcal{S}_{\ell}^{\text{Bipart}(D),*}$  via the conditional risk requires  $\mathcal{S}_{\ell}^{\text{Bipart}(D),*} \subseteq \mathcal{S}_{\text{Decomp}}$ . **When does this happen?**

## Decomposable solutions

An innocuous lemma will prove important. Let  $\sigma(\cdot)$  denote the sigmoid.

**Lemma.** The observation-conditional density for  $\text{Bipart}(D)$  is

$$\eta_{\text{Pair}} = \sigma \circ \text{Diff}(\sigma^{-1} \circ \eta).$$

Consequently, for a proper composite  $\ell$ , the optimal pair-scorer must be

$$\mathcal{S}_{\ell}^{\text{Bipart}(D),*} = \Psi \circ \sigma \circ \text{Diff}(\sigma^{-1} \circ \eta).$$

Thus, **under a condition on the link function**, we can seamlessly import results from classification.

**Proposition.** Given any strictly proper composite loss  $\ell$  with a differentiable, invertible link function  $\Psi$  such that  $(\exists a \in \mathbb{R}_+) \Psi^{-1} : v \mapsto \frac{1}{1+e^{-av}}$ ,

$$(A) \mathcal{S}_{\text{Bipart}, \ell}^{D,*} = \{\Psi \circ \eta + b : b \in \mathbb{R}\} \subseteq \mathcal{S}_{\text{Bipart}, 01}^{D,*}$$

$$(B) \exists \text{ convex } F_{\ell} : [0, 1] \rightarrow \mathbb{R}_+ \text{ such that}$$

$$(\forall D, s : \mathcal{X} \rightarrow \mathbb{R}) F_{\ell}(\text{regret}_{\text{Bipart}, 01}^D(s)) \leq \text{regret}_{\text{Bipart}, \ell}^D(s).$$

## Non-decomposable solutions

When  $\mathcal{S}_{\ell}^{\text{Bipart}(D),*} \cap \mathcal{S}_{\text{Decomp}} = \emptyset$ , the Bayes-optimal scorers may still be computed by explicitly differentiating the risk.

**Proposition.** Given any  $D$  and strictly proper composite loss  $\ell(y, v) = \phi(yv)$  with  $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$  convex, if  $\phi'$  is bounded, or  $D$  has finite support,

$$\mathcal{S}_{\text{Bipart}, \ell}^{D,*} = \{s^* : \mathcal{X} \rightarrow \mathbb{R} : \eta = f_{s^*}^D \circ s^*\},$$

$$f_{s^*}^D : v \mapsto \frac{\pi \mathbb{E}_{X \sim P} [\ell'_{-1}(v - s^*(X))]}{\pi \mathbb{E}_{X \sim P} [\ell'_{-1}(v - s^*(X))] - (1 - \pi) \mathbb{E}_{X' \sim Q} [\ell'_{-1}(v - s^*(X'))]}.$$

Further,  $f_{s^*}^D$  is invertible if  $(\forall v \in \mathcal{V}) \phi'(v) = 0 \iff \phi'(-v) \neq 0$ .

Here, the optimal scorer is a monotone transform of  $\eta$ , albeit with a **distribution dependent** link function.

## The p-norm push risk

The **p-norm push risk** aims to focus effort on highly ranked instances:

$$\mathbb{L}_{\text{push}, \ell, g}^D(s) = \mathbb{E}_{X' \sim Q} [(\mathbb{E}_{X \sim P} [\ell_1(s(X) - s(X'))])^p], p \in [1, \infty).$$

Under exponential loss, the optimal scorer is a scaling of  $\sigma^{-1} \circ \eta$ .

**Proposition.** For any  $D$  with finite support, with  $\ell^{\text{exp}}(y, v) = e^{-yv}$ ,

$$\mathcal{S}_{\text{push}, \ell^{\text{exp}}, g^p}^{D,*} = \left\{ \frac{1}{p+1} (\sigma^{-1} \circ \eta) + b : b \in \mathbb{R} \right\}.$$

This is identical to the optimal scorer for the **p-classification loss**,  $\ell(v) = (e^{vp}/p, e^{-v})$ , which has asymmetric **weight function** over misclassification costs  $w : c \mapsto \left( (p+1) \cdot c^{1+\frac{1}{p+1}} (1-c)^{2-\frac{1}{p+1}} \right)^{-1}$ . Thus, the loss is seen to **focus on instances with high  $\eta$** .

## Risk equivalences

Our results imply that the following risk minimisers are **equivalent**, being the same monotone transform of  $\eta$ . This highlights the close connections between class-probability estimation and bipartite ranking.

$$(1) \text{Diff} \left( \underset{s : \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{(X, Y) \sim D} [e^{-Ys(X)}] \right) \quad (2) \text{Diff} \left( \underset{s : \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{X \sim P, X' \sim Q} [e^{-(s(X) - s(X'))}] \right) \\ (3) \underset{s_{\text{Pair}} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{X \sim P, X' \sim Q} [e^{-s_{\text{Pair}}(X, X')}] \quad (4) \text{Diff} \left( \underset{s : \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{X' \sim Q} \left[ \left( \mathbb{E}_{X \sim P} [e^{-(s(X) - s(X'))}] \right)^p \right] \right)$$