

Bipartite Ranking: Risk, Optimality, and Equivalences

Aditya Krishna Menon Robert C. Williamson

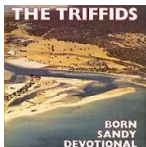
National ICT Australia and The Australian National University



Australian
National
University

August 7, 2014

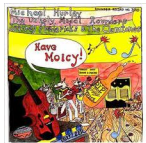
Binary classification



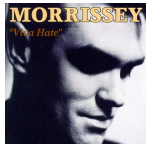
+



-



-

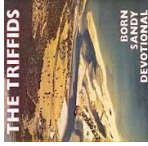


+



+

Binary class-probability estimation



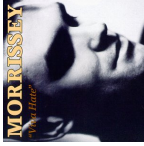
+



-



-



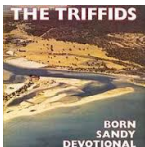
+



0.8



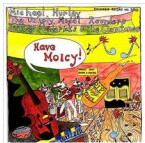
Bipartite ranking



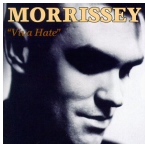
+



-



-



+



<

>

>

<



Take-home messages

Bipartite ranking = classification over pairs

Decomposability \rightarrow Bayes-optimal scorers, risk equivalences

Some risk equivalences hold for restricted function classes

- Algorithmic implications for bipartite ranking and its generalisations

Outline

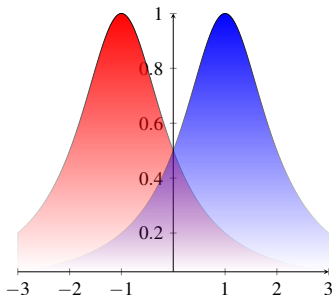
- 1 The classification risk
- 2 The bipartite risk
- 3 Decomposability and risk minimisers
- 4 Risk equivalences and algorithmic implications
- 5 Conclusion

Distributions for learning with binary labels

Instance space \mathcal{X} (e.g. \mathbb{R}^N)

Let $D = D_{P,Q,\pi}$ be a distribution over $\mathcal{X} \times \{\pm 1\}$, where

$$(P(x), Q(x), \pi) = (\Pr[X = x | Y = 1], \Pr[X = x | Y = -1], \Pr[Y = 1])$$



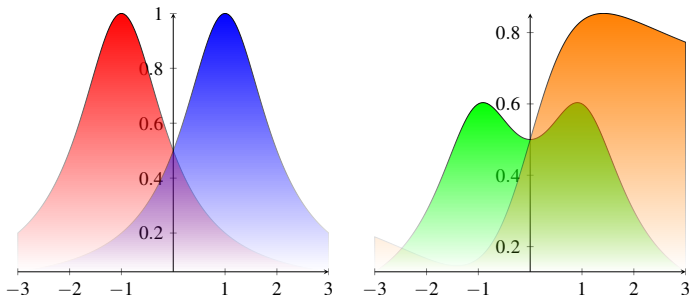
Distributions for learning with binary labels

Instance space \mathcal{X} (e.g. \mathbb{R}^N)

Let $D = D_{P,Q,\pi} = D_{M,\eta}$ be a distribution over $\mathcal{X} \times \{\pm 1\}$, where

$$(P(x), Q(x), \pi) = (\Pr[X = x | Y = 1], \Pr[X = x | Y = -1], \Pr[Y = 1])$$

$$(M(x), \eta(x)) = (\Pr[X = x], \Pr[Y = 1 | X = x])$$



Binary classification

Input IID samples from D over $\mathcal{X} \times \{\pm 1\}$

Output Classifier $c : \mathcal{X} \rightarrow \{\pm 1\}$

Risk Misclassification rate:

$$\mathbb{L}_{\text{Class}}^D(c) = \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim D} [\mathbb{I}[\mathbf{Y} \neq c(\mathbf{X})]]$$

Binary classification with scorers

Input IID samples from D over $\mathcal{X} \times \{\pm 1\}$

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

Risk Misclassification rate:

$$\mathbb{L}_{\text{Class}}^D(s) = \mathbb{E}_{(\mathbf{X}, Y) \sim D} \left[\mathbb{I}[Y \cdot s(\mathbf{X}) < 0] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = 0] \right]$$

Surrogate classification risk

Classification risk is

$$\begin{aligned}\mathbb{L}_{\text{Class}}^D(s) &= \mathbb{E}_{(X,Y) \sim D} \left[\mathbb{I}[Y \cdot s(X) < 0] + \frac{1}{2} \mathbb{I}[s(X) = 0] \right] \\ &= \mathbb{E}_{(X,Y) \sim D} \left[\ell^{01}(Y, s(X)) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Surrogate classification risk

Classification risk is

$$\begin{aligned}\mathbb{L}_{\text{Class}}^D(s) &= \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim D} \left[\mathbb{I}[\mathbf{Y} \cdot s(\mathbf{X}) < 0] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = 0] \right] \\ &= \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim D} \left[\ell^{01}(\mathbf{Y}, s(\mathbf{X})) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Solution: for **surrogate loss** $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$, minimise

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

Surrogate classification risk

Classification risk is

$$\begin{aligned}\mathbb{L}_{\text{Class}}^D(s) &= \mathbb{E}_{(X,Y) \sim D} \left[\mathbb{I}[Y \cdot s(X) < 0] + \frac{1}{2} \mathbb{I}[s(X) = 0] \right] \\ &= \mathbb{E}_{(X,Y) \sim D} \left[\ell^{01}(Y, s(X)) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Solution: for **surrogate loss** $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$, minimise

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

What is a suitable surrogate loss?

Bayes-optimal scorers

Bayes-optimal scorers for the surrogate classification risk:

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

Minimally, surrogate should preserve optimal solutions of ℓ^{01} :

$$\mathcal{S}_{\text{Class},\ell}^{D,*} \subseteq \mathcal{S}_{\text{Class},01}^{D,*}$$

Bayes-optimal scorers: ℓ^{01}

Bayes-optimal scorer for ℓ^{01} :

$$\mathcal{S}_{\text{Class},01}^{D,*} = \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Class},01}^D(s)$$

Bayes-optimal scorers: ℓ^{01}

Bayes-optimal scorer for ℓ^{01} :

$$\begin{aligned} \mathcal{S}_{\text{Class},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Class},01}^D(s) \\ &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{E}_{\mathbf{X} \sim M} [L(\boldsymbol{\eta}(\mathbf{X}), s(\mathbf{X}))] \end{aligned}$$

where

$$L(\boldsymbol{\eta}, s) = \eta \llbracket s < 0 \rrbracket + (1 - \eta) \llbracket s > 0 \rrbracket + \frac{1}{2} \llbracket s = 0 \rrbracket$$

Bayes-optimal scorers: ℓ^{01}

Bayes-optimal scorer for ℓ^{01} :

$$\begin{aligned} \mathcal{S}_{\text{Class},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Class},01}^D(s) \\ &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{E}_{\mathbf{X} \sim M} [L(\boldsymbol{\eta}(\mathbf{X}), s(\mathbf{X}))] \\ &= \left\{ s: \mathcal{X} \rightarrow \mathbb{R} \mid s: x \mapsto \underset{v}{\text{Argmin}} L(\boldsymbol{\eta}(x), v) \right\} \end{aligned}$$

where

$$L(\boldsymbol{\eta}, s) = \boldsymbol{\eta} \llbracket s < 0 \rrbracket + (1 - \boldsymbol{\eta}) \llbracket s > 0 \rrbracket + \frac{1}{2} \llbracket s = 0 \rrbracket$$

Bayes-optimal scorers: ℓ^{01}

Bayes-optimal scorer for ℓ^{01} :

$$\begin{aligned} S_{\text{Class},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Class},01}^D(s) \\ &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{E}_{\mathbf{X} \sim M} [L(\boldsymbol{\eta}(\mathbf{X}), s(\mathbf{X}))] \\ &= \left\{ s: \mathcal{X} \rightarrow \mathbb{R} \mid s: x \mapsto \underset{v}{\text{Argmin}} L(\boldsymbol{\eta}(x), v) \right\} \\ &= \{s: \mathcal{X} \rightarrow \mathbb{R} \mid \text{sign}(s(x)) = \text{sign}(2\boldsymbol{\eta}(x) - 1)\} \end{aligned}$$

where

$$L(\boldsymbol{\eta}, s) = \boldsymbol{\eta} \llbracket s < 0 \rrbracket + (1 - \boldsymbol{\eta}) \llbracket s > 0 \rrbracket + \frac{1}{2} \llbracket s = 0 \rrbracket$$

Decision boundary is determined by $\boldsymbol{\eta}$

- Motivation for focussing on instances with $\boldsymbol{\eta}(x) \approx \frac{1}{2}$

Bayes-optimal scorers: proper composite ℓ

Call ℓ **strictly proper composite** if

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

for some invertible **link function** $\Psi : [0, 1] \rightarrow \mathbb{R}$

- Logistic loss: $\Psi^{-1} : v \mapsto \sigma(v)$
- Exponential loss: $\Psi^{-1} : v \mapsto \sigma(2v)$

Outline

- 1 The classification risk
- 2 The bipartite risk**
- 3 Decomposability and risk minimisers
- 4 Risk equivalences and algorithmic implications
- 5 Conclusion

Bipartite ranking

Input IID samples from D over $\mathcal{X} \times \{\pm 1\}$

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

Risk Fraction of **discordant pairs**:

$$\mathbb{L}_{\text{Bipart}}^D(s) = \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right]$$

where $P = \Pr[\mathbf{X} | Y = 1]$, $Q = \Pr[\mathbf{X} | Y = -1]$

Bipartite ranking

Input IID samples from D over $\mathcal{X} \times \{\pm 1\}$

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

Risk Fraction of **discordant pairs**:

$$\mathbb{L}_{\text{Bipart}}^D(s) = \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right]$$

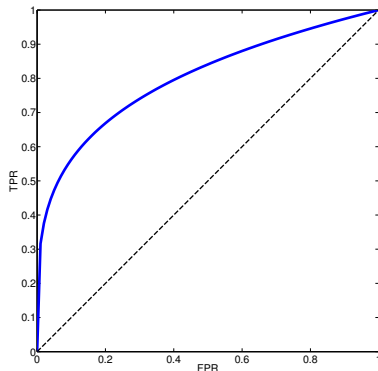
where $P = \Pr[\mathbf{X} | Y = 1]$, $Q = \Pr[\mathbf{X} | Y = -1]$

Intuitively, s ranks instances by “how positive” they are

Bipartite risk and AUC

$$\mathbb{L}_{\text{Bipart}}^D(s) = 1 - \text{AUC}^D(s)$$

- Minimising bipartite risk \rightarrow maximising AUC



Surrogate bipartite risk

Bipartite risk is

$$\begin{aligned}\mathbb{L}_{\text{Bipart}}^D(s) &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1^{01}(s(\mathbf{X}) - s(\mathbf{X}')) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Surrogate bipartite risk

Bipartite risk is

$$\begin{aligned}\mathbb{L}_{\text{Bipart}}^D(s) &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1^{01}(s(\mathbf{X}) - s(\mathbf{X}')) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Solution: for **surrogate loss** $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$ minimise,

$$\mathbb{L}_{\text{Bipart}, \ell}^D(s) = \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1(s(\mathbf{X}) - s(\mathbf{X}')) \right]$$

Surrogate bipartite risk

Bipartite risk is

$$\begin{aligned}\mathbb{L}_{\text{Bipart}}^D(s) &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1^{01}(s(\mathbf{X}) - s(\mathbf{X}')) \right]\end{aligned}$$

Problem: $\ell^{01} \rightarrow$ discontinuous, non-convex

Solution: for **surrogate loss** $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$ minimise,

$$\mathbb{L}_{\text{Bipart}, \ell}^D(s) = \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1(s(\mathbf{X}) - s(\mathbf{X}')) \right]$$

What is a suitable surrogate loss?

Bayes-optimal bipartite scorers

Bayes-optimal scorers for the bipartite risk wrt loss ℓ :

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

Minimally, would like agreement with optimal scorers for ℓ^{01} :

$$\mathcal{S}_{\text{Bipart},\ell}^{D,*} \subseteq \mathcal{S}_{\text{Bipart},01}^{D,*}$$

Bayes-optimal bipartite scorers

Bayes-optimal scorers for the bipartite risk:

$$\begin{aligned} \mathcal{S}_{\text{Bipart},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Bipart},01}^D(s) \\ &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1^{01}(s(\mathbf{X}) - s(\mathbf{X}')) \right] \end{aligned}$$

Bayes-optimal bipartite scorers

Bayes-optimal scorers for the bipartite risk:

$$\begin{aligned}\mathcal{S}_{\text{Bipart},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{L}_{\text{Bipart},01}^D(s) \\ &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmin}} \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\ell_1^{01}(s(\mathbf{X}) - s(\mathbf{X}')) \right] \\ &= ?\end{aligned}$$

Pointwise analysis?

Bayes-optimal bipartite scorers

AUC connection obviates need for conditional risk:

$$\begin{aligned} S_{\text{Bipart},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmax}} \text{AUC}^D(s) \\ &= \{s: \mathcal{X} \rightarrow \mathbb{R} \mid \eta = \phi \circ s, \phi \text{ monotone increasing} \} \end{aligned}$$

by **Neyman-Pearson** lemma

“Lax” class-probability estimation

- Only care about **ordering** induced by η

Bayes-optimal bipartite scorers

AUC connection obviates need for conditional risk:

$$\begin{aligned} S_{\text{Bipart},01}^{D,*} &= \underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{Argmax}} \text{AUC}^D(s) \\ &= \{s: \mathcal{X} \rightarrow \mathbb{R} \mid \eta = \phi \circ s, \phi \text{ monotone increasing} \} \end{aligned}$$

by **Neyman-Pearson** lemma

“Lax” class-probability estimation

- Only care about **ordering** induced by η

General ℓ ?

The Bipart(D) distribution

Given D , define a distribution over pairs, $\text{Bipart}(D)$ via:

- $(X, Y) \sim D$
- $(X', Y') \sim D$
- If $Y = Y'$, reject and repeat; else, $Z = 2\mathbb{I}[Y > Y'] - 1$.

Class-conditionals and base rate are

$$(P_{\text{pair}}, Q_{\text{pair}}, \pi_{\text{pair}}) = \left(P \times Q, Q \times P, \frac{1}{2} \right)$$

From scorers to pair-scorers

Given some $s: \mathcal{X} \rightarrow \mathbb{R}$, let

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

- Converts a scorer to a **pair-scorer**

The set of **decomposable** pair-scorers:

$$\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) \mid s: \mathcal{X} \rightarrow \mathbb{R}\}$$

The bipartite risk revisited

We can rewrite the bipartite risk as

$$\mathbb{L}_{\text{Bipart},\ell}^D(s) = \mathbb{E}_{\mathbf{X} \sim P, \mathbf{X}' \sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right]$$

The bipartite risk revisited

We can rewrite the bipartite risk as

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^D(s) &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\ell_1^{01}(\text{Diff}(s)(\mathbf{X}, \mathbf{X}')) \right]\end{aligned}$$

The bipartite risk revisited

We can rewrite the bipartite risk as

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^D(s) &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\ell_1^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right] \\ &= \frac{1}{2} \mathbb{E}_{(\mathbf{X}, \mathbf{X}') \sim (P \times Q)} \left[\ell_1^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right] + \\ &\quad \frac{1}{2} \mathbb{E}_{(\mathbf{X}, \mathbf{X}') \sim (Q \times P)} \left[\ell_{-1}^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right]\end{aligned}$$

The bipartite risk revisited

We can rewrite the bipartite risk as

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^D(s) &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')] + \frac{1}{2} \mathbb{I}[s(\mathbf{X}) = s(\mathbf{X}')] \right] \\ &= \mathbb{E}_{\mathbf{X}\sim P, \mathbf{X}'\sim Q} \left[\ell_1^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right] \\ &= \frac{1}{2} \mathbb{E}_{(\mathbf{X}, \mathbf{X}')\sim (P \times Q)} \left[\ell_1^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right] + \\ &\quad \frac{1}{2} \mathbb{E}_{(\mathbf{X}, \mathbf{X}')\sim (Q \times P)} \left[\ell_{-1}^{01}((\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right] \\ &= \mathbb{E}_{((\mathbf{X}, \mathbf{X}'), \mathbf{Z})\sim \text{Bipart}(D)} \left[\ell^{01}(\mathbf{Z}, (\text{Diff}(s))(\mathbf{X}, \mathbf{X}')) \right]\end{aligned}$$

Reduction to classification

This generalises for any surrogate loss ℓ :

$$\mathbb{L}_{\text{Bipart},\ell}^D(s) = \mathbb{L}_{\text{Class},\ell}^{\text{Bipart}(D)}(\text{Diff}(s))$$

Equivalence: Bipartite ranking = binary classification over pairs

- Can transport all results over to bipartite setting!

Bayes-optimal surrogate bipartite scorers

Surrogate bipartite **Bayes risk**:

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

Bayes-optimal surrogate bipartite scorers

Surrogate bipartite **Bayes risk**:

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^{D,*} &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\text{Bipart},\ell}^D(s) \\ &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(\text{Diff}(s))\end{aligned}$$

Bayes-optimal surrogate bipartite scorers

Surrogate bipartite **Bayes risk**:

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^{D,*} &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\text{Bipart},\ell}^D(s) \\ &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(\text{Diff}(s)) \\ &= \min_{s_{\text{Pair}} \in \mathcal{S}_{\text{Decomp}}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(s_{\text{Pair}})\end{aligned}$$

where

$$\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) \mid s: \mathcal{X} \rightarrow \mathbb{R}\}$$

Bayes-optimal surrogate bipartite scorers

Surrogate bipartite **Bayes risk**:

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^{D,*} &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\text{Bipart},\ell}^D(s) \\ &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(\text{Diff}(s)) \\ &= \min_{s_{\text{Pair}} \in \mathcal{S}_{\text{Decomp}}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(s_{\text{Pair}}) \\ &\neq \min_{s_{\text{Pair}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(s_{\text{Pair}})\end{aligned}$$

where

$$\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) \mid s: \mathcal{X} \rightarrow \mathbb{R}\}$$

Bayes-optimal surrogate bipartite scorers

Surrogate bipartite **Bayes risk**:

$$\begin{aligned}\mathbb{L}_{\text{Bipart},\ell}^{D,*} &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\text{Bipart},\ell}^D(s) \\ &= \min_{s: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(\text{Diff}(s)) \\ &= \min_{s_{\text{Pair}} \in \mathcal{S}_{\text{Decomp}}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(s_{\text{Pair}}) \\ &\neq \min_{s_{\text{Pair}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}} \mathbb{L}_{\ell}^{\text{Bipart}(D)}(s_{\text{Pair}}) \\ &= \mathbb{L}_{\ell}^{\text{Bipart}(D),*}\end{aligned}$$

where

$$\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) \mid s: \mathcal{X} \rightarrow \mathbb{R}\}$$

An inconvenient truth

Catch: $\text{Bipart}(D)$ operates on **decomposable** pair-scorers:

$$\mathcal{S}_{\text{Decomp}} = \{\text{Diff}(s) \mid s: \mathcal{X} \rightarrow \mathbb{R}\}$$

- Effectively a restricted function class

In general,

$$\begin{aligned} \text{Diff}\left(\mathcal{S}_{\text{Bipart},\ell}^{D,*}\right) &\neq \mathcal{S}_{\ell}^{\text{Bipart}(D),*} \\ \text{regret}_{\text{Bipart},\ell}^D(s) &\neq \text{regret}_{\ell}^{\text{Bipart}(D)}(\text{Diff}(s)) \end{aligned}$$

Outline

- 1 The classification risk
- 2 The bipartite risk
- 3 Decomposability and risk minimisers**
- 4 Risk equivalences and algorithmic implications
- 5 Conclusion

Decomposable solutions

Suppose ℓ is such that

$$\mathcal{S}_{\ell}^{\text{Bipart}(D),*} \subseteq \mathcal{S}_{\text{Decomp}}$$

i.e. the **optimal pair-scorer** is **decomposable**

For such losses,

$$\begin{aligned}\mathcal{S}_{\text{Bipart},\ell}^{D,*} &= \mathcal{S}_{\ell}^{\text{Bipart}(D),*} \\ \text{regret}_{\text{Bipart},\ell}^D(s) &= \text{regret}_{\ell}^{\text{Bipart}(D)}(s)\end{aligned}$$

Decomposable solutions

Suppose ℓ is such that

$$\mathcal{S}_{\ell}^{\text{Bipart}(D),*} \subseteq \mathcal{S}_{\text{Decomp}}$$

i.e. the **optimal pair-scorer** is **decomposable**

For such losses,

$$\begin{aligned}\mathcal{S}_{\text{Bipart},\ell}^{D,*} &= \mathcal{S}_{\ell}^{\text{Bipart}(D),*} \\ \text{regret}_{\text{Bipart},\ell}^D(s) &= \text{regret}_{\ell}^{\text{Bipart}(D)}(s)\end{aligned}$$

Which ℓ induce this?

Characterising decomposability

For **strictly proper composite** ℓ with link function Ψ ,

$$\mathcal{S}_{\ell}^{\text{Bipart}(D),*} = \Psi \circ \eta_{\text{Pair}}$$

- $\eta_{\text{Pair}} : (x, x') \mapsto \Pr[Z = 1 | X = x, X' = x']$
- Observation-conditional density for $\text{Bipart}(D)$

$\mathcal{S}_{\ell}^{\text{Bipart}(D),*}$ decomposable \rightarrow interplay of Ψ and η_{Pair}

Characterising decomposability

For **strictly proper composite** ℓ with link function Ψ ,

$$\mathcal{S}_{\ell}^{\text{Bipart}(D),*} = \Psi \circ \eta_{\text{Pair}}$$

- $\eta_{\text{Pair}} : (x, x') \mapsto \Pr[Z = 1 | X = x, X' = x']$
- Observation-conditional density for $\text{Bipart}(D)$

$\mathcal{S}_{\ell}^{\text{Bipart}(D),*}$ decomposable \rightarrow interplay of Ψ and η_{Pair}

What is η_{Pair} ?

An innocuous lemma

Lemma

For $\sigma(\cdot)$ being the sigmoid function,

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

An innocuous lemma

Lemma

For $\sigma(\cdot)$ being the sigmoid function,

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

Peculiar re-expression of Bayes' rule:

$$\begin{aligned}\eta_{\text{Pair}}(x, x') &= \frac{\Pr[\mathbf{X} = x, \mathbf{X}' = x' | \mathbf{Z} = 1] \cdot \Pr[\mathbf{Z} = 1]}{\Pr[\mathbf{X} = x, \mathbf{X}' = x']} \\ &= \frac{1}{1 + \frac{Q(x)}{P(x)} \cdot \frac{P(x')}{Q(x')}} \\ &= \frac{1}{1 + \frac{1 - \eta(x)}{\eta(x)} \cdot \frac{\eta(x')}{1 - \eta(x')}}.\end{aligned}$$

Back to decomposability

For strictly proper composite ℓ ,

$$\begin{aligned}\mathcal{S}_\ell^{\text{Bipart}(D),*} &= \Psi \circ \eta_{\text{Pair}} \\ &= \Psi \circ \sigma \circ \text{Diff}(\sigma^{-1} \circ \eta)\end{aligned}$$

i.e. monotone transform of **decomposable** pair-scorer

$\mathcal{S}_\ell^{\text{Bipart}(D),*}$ decomposable \rightarrow **Ψ “cancelling” σ**

Characterising decomposability

Let

$$\Sigma = \{f : v \mapsto \sigma(av) \mid a \in \mathbb{R} - \{0\}\}$$

Lemma

Given any strictly proper composite loss ℓ with a differentiable, invertible link function Ψ ,

$$(\forall D) \mathcal{S}_{\ell}^{\text{Bipart}(D),*} \subseteq \mathcal{S}_{\text{Decomp}} \iff \Psi^{-1} \in \Sigma.$$

Inverse link must be **scaled sigmoid**

- Holds for logistic, exponential loss

Bayes-optimal scorers

Proposition

Given any strictly proper composite loss ℓ with a differentiable, invertible link function Ψ ,

$$\Psi^{-1} \in \Sigma \implies \mathcal{S}_{\text{Bipart},\ell}^{D,*} = \{\Psi \circ \eta + b : b \in \mathbb{R}\} \subseteq \mathcal{S}_{\text{Bipart},01}^{D,*}.$$

Follows because

$$\begin{aligned}\mathcal{S}_{\text{Bipart},\ell}^{D,*} &= \Psi \circ \sigma \circ \eta_{\text{Pair}} \\ &= \frac{1}{a} \sigma^{-1} \circ \sigma \circ \text{Diff}(\sigma^{-1} \circ \eta) \\ &= \text{Diff}(\Psi \circ \eta)\end{aligned}$$

Surrogate regret bound

Surrogate regret bound also follows immediately

Proposition

Given any strictly proper composite loss ℓ with a differentiable, invertible link function Ψ ,

$$\Psi^{-1} \in \Sigma \implies (\exists F_\ell) (\forall D, s) F_\ell \left(\text{regret}_{\text{Bipart}, 01}^D(s) \right) \leq \text{regret}_{\text{Bipart}, \ell}^D(s).$$

F_ℓ identical to that in surrogate bounds for classification

- Implies **Bayes-consistency** of suitable pairwise surrogate minimisation

Comments

Decomposability is **sufficient** for consistency

- Non-decomposable loss can be infinite-sample consistent
- Hinge-loss \rightarrow **inconsistent**

What is special about the link functions in Σ ?

- Boils down to form of η_{Pair}
- **Strict utility representation** for probabilistic binary relations

Outline

- 1 The classification risk
- 2 The bipartite risk
- 3 Decomposability and risk minimisers
- 4 Risk equivalences and algorithmic implications**
- 5 Conclusion

Theoretical equivalences of risks

For proper composite ℓ with inverse link in Σ ,

$$\text{Diff}(\mathcal{S}_{\text{Class},\ell}^{D,*}) = \text{Diff}(\mathcal{S}_{\text{Bipart},\ell}^{D,*}) = \mathcal{S}_{\ell}^{\text{Bipart}(D),*}$$

Disparate risks have **identical minimisers**:

$$\begin{aligned} & \underset{s_{\text{Pair}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{(\mathbf{X}, \mathbf{X}') \sim P \times Q} \left[e^{-s_{\text{Pair}}(\mathbf{X}, \mathbf{X}')} \right] \\ &= \text{Diff} \left(\underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{(\mathbf{X}, \mathbf{X}') \sim P \times Q} \left[e^{-(s(\mathbf{X}) - s(\mathbf{X}'))} \right] \right) \\ &= \text{Diff} \left(\underset{s: \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim D} \left[e^{-\mathbf{Y}s(\mathbf{X})} \right] \right) \end{aligned}$$

Practical equivalences of risks

To be “practically equivalent”, for $\mathcal{F} \subset \{s: \mathcal{X} \rightarrow \mathbb{R}\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell}^D(s) \stackrel{?}{=} \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Bipart}, \ell}^D(s)$$

- Failing which, surrogate regret bounds

Practical equivalences of risks

To be “practically equivalent”, for $\mathcal{F} \subset \{s: \mathcal{X} \rightarrow \mathbb{R}\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell}^D(s) \stackrel{?}{=} \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Bipart}, \ell}^D(s)$$

- Failing which, surrogate regret bounds

Remarkably, for $\ell^{\text{exp}}(\mathbf{y}, \mathbf{v}) = e^{-\mathbf{y}\mathbf{v}}$ and $\mathcal{F} = \{x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^N\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell^{\text{exp}}}^D(s) = \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Bipart}, \ell^{\text{exp}}}^D(s)$$

Practical equivalences of risks

To be “practically equivalent”, for $\mathcal{F} \subset \{s: \mathcal{X} \rightarrow \mathbb{R}\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell}^D(s) \stackrel{?}{=} \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Bipart}, \ell}^D(s)$$

- Failing which, surrogate regret bounds

Remarkably, for $\ell^{\text{exp}}(y, v) = e^{-yv}$ and $\mathcal{F} = \{x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^N\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell^{\text{exp}}}^D(s) = \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Bipart}, \ell^{\text{exp}}}^D(s)$$

Pointwise versus pairwise bipartite ranking

- Other “practical equivalences”?

The p -norm push risk

For $p \in [1, \infty)$, the p -norm push risk (Rudin, 2009) is

$$\mathbb{L}_{\text{push}}^D(s) = \mathbb{E}_{\mathbf{X}' \sim Q} \left[\left(\mathbb{E}_{\mathbf{X} \sim P} [\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')]] \right)^p \right]$$

- $p = 1 \rightarrow$ standard bipartite risk
- $p > 1 \rightarrow$ penalises high false negative rates
 - ▶ Suitable for “ranking the best”

The p -norm push risk

For $p \in [1, \infty)$, the p -norm push risk (Rudin, 2009) is

$$\mathbb{L}_{\text{push}}^D(s) = \mathbb{E}_{\mathbf{X}' \sim Q} \left[\left(\mathbb{E}_{\mathbf{X} \sim P} [\mathbb{I}[s(\mathbf{X}) < s(\mathbf{X}')]] \right)^p \right]$$

- $p = 1 \rightarrow$ standard bipartite risk
- $p > 1 \rightarrow$ penalises **high false negative** rates
 - ▶ Suitable for “ranking the best”

For $p \in [1, \infty)$, and **surrogate loss** $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$, define

$$\mathbb{L}_{\text{push}, \ell}^D(s) = \mathbb{E}_{\mathbf{X}' \sim Q} \left[\left(\mathbb{E}_{\mathbf{X} \sim P} [\ell_1(s(\mathbf{X}) - s(\mathbf{X}'))] \right)^p \right]$$

Bayes-optimal scorers for p -norm push

Can show (less easily than before!):

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

Bayes-optimal scorers for p -norm push

Can show (less easily than before!):

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

If ℓ is strictly proper composite with $\Psi = \frac{1}{p+1} \sigma^{-1}$,

$$\text{Diff}(\mathcal{S}_{\text{Class},\ell}^{D,*}) = \text{Diff}(\mathcal{S}_{\text{Bipart},\ell}^{D,*}) = \text{Diff}(\mathcal{S}_{\text{push,exp}}^{D,*})$$

- $\ell(y, v) = \frac{1}{p+1} \log(1 + e^{-y(p+1)v})$

Bayes-optimal scorers for p -norm push

Can show (less easily than before!):

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

If ℓ is strictly proper composite with $\Psi = \frac{1}{p+1} \sigma^{-1}$,

$$\text{Diff}(\mathcal{S}_{\text{Class},\ell}^{D,*}) = \text{Diff}(\mathcal{S}_{\text{Bipart},\ell}^{D,*}) = \text{Diff}(\mathcal{S}_{\text{push,exp}}^{D,*})$$

- $\ell(y, v) = \frac{1}{p+1} \log(1 + e^{-y(p+1)v})$

Restricted function class?

The p -classification loss

For ℓ being the p -classification loss,

$$\ell^{\text{pc}}(y, v) = \llbracket y = 1 \rrbracket e^{-v} + \llbracket y = -1 \rrbracket \frac{1}{p} e^{vp}$$

and for $\mathcal{F} = \{x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^N\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell^{\text{pc}}}^D(s) = \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{push}, \ell^{\text{pc}}}^D(s)$$

The p -classification loss

For ℓ being the p -classification loss,

$$\ell^{\text{pc}}(y, v) = \llbracket y = 1 \rrbracket e^{-v} + \llbracket y = -1 \rrbracket \frac{1}{p} e^{vp}$$

and for $\mathcal{F} = \{x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^N\}$,

$$\operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{Class}, \ell^{\text{pc}}}^D(s) = \operatorname{argmin}_{s \in \mathcal{F}} \mathbb{L}_{\text{push}, \ell^{\text{pc}}}^D(s)$$

How does this loss help in “ranking the best”?

Weight function for proper losses

Any proper composite loss is expressible as a **weighted combination** of **cost-sensitive losses**:

$$\ell(y, v) = \int_0^1 w(c) \cdot \ell^{\text{CS}(c)}(y, v) dc$$

where

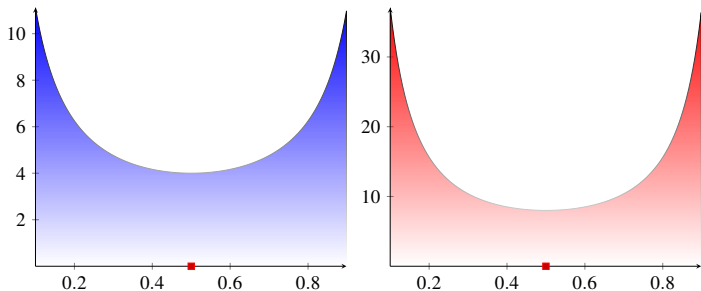
$w(c)$ = **weight function** over misclassification costs

$$\ell^{\text{CS}(c)}(y, v) = \llbracket y = 1 \wedge \Psi^{-1}(v) < 0 \rrbracket \cdot (1 - c) + \llbracket y = -1 \wedge \Psi^{-1}(v) > 0 \rrbracket \cdot c$$

Weight functions: Logistic and Exponential Loss

$$\log(1 + e^{-yv}) = \int_0^1 \frac{1}{c(1-c)} \cdot \ell^{\text{CS}(c)}(y, \sigma(v)) dc$$

$$e^{-yv} = \int_0^1 \frac{1}{c^{3/2}(1-c^{3/2})} \cdot \ell^{\text{CS}(c)}(y, \sigma(2v)) dc$$

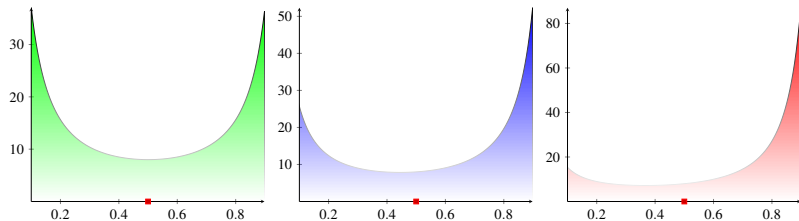


Weight function for p -classification

ℓ^{pc} has asymmetric weight function

$$w(c) = \frac{1}{c^{1+\frac{1}{p+1}}(1-c)^{2-\frac{1}{p+1}}}$$

Increase $c \rightarrow$ focus on high cost ratios



Alternate asymmetric losses

Can consider other losses with asymmetric weights, e.g.

$$w(c) = \frac{1}{c(1-c)^{3/2}}$$

corresponding to

$$\ell(v) = \left(\frac{1}{\sqrt{\sigma(-v)}}, \tanh^{-1}(\sqrt{\sigma(-v)}) \right)$$

Alternate asymmetric losses

Can consider other losses with asymmetric weights, e.g.

$$w(c) = \mathbb{I}[2c < 1] \frac{1}{c(1-c)} + \mathbb{I}[2c > 1] \frac{1}{2c^{3/2}(1-c)^{3/2}}$$

corresponding to

$$\ell(v) = \left(\begin{array}{cc} \left\{ \begin{array}{ll} \log(1 + e^v) & \text{if } v < 0 \\ e^{v/2} + \log 2 - 1 & \text{if } v \geq 0 \end{array} \right. & \left\{ \begin{array}{ll} \lg(1 + e^{-v}) + 1 & \text{if } v < 0 \\ e^{-v/2} & \text{if } v \geq 0 \end{array} \right. \end{array} \right)$$

Empirical Performance

Loss	ℓ^{\log} , ℓ^{\exp} , ℓ^{pc} , hybrid
Risk	Classification, bipartite, p -norm
Datasets	ionosphere, housing, german, car
Performance	AUC, MRR, DCG, AP, PTop
Caveat	Assessing viability of our “recipe” (Not that we “rank the best” “the best”)

Empirical Performance

Method	AUC	MRR	DCG	AP	P _{Top}
Proper Logistic	6.0000	7.7500	8.0000	7.7500	3.2500
Proper Exponential	5.2500	5.5000	5.7500	7.2500	4.5000
Proper P-Classification	7.0000	8.7500	8.5000	7.7500	4.5000
Proper Asymmetric A	5.2500	7.5000	7.5000	5.0000	1.5000
Proper Asymmetric B	4.7500	7.7500	7.5000	9.0000	6.2500
Bipartite Logistic	4.5000	7.0000	7.7500	6.2500	2.5000
Bipartite Exponential	6.7500	5.5000	6.2500	8.2500	4.0000
Bipartite P-Classification	5.2500	7.2500	7.5000	5.7500	3.0000
Bipartite Asymmetric A	3.0000	7.0000	6.7500	3.7500	2.5000
Bipartite Asymmetric B	8.0000	7.7500	9.0000	7.0000	3.2500
P-Norm Logistic	7.5000	9.0000	10.0000	7.0000	2.2500
P-Norm Exponential	6.7500	7.5000	7.2500	8.7500	4.7500
P-Norm Asymmetric A	7.0000	7.2500	7.7500	9.2500	3.7500
P-Norm Asymmetric B	3.2500	5.7500	5.5000	7.2500	5.2500

Comments

Bayes-optimal scorers for “ranking the best”

- Non-strict proper losses
- Cannot be made convex!

Why exponential loss?

- Bregman divergence perspective

P_{Top} regret bounds?

Outline

- 1 The classification risk
- 2 The bipartite risk
- 3 Decomposability and risk minimisers
- 4 Risk equivalences and algorithmic implications
- 5 Conclusion**

Take-home messages

Bipartite ranking = classification problem over pairs

Decomposability \rightarrow Bayes-optimal scorers, risk equivalences

Some risk equivalences hold for restricted function classes

- Algorithmic implications for bipartite ranking and its generalisations

Thanks!