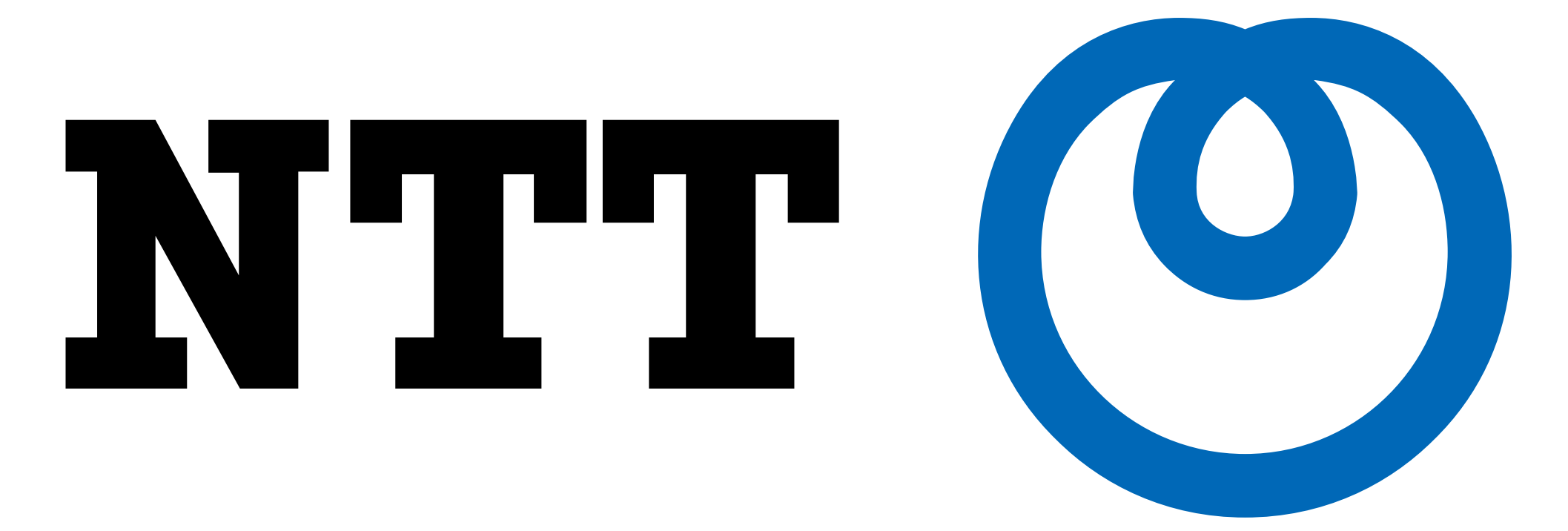


Structured Prediction with Projection Oracles



Mathieu Blondel, NTT Communication Science laboratories, Kyoto, Japan

Structured prediction

Target loss

- ▶ Goal: Learn a mapping $f: \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ **Target loss** $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$
- ▶ Minimize target expected risk

$$\mathcal{L}(f) := \mathbb{E}_{(X,Y) \sim \rho} L(f(X), Y)$$

Surrogate loss

- ▶ Vector space $\Theta \subseteq \mathbb{R}^p$
- ▶ Label encoding (embedding) $\varphi: \mathcal{Y} \rightarrow \Theta$
- ▶ Model $g: \mathcal{X} \rightarrow \Theta$, $\theta = g(x)$
- ▶ **Surrogate loss** $S: \Theta \times \Theta \rightarrow \mathbb{R}$
- ▶ Minimize surrogate expected risk

$$\mathcal{S}(g) := \mathbb{E}_{(X,Y) \sim \rho} S(g(X), \varphi(Y))$$
- ▶ Pull-back: **decoder** $d: \Theta \rightarrow \mathcal{Y}$
- ▶ Example: Maximum A-Posteriori decoder

$$\text{MAP}(\theta) := \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \langle \theta, \varphi(y) \rangle$$

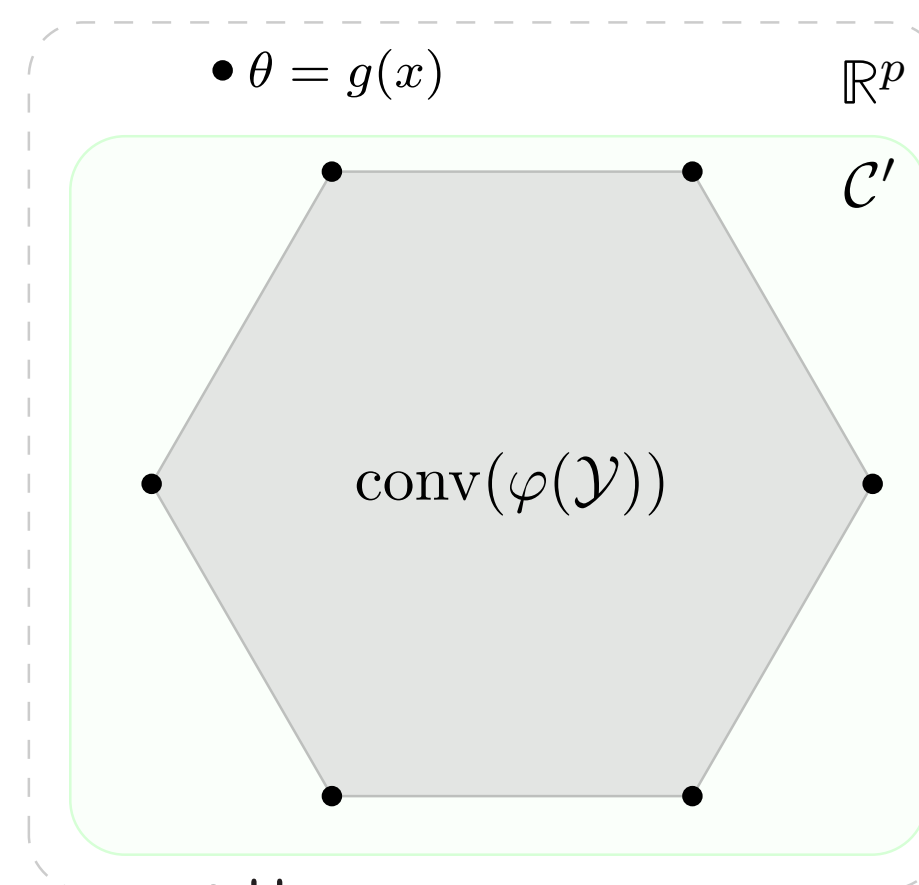
Why projections?

- ▶ Euclidean **projection** onto a closed convex set \mathcal{C}

$$P_{\mathcal{C}}(\theta) := \underset{u \in \mathcal{C}}{\operatorname{argmin}} \|\theta - u\|$$
- ▶ As long as $\varphi(y) \in \mathcal{C}$, we have

$$\|\varphi(y) - P_{\mathcal{C}}(\theta)\|^2 \leq \|\varphi(y) - \theta\|^2$$
 → projection achieves **smaller** squared loss!
- ▶ Projections as an **output layer**

$$x \in \mathcal{X} \xrightarrow[\text{model}]{g} \theta \in \Theta \xrightarrow[\text{projection}]{P_{\mathcal{C}}} u \in \mathcal{C}$$



Loss functions

Composite loss

- ▶ $\ell_{\mathcal{C}}(\theta, y) := \frac{1}{2} \|P_{\mathcal{C}}(\theta) - \varphi(y)\|^2$ **Non-convex!**

Fenchel-Young losses

- ▶ Conjugate $\Omega^*(\theta) := \max_{\mu \in \operatorname{dom}(\Omega)} \langle \mu, \theta \rangle - \Omega(\mu)$
- ▶ Fenchel-Young loss generated by Ω

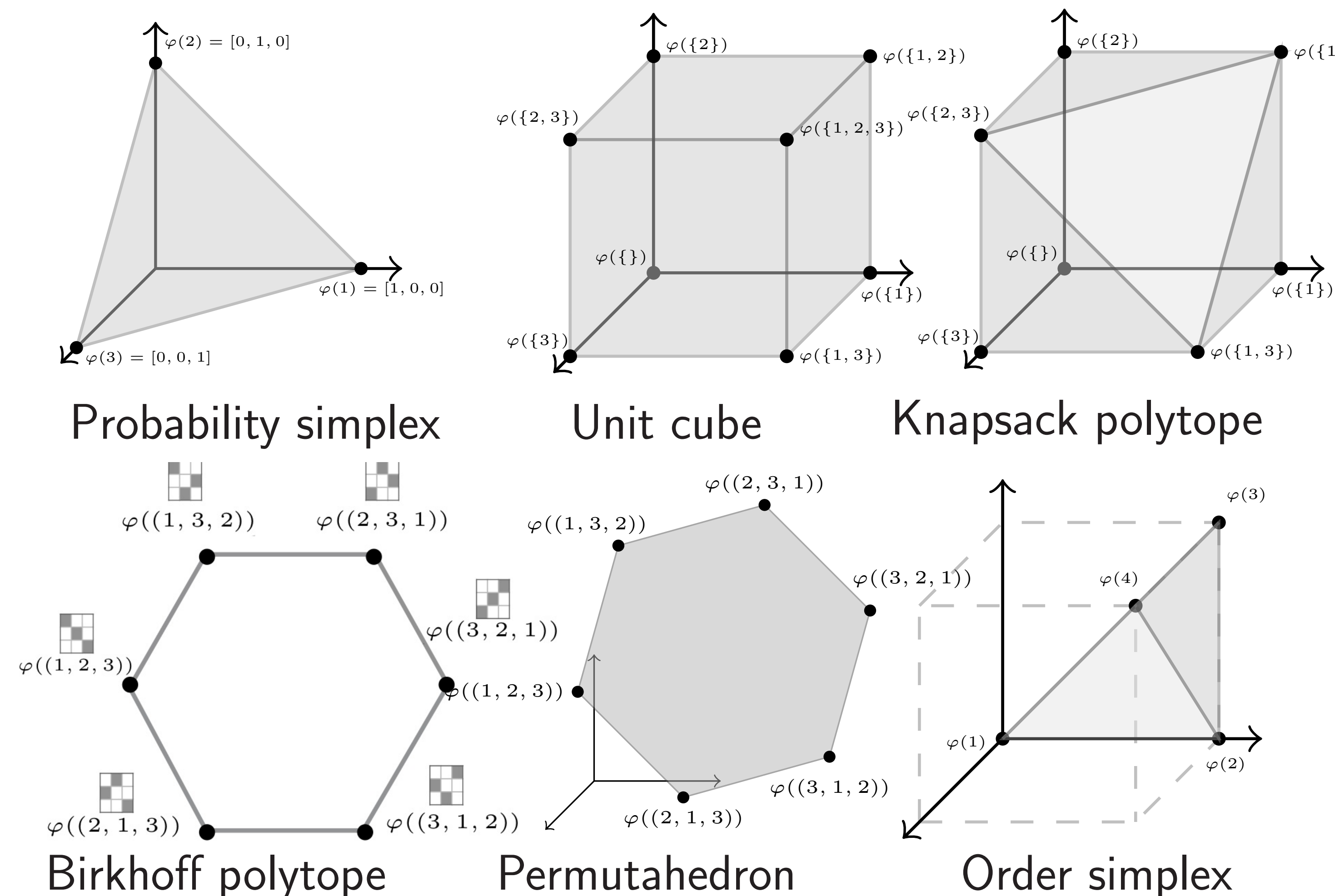
$$S_{\Omega}(\theta, y) := \Omega^*(\theta) + \Omega(\varphi(y)) - \langle \theta, \varphi(y) \rangle$$
- ▶ Convex, smooth if Ω strongly convex, non-negative

Projection loss

- ▶ Choose $\Omega = \frac{1}{2} \|\cdot\|^2 + I_{\mathcal{C}}$ and define shorthand

$$S_{\mathcal{C}}(\theta, y) := S_{\frac{1}{2}\|\cdot\|^2 + I_{\mathcal{C}}}(\theta, y)$$
- ▶ Zero loss: $S_{\mathcal{C}}(\theta, y) = 0 \Leftrightarrow P_{\mathcal{C}}(\theta) = \varphi(y)$
- ▶ **Convex upper-bound**: $\ell_{\mathcal{C}}(\theta, y) \leq S_{\mathcal{C}}(\theta, y)$
- ▶ **Smaller sets enjoy smaller loss**

$$S_{\mathcal{C}}(\theta, y) \leq S_{\mathcal{C}'}(\theta, y) \quad \forall \mathcal{C} \subseteq \mathcal{C}'$$
- ▶ Smallest convex set: the convex hull of $\varphi(\mathcal{Y})$



Calibration analysis

- ▶ Affine decomposition

$$L(\hat{y}, y) = \langle \varphi(\hat{y}), V\varphi(y) + b \rangle + c(y)$$
- ▶ Decoding **calibrated** for loss L

$$\hat{y}_L(u) := \underset{y' \in \mathcal{Y}}{\operatorname{argmin}} \langle \varphi(y'), Vu + b \rangle = \text{MAP}(-Vu - b)$$
- ▶ Analyzed prediction pipeline

$$x \in \mathcal{X} \xrightarrow[\text{model}]{g} \theta \in \Theta \xrightarrow[\text{projection}]{P_{\mathcal{C}}} u \in \mathcal{C} \xrightarrow[\text{decoding}]{\hat{y}_L} \hat{y} \in \mathcal{Y}$$
- ▶ Excess of risks

$$\delta \mathcal{L}(f) := \mathcal{L}(f) - \mathcal{L}(f^*), \delta \mathcal{S}_{\mathcal{C}}(g) := \mathcal{S}_{\mathcal{C}}(g) - \mathcal{S}_{\mathcal{C}}(g^*)$$
- ▶ **Theorem**

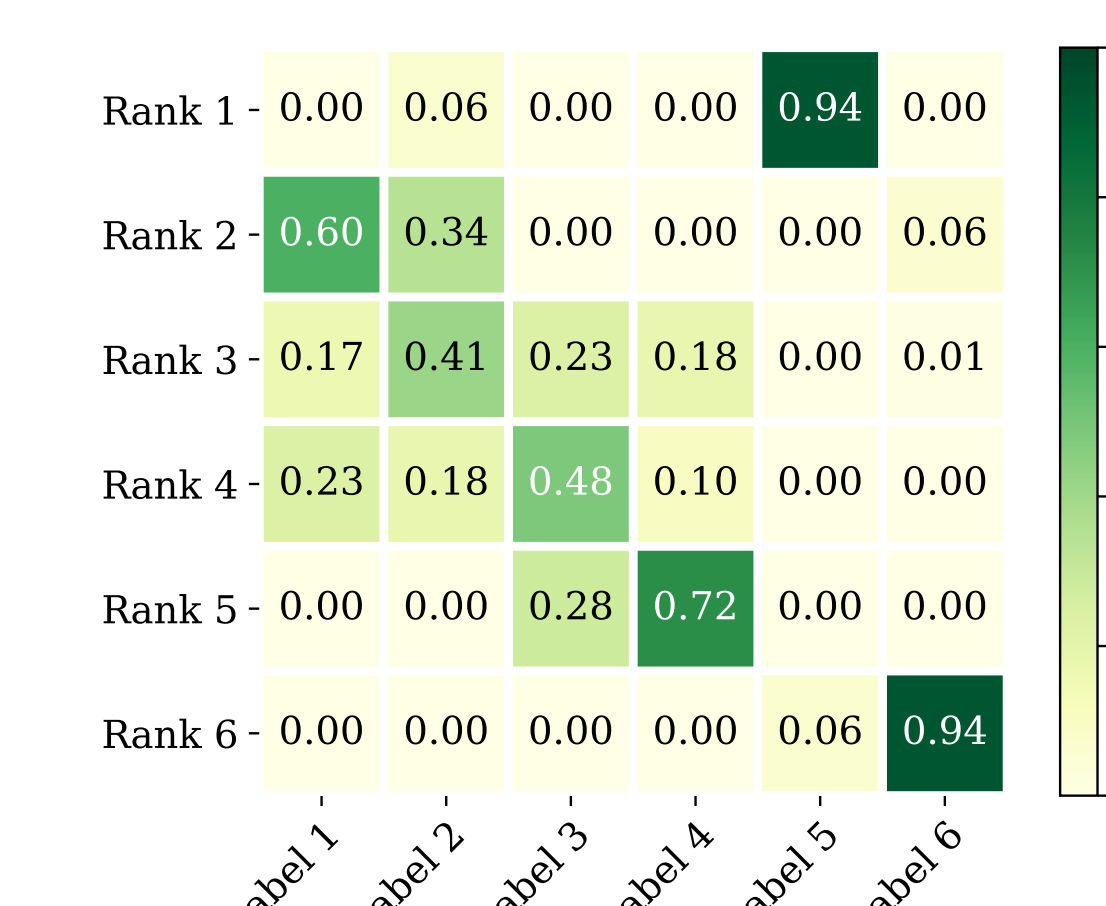
$$\forall g: \mathcal{X} \rightarrow \Theta: \frac{\delta \mathcal{L}(\hat{y}_L \circ P_{\mathcal{C}} \circ g)^2}{8\beta\sigma^2} \leq \delta \mathcal{S}_{\mathcal{C}}(g)$$

$1/\beta$: strong-convexity constant of Ω w.r.t. $\|\cdot\|$
 $\sigma := \sup_{\hat{y} \in \mathcal{Y}} \|V^{\top} \varphi(\hat{y})\|_*$

Experiment: label ranking

Projection	$[0, 1]^{k \times k}$	$\Delta^{k \times k}$	$\mathbb{R}^{k \times k}$	$[0, 1]^{k \times k}$	$\Delta^{k \times k}$	\mathcal{M}
Decoding	$[0, 1]^{k \times k}$	$\Delta^{k \times k}$	\mathcal{M}	\mathcal{M}	\mathcal{M}	\mathcal{M}
Authorship	12.83	5.62	5.70	5.18	5.70	5.10
Glass	24.35	5.43	7.11	5.68	5.04	4.65
Iris	27.78	10.37	19.26	4.44	1.48	2.96
Vehicle	26.36	7.43	9.04	7.57	6.99	5.88
Vowel	43.71	9.65	10.57	9.56	9.18	8.76
Wine	10.19	1.85	1.23	1.85	1.85	1.85

\mathcal{M} : Birkhoff polytope



Multi-label classification and ordinal regression experiments in the paper!

github.com/mbondel/projection-losses