# Structured Prediction with Projection Oracles

Mathieu Blondel

NTT, Kyoto, Japan

June 19th, 2019

# Outline

1. Background

2. Proposed framework

3. Experiments

# Outline

## 1. Background

## 2. Proposed framework

## 3. Experiments

# Structured prediction

Learn a mapping from input space to output space

$$f: \mathcal{X} \to \mathcal{Y}$$

# Structured prediction

Learn a mapping from input space to output space

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

exponentially large

# Structured prediction

Learn a mapping from input space to output space

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

exponentially large

Typically assume $f = dec \circ g$

$$x \in \mathcal{X} \xrightarrow[\text{model}]{g} \theta \in \Theta = \mathbb{R}^d \xrightarrow[\text{decoding}]{dec} \hat{y} \in \mathcal{Y}$$

# Structured prediction

Learn a mapping from input space to output space

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

exponentially large

Typically assume $f = dec \circ g$

$$x \in \mathcal{X} \xrightarrow[\text{model}]{g} \theta \in \Theta = \mathbb{R}^d \xrightarrow[\text{decoding}]{dec} \hat{y} \in \mathcal{Y}$$

$$\mathcal{L}(f) \triangleq \mathbb{E}_{(X,Y) \sim p} L(f(X), Y) \qquad L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$$

# Structured prediction

**Goal**

Learn a mapping from input space to output space
$$f: \mathscr{X} \to \mathscr{Y} \leftarrow \text{exponentially large}$$

**Decomposition**

Typically assume $f = dec \circ g$

$$x \in \mathscr{X} \xrightarrow[\text{model}]{g} \theta \in \Theta = \mathbb{R}^d \xrightarrow[\text{decoding}]{dec} \hat{y} \in \mathscr{Y}$$

**Target loss risk**

$$\mathscr{L}(f) \triangleq \mathbb{E}_{(X,Y)\sim p} L(f(X), Y) \qquad L: \mathscr{Y} \times \mathscr{y}$$

**Non-convex, discontinuous!**

# Surrogate losses

**Surrogate loss risk**

$$\mathcal{S}(g) \triangleq \mathbb{E}_{(X,Y)\sim p}\, S(g(X), Y) \qquad S: \Theta \times \mathcal{Y} \to \mathbb{R}_+$$

# Surrogate losses

$$\mathcal{S}(g) \triangleq \mathbb{E}_{(X,Y)\sim p}\, S(g(X), Y) \qquad S: \Theta \times \mathcal{Y} \to \mathbb{R}_+$$

**Fisher consistency**

$$\mathcal{S}(g_n) \to \inf_{g\in\mathcal{G}} \mathcal{S}(g) \quad\longrightarrow\quad \mathscr{L}(dec \circ g_n) \to \inf_{g\in\mathcal{G}} \mathscr{L}(dec \circ g)$$

# Surrogate losses

**Surrogate loss risk**

$$\mathcal{S}(g) \triangleq \mathbb{E}_{(X,Y) \sim p} \, S(g(X), Y) \qquad S : \Theta \times \mathscr{Y} \to \mathbb{R}_+$$

**Fisher consistency**

$$\mathcal{S}(g_n) \to \inf_{g \in \mathscr{G}} \mathcal{S}(g) \longrightarrow \mathscr{L}(dec \circ g_n) \to \inf_{g \in \mathscr{G}} \mathscr{L}(dec \circ g)$$

Extensively studied in the multiclass setting [Zhang 2004, Bartlett et al. 2006]

Only recently studied in the structured setting

[Ciliberto et al 2016,
Osokin et al. 2017,
Nowak-Vila et al. 2019]

# Structured perceptron loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathcal{Y}} \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle \qquad \varphi \colon \mathcal{Y} \to \mathbb{R}^d$$

# Structured perceptron loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathcal{Y}} \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle \qquad \varphi : \mathcal{Y} \to \mathbb{R}^d$$

**Training oracle**

$$MAP(\theta) \triangleq \arg\max_{y \in \mathcal{Y}} \langle \theta, \varphi(y) \rangle$$

# Structured perceptron loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathscr{Y}} \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle$$

$$\varphi : \mathscr{Y} \to \mathbb{R}^d$$

**Training oracle**

$$MAP(\theta) \triangleq \arg\max_{y \in \mathscr{Y}} \langle \theta, \varphi(y) \rangle$$

**Decoding oracle**

$$dec = MAP$$

# Structured perceptron loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathscr{Y}} \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle \qquad \varphi : \mathscr{Y} \to \mathbb{R}^d$$

**Training oracle**

$$MAP(\theta) \triangleq \arg\max_{y \in \mathscr{Y}} \langle \theta, \varphi(y) \rangle$$

**Decoding oracle**

$$dec = MAP$$

✕ **Not smooth**
✕ **Not consistent**

# Structured hinge loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathcal{Y}} \; L(y', y) + \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle$$

# Structured hinge loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathcal{Y}} \ L(y', y) + \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$\arg\max_{y' \in \mathcal{Y}} \ L(y, y') + \langle \theta, \varphi(y') \rangle$$

# Structured hinge loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathcal{Y}} \ L(y', y) + \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$\arg\max_{y' \in \mathcal{Y}} \ L(y, y') + \langle \theta, \varphi(y') \rangle$$

**Decoding oracle**

$$dec = MAP$$

# Structured hinge loss

**Loss**

$$S(\theta, y) \triangleq \max_{y' \in \mathscr{Y}} \; L(y', y) + \langle \theta, \varphi(y') \rangle - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$\arg\max_{y' \in \mathscr{Y}} \; L(y, y') + \langle \theta, \varphi(y') \rangle$$

**Decoding oracle**

$$dec = MAP$$

$\times$ **Not smooth**
$\times$ **Not consistent**

# Conditional Random Field (CRF) loss

[Lafferty et al., 2001]

**Loss**

$$S(\theta, y) \triangleq \log \sum_{y' \in \mathcal{Y}} e^{\langle \theta, \varphi(y') \rangle} - \langle \theta, \varphi(y) \rangle$$

# Conditional Random Field (CRF) loss

**Loss**

$$S(\theta, y) \triangleq \log \sum_{y' \in \mathcal{Y}} e^{\langle \theta, \varphi(y') \rangle} - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$marginal(\theta) \triangleq \mathbb{E}_{Y \sim p(\cdot; \theta)}[\varphi(Y)]$$

$$p(y; \theta) \propto \exp\langle \varphi(y), \theta \rangle$$

# Conditional Random Field (CRF) loss

[Lafferty et al., 2001]

**Loss**

$$S(\theta, y) \triangleq \log \sum_{y' \in \mathcal{Y}} e^{\langle \theta, \varphi(y') \rangle} - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$marginal(\theta) \triangleq \mathbb{E}_{Y \sim p(\cdot; \theta)}[\varphi(Y)]$$

$$p(y; \theta) \propto \exp\langle \varphi(y), \theta \rangle$$

**Decoding oracle**

$$MAP$$

Calibrated decoding

# Conditional Random Field (CRF) loss

[Lafferty et al., 2001]

**Loss**

$$S(\theta, y) \triangleq \log \sum_{y' \in \mathcal{Y}} e^{\langle \theta, \varphi(y') \rangle} - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$marginal(\theta) \triangleq \mathbb{E}_{Y \sim p(\cdot; \theta)}[\varphi(Y)]$$

$$p(y; \theta) \propto \exp\langle \varphi(y), \theta \rangle$$

**Decoding oracle**

$$MAP$$

Calibrated decoding

✓ **Smooth**
✓ **Consistent (w/ calibrated decoding) [Nowak-Vila et al., 2019]**

# Conditional Random Field (CRF) loss

[Lafferty et al., 2001]

**Loss**

$$S(\theta, y) \triangleq \log \sum_{y' \in \mathcal{Y}} e^{\langle \theta, \varphi(y') \rangle} - \langle \theta, \varphi(y) \rangle$$

**Training oracle**

$$marginal(\theta) \triangleq \mathbb{E}_{Y \sim p(\cdot;\theta)}[\varphi(Y)]$$

$$p(y;\theta) \propto \exp\langle \varphi(y), \theta \rangle$$

**Decoding oracle**

$$MAP$$

Calibrated decoding

✓ **Smooth**

✓ **Consistent (w/ calibrated decoding) [Nowak-Vila et al., 2019]**

✗ **Marginal inference is intractable for some tasks**

# Squared loss

**Loss**

$$S(\theta, y) \triangleq \frac{1}{2} \|\varphi(y) - \theta\|^2$$

# Squared loss

**Loss**

$$S(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - \theta\|^2$$

**Training oracle**

None!

# Squared loss

**Loss**

$$S(\theta, y) \triangleq \frac{1}{2} \|\varphi(y) - \theta\|^2$$

**Training oracle**

None!

**Decoding oracle**

Calibrated decoding

# Squared loss

**Loss**

$$S(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - \theta\|^2$$

**Training oracle**

None!

**Decoding oracle**

Calibrated decoding

✓ **Smooth**
✓ **Consistent (when using calibrated decoding)**

# Squared loss

**Loss**

$$S(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - \theta\|^2$$

**Training oracle**

None!

**Decoding oracle**

Calibrated decoding

✓ **Smooth**
✓ **Consistent (when using calibrated decoding)**
✗ **Ignores structural information at training time**

# Summary

| Loss | Training oracle | Decoding | Smooth | Consistent |
|---|---|---|---|---|
| Perceptron | MAP | MAP | No | No |
| SVM | Loss-augmented MAP | MAP | No | No |
| CRF | Marginal | MAP<br>Calibrated decoding | Yes | No<br>Yes |
| Squared | None | Calibrated decoding | Yes | Yes |

# Summary

| Loss | Training oracle | Decoding | Smooth | Consistent |
|---|---|---|---|---|
| Perceptron | MAP | MAP | No | No |
| SVM | Loss-augmented MAP | MAP | No | No |
| CRF | Marginal | MAP<br>Calibrated decoding | Yes | No<br>Yes |
| Squared | None | Calibrated decoding | Yes | Yes |
| **Proposed** | **Projection** | **Calibrated decoding** | **Yes** | **Yes** |

# Outline

1. Background

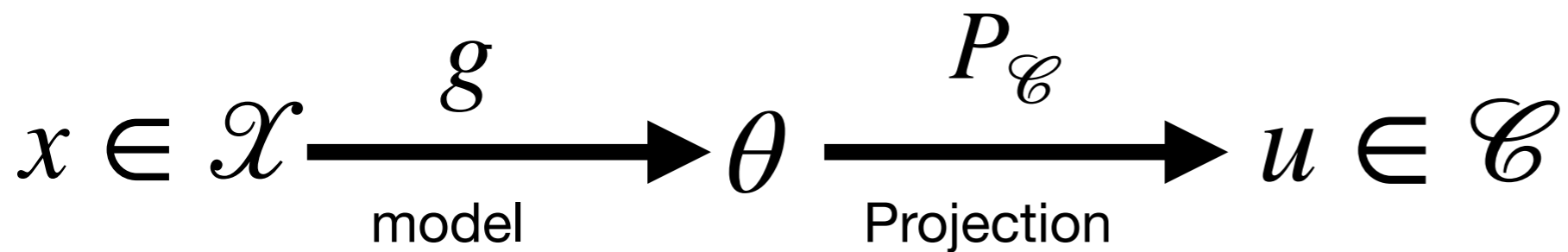2. Proposed framework

3. Experiments

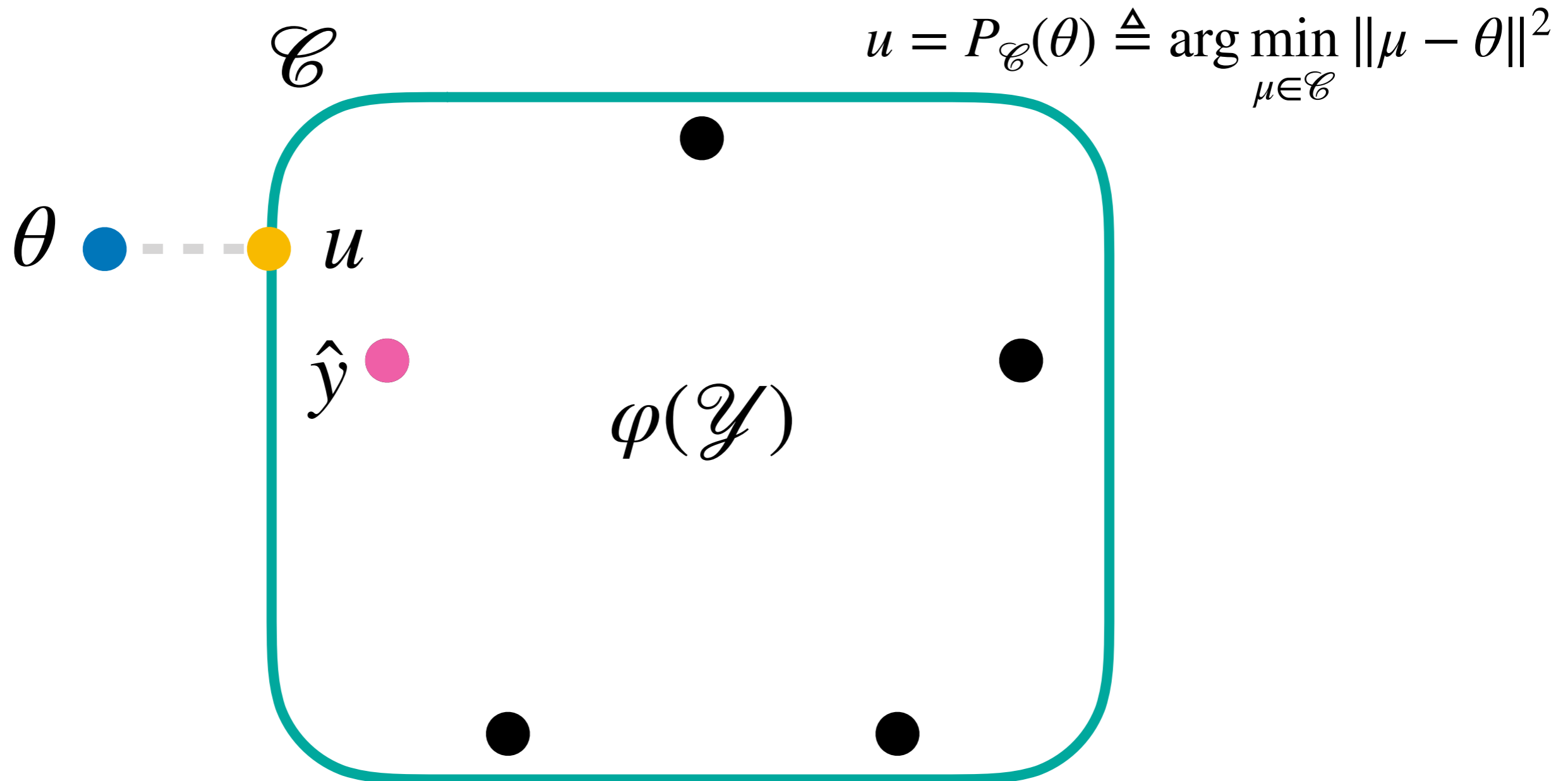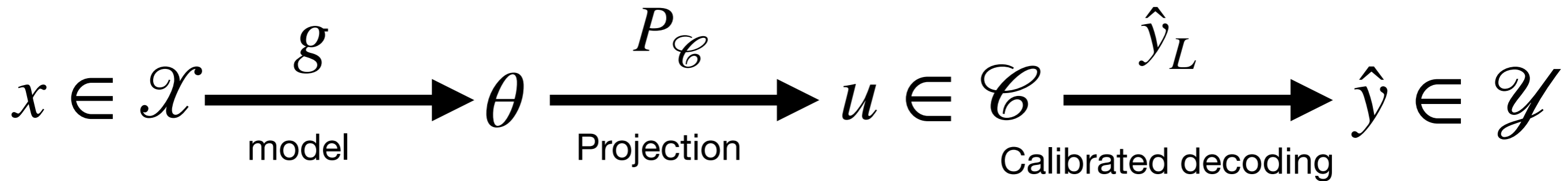$$x \in \mathscr{X} \xrightarrow[\text{model}]{g} \theta$$

$\theta$ •

$\varphi(\mathscr{Y})$

# Proposed inference pipeline

$$x \in \mathcal{X} \xrightarrow{\underset{\text{model}}{g}} \theta \xrightarrow{\underset{\text{Projection}}{P_{\mathscr{C}}}} u \in \mathscr{C}$$

$$u = P_{\mathscr{C}}(\theta) \triangleq \arg\min_{\mu \in \mathscr{C}} \|\mu - \theta\|^2$$

# Proposed inference pipeline

$$x \in \mathscr{X} \xrightarrow[\text{model}]{g} \theta \xrightarrow[\text{Projection}]{P_{\mathscr{C}}} u \in \mathscr{C} \xrightarrow[\text{Calibrated decoding}]{\hat{y}_L} \hat{y} \in \mathscr{Y}$$

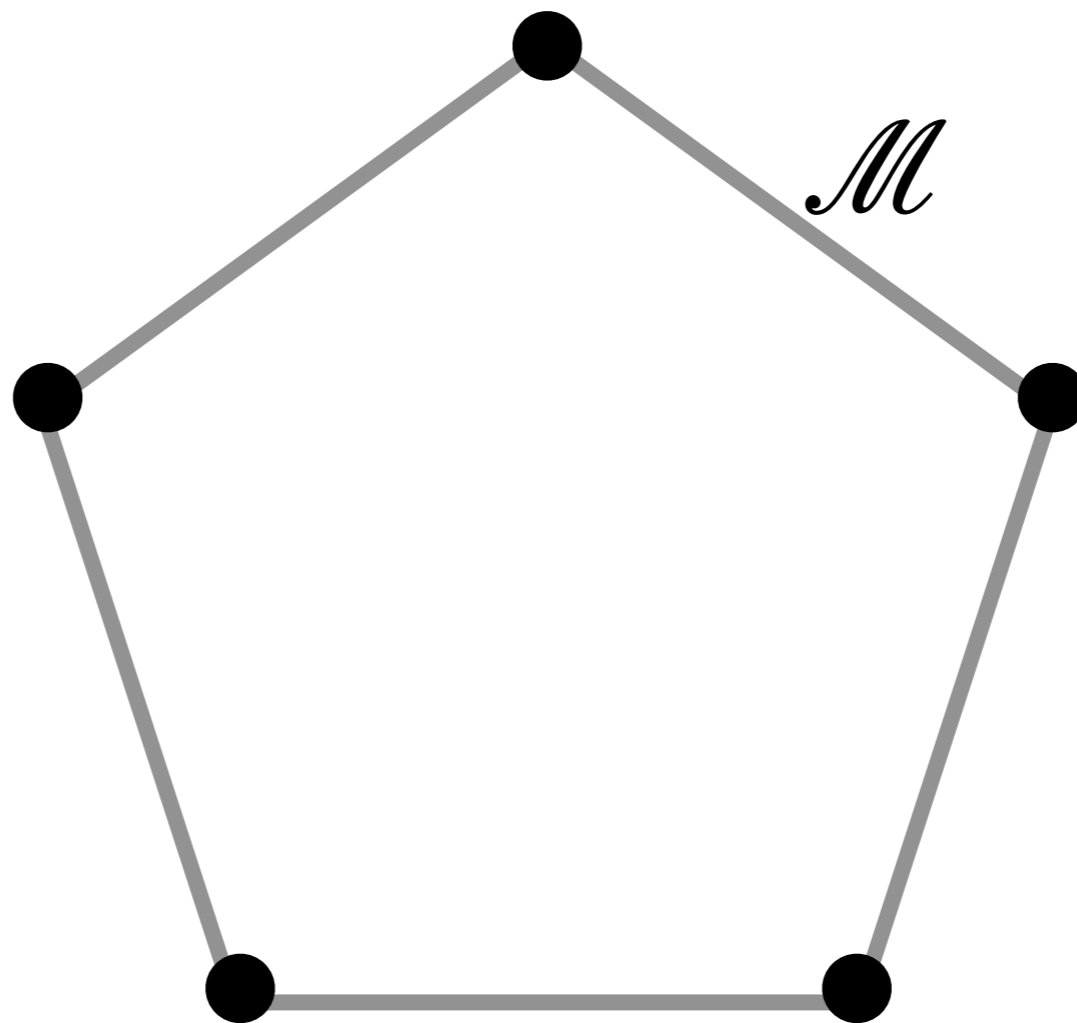$$u = P_{\mathscr{C}}(\theta) \triangleq \arg\min_{\mu \in \mathscr{C}} \|\mu - \theta\|^2$$

# Choice of the convex set
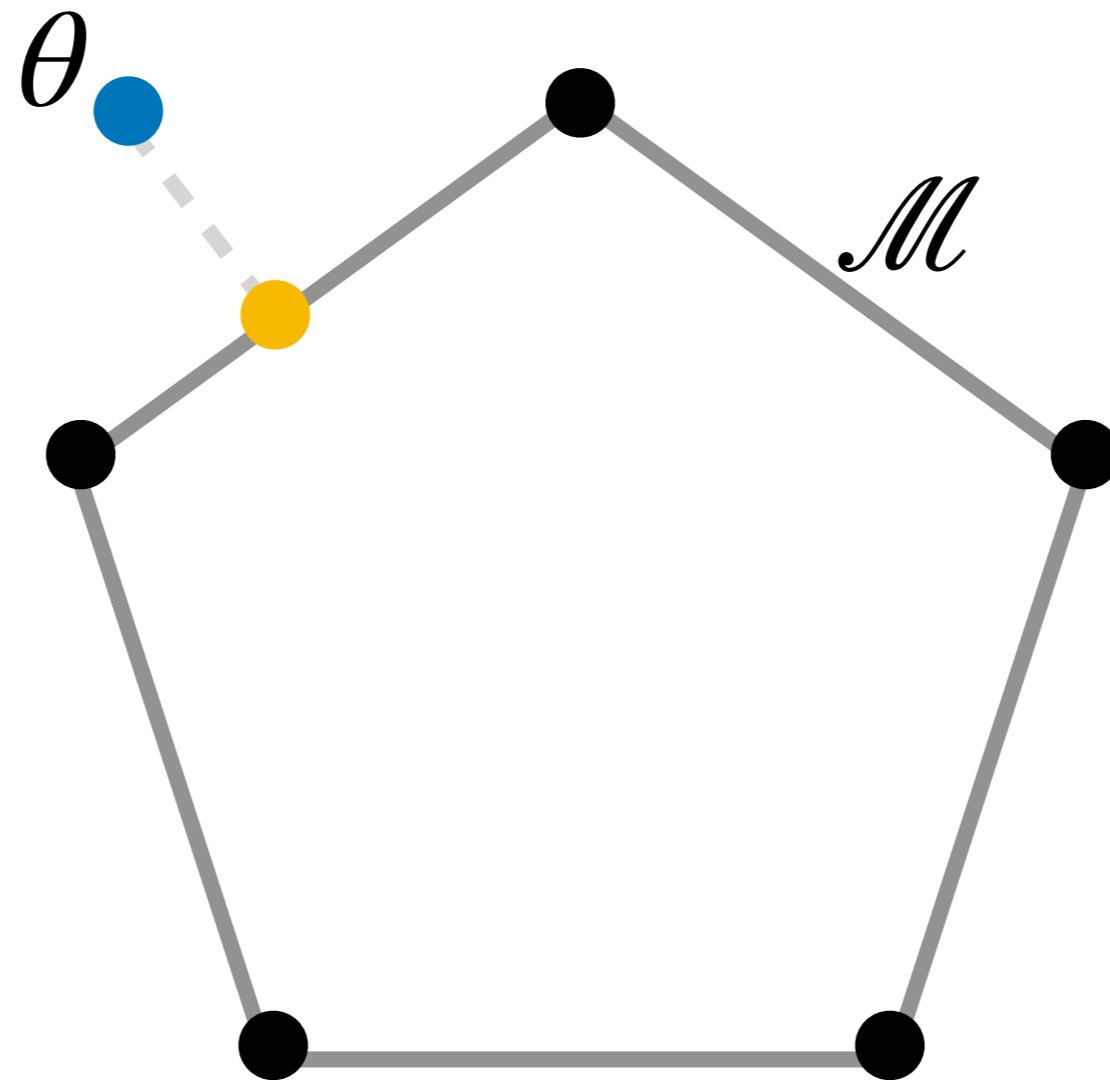
Smallest convex set = **convex hull** (a.k.a. marginal polytope)



$$\mathcal{M} \triangleq conv(\varphi(\mathcal{Y}))$$

# Choice of the convex set

Smallest convex set = **convex hull** (a.k.a. marginal polytope)



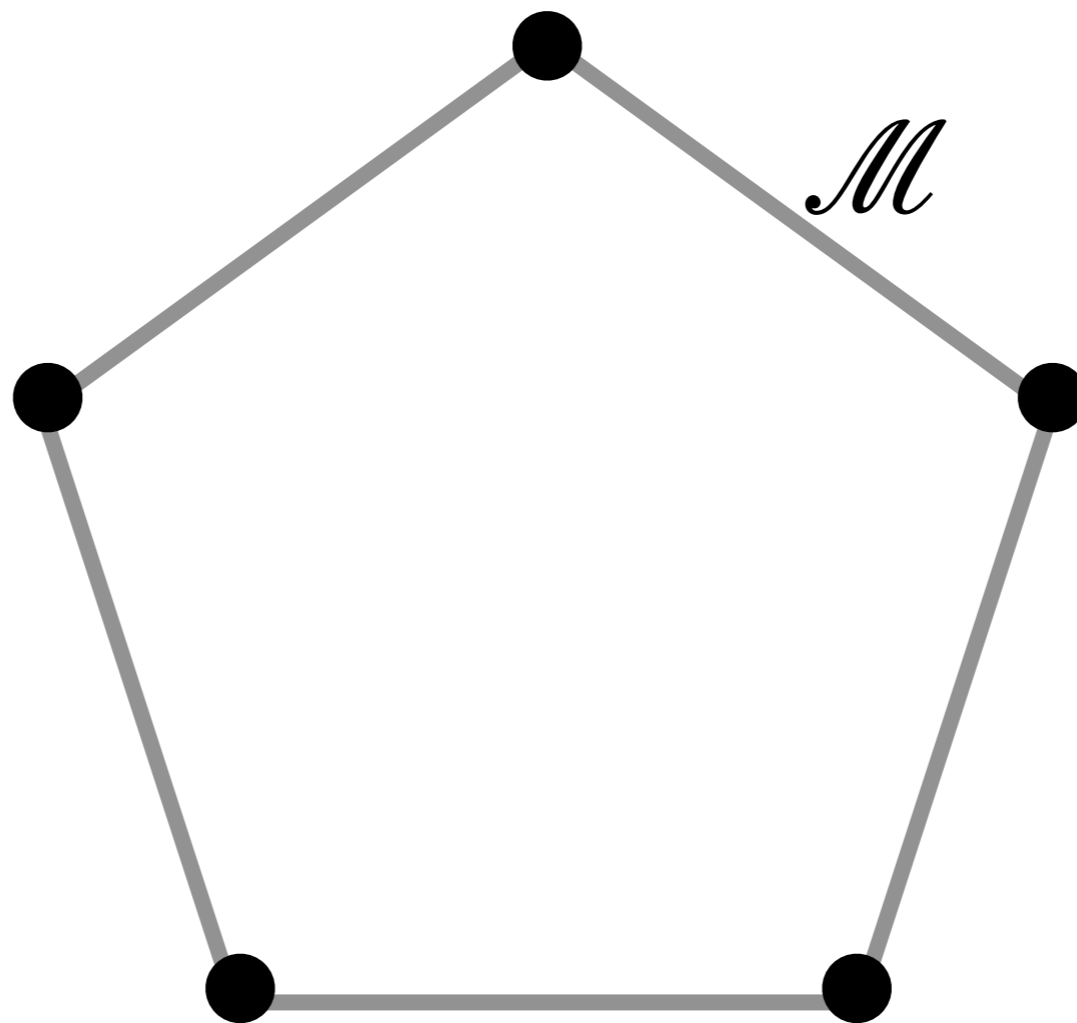$$\mathcal{M} \triangleq conv(\varphi(\mathcal{Y}))$$

# Choice of the convex set

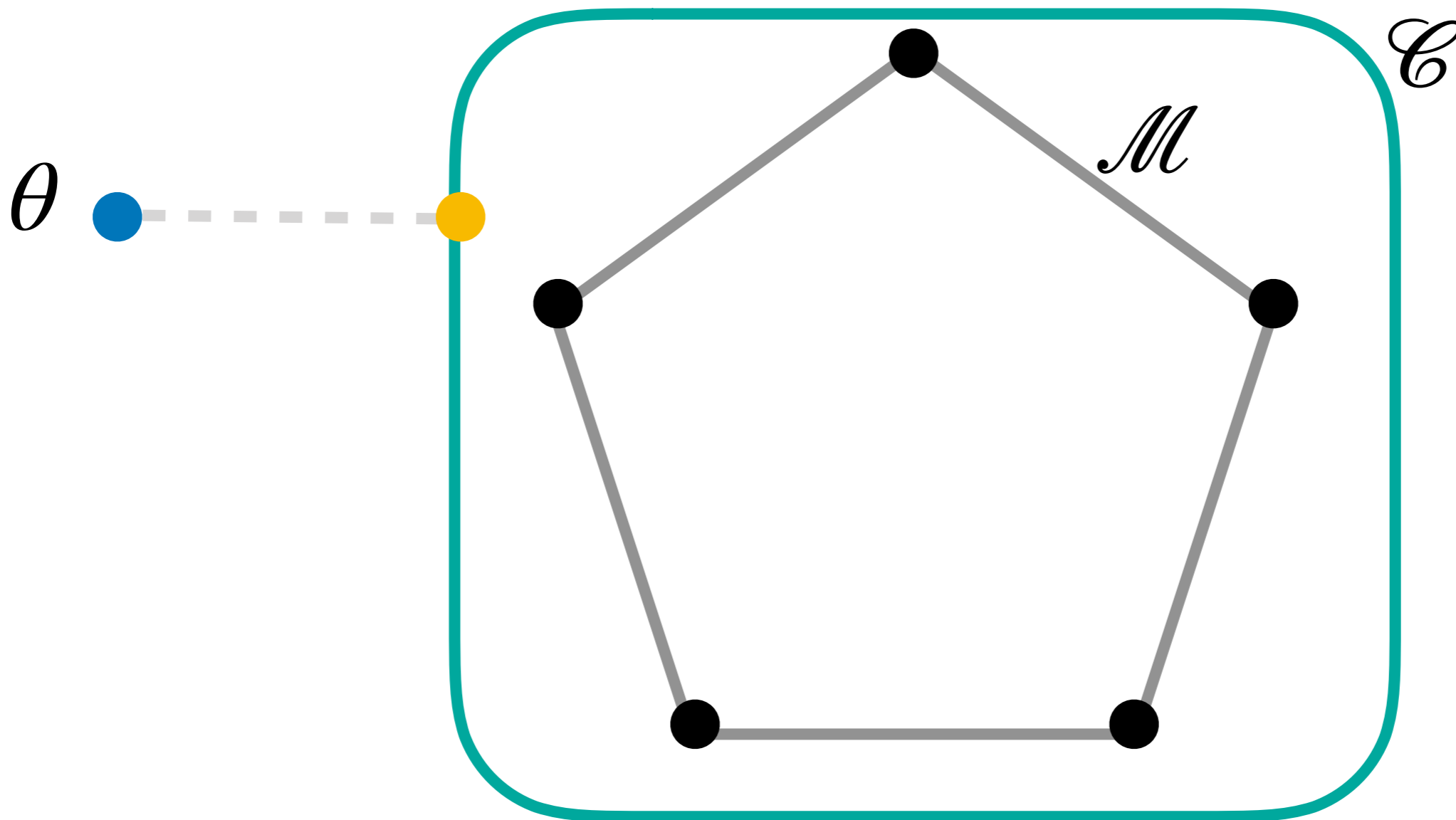Smallest convex set = **convex hull** (a.k.a. marginal polytope)



$$\mathcal{M} \triangleq conv(\varphi(\mathcal{Y}))$$

# Choice of the convex set

Smallest convex set = **convex hull** (a.k.a. marginal polytope)

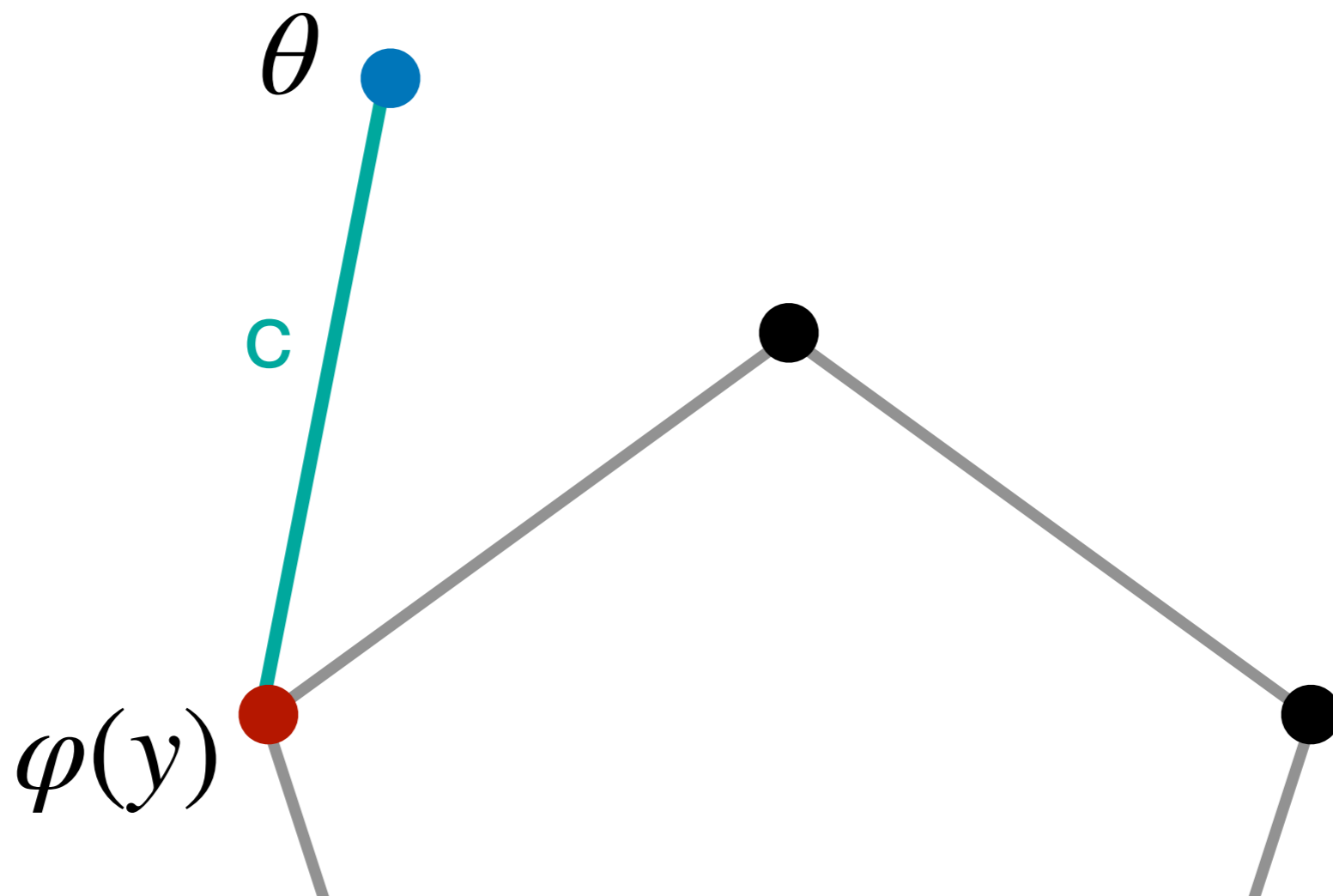Can use any **superset** with cheaper to compute projection



$$\mathcal{M} \triangleq conv(\varphi(\mathcal{Y}))$$

# Associated loss function

# Associated loss function
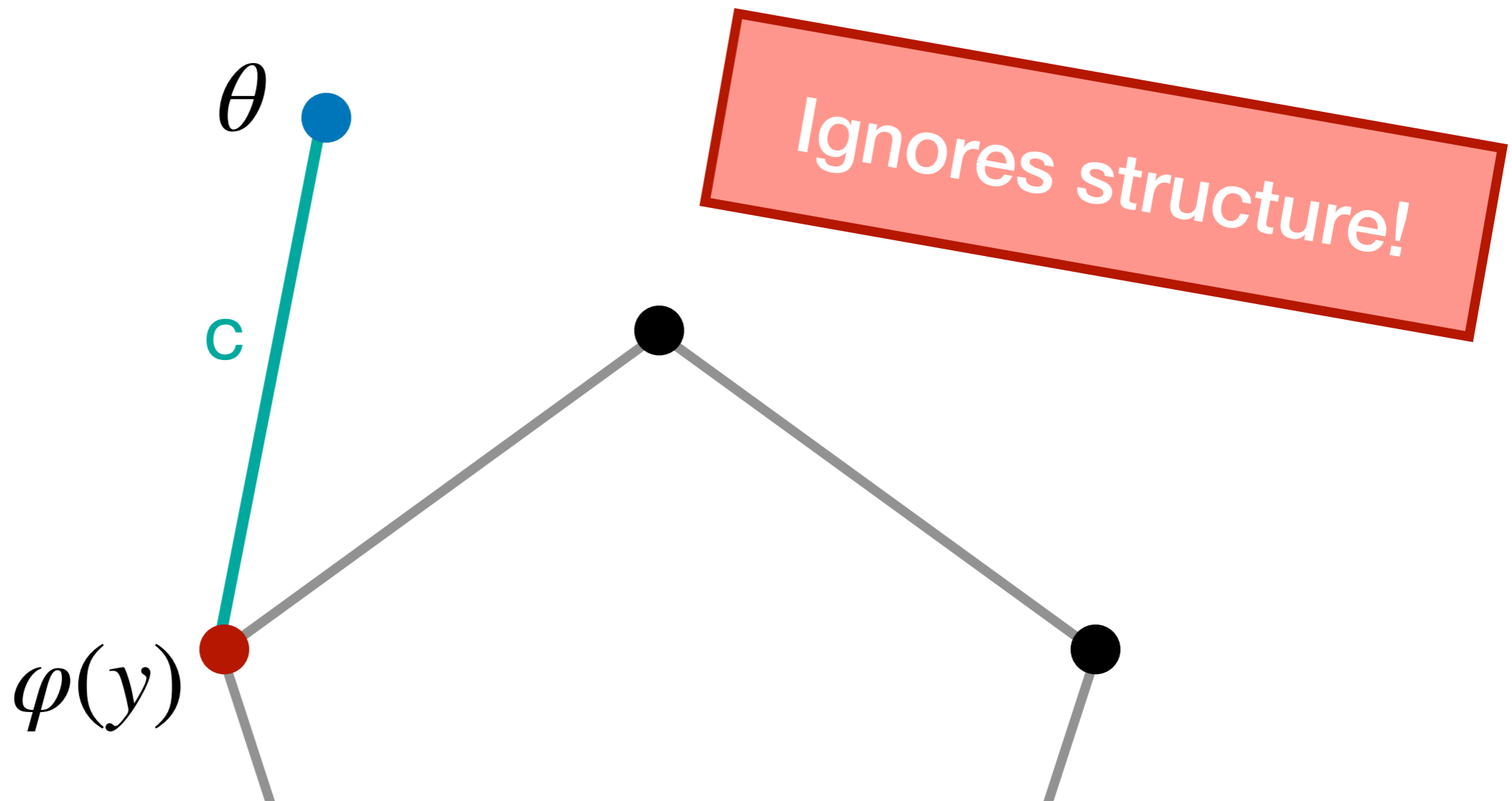
**Squared loss**

$$SQ(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - \theta\|^2 = c$$

# Associated loss function

**Squared loss**

$$SQ(\theta, y) \triangleq \frac{1}{2} \|\varphi(y) - \theta\|^2 = c$$

$\theta$

$c$

$\varphi(y)$

**Ignores structure!**

# Associated loss function
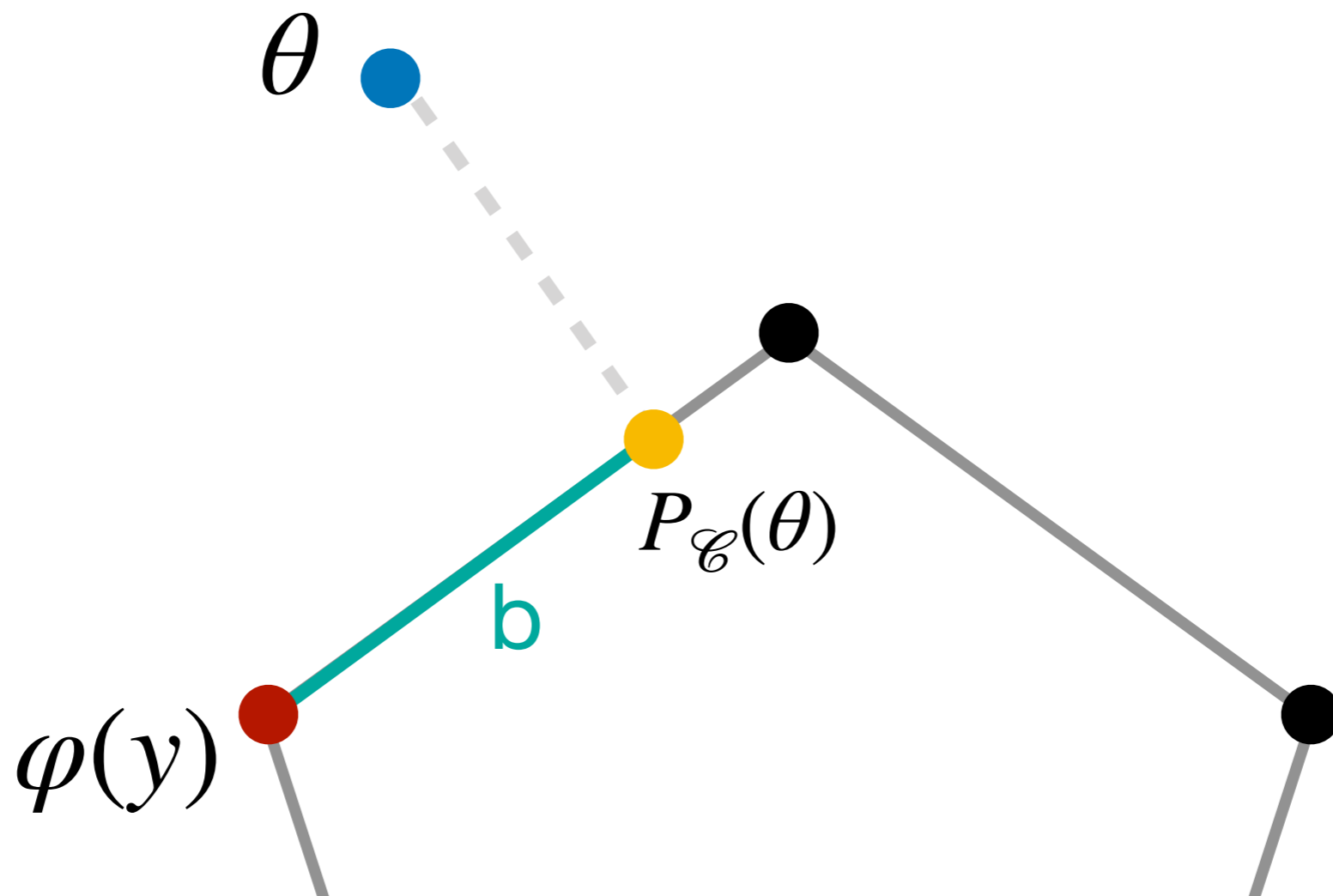
$$NC_{\mathscr{C}}(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - P_{\mathscr{C}}(\theta)\|^2 = b$$
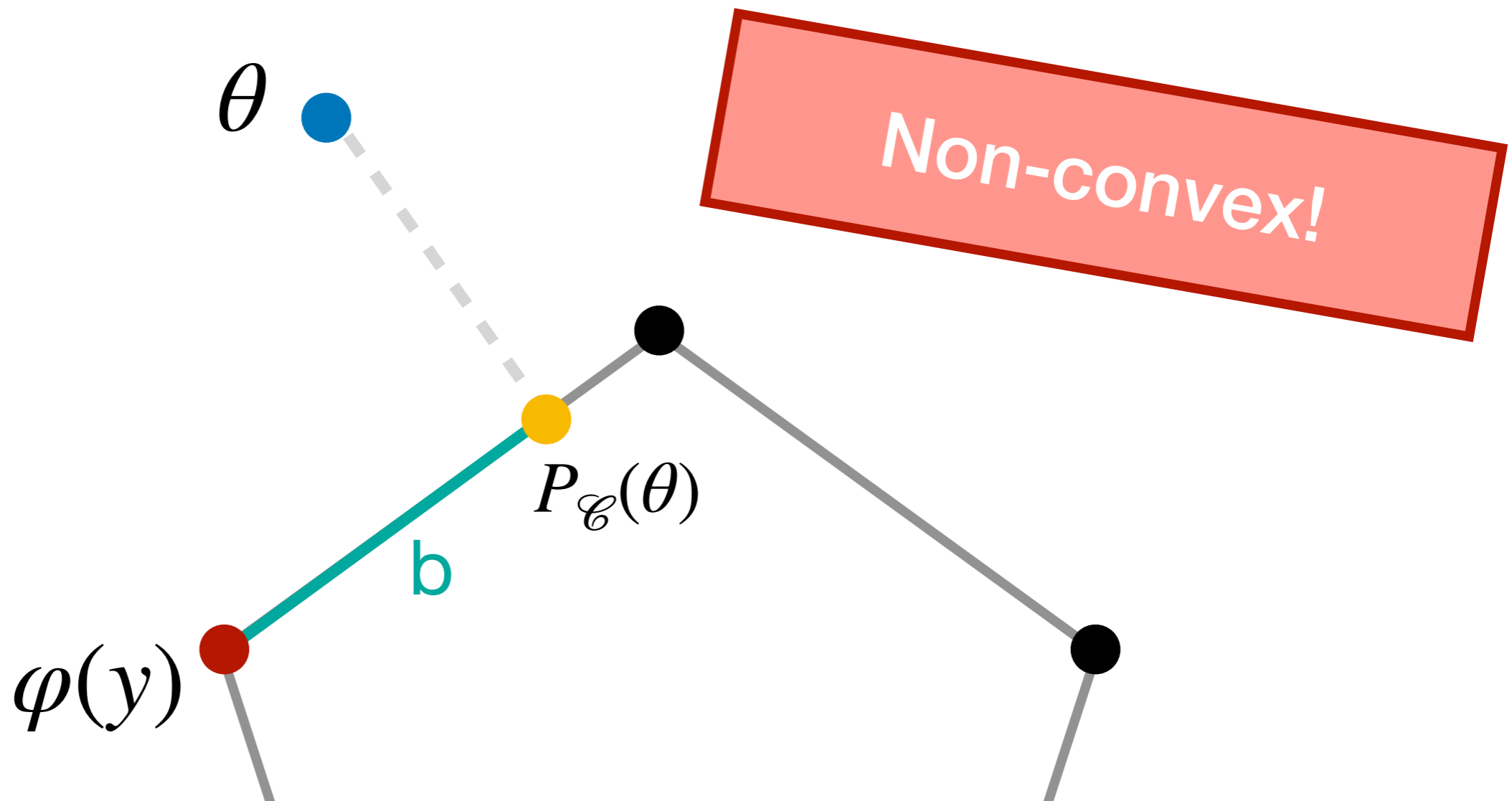
# Associated loss function

$$NC_{\mathscr{C}}(\theta, y) \triangleq \frac{1}{2}\|\varphi(y) - P_{\mathscr{C}}(\theta)\|^2 = b$$

# Associated loss function

**Proposed loss**

$$S_{\mathscr{C}}(\theta, y) \triangleq SQ(\theta, y) - \frac{1}{2}\|\theta - P_{\mathscr{C}}(\theta)\|^2 = c - a$$

# Associated loss function

**Proposed loss**

$$S_{\mathscr{C}}(\theta, y) \triangleq SQ(\theta, y) - \frac{1}{2}\|\theta - P_{\mathscr{C}}(\theta)\|^2 = c - a$$

**Generalized Pythagorean theorem**

$$a + b \leq c$$

$$\downarrow$$

$$NC_{\mathscr{C}}(\theta, y) \leq S_{\mathscr{C}}(\theta, y)$$

# Properties

1. $S_{\mathscr{C}}(\theta, y)$ is convex w.r.t. θ

2. $S_{\mathscr{C}}(\theta, y)$ is smooth w.r.t. θ (gradient is Lipschitz cont.)

3. $S_{\mathscr{C}}(\theta, y) \geq 0$

4. $S_{\mathscr{C}}(\theta, y) = 0 \Leftrightarrow P_{\mathscr{C}}(\theta) = \varphi(y)$

# Upper bounds

**Convex upper bound**

$$NC_{\mathscr{C}}(\theta, y) \leq S_{\mathscr{C}}(\theta, y) \quad \forall \theta, \varphi(y) \in \mathscr{C}$$

# Upper bounds

$$NC_{\mathscr{C}}(\theta, y) \leq S_{\mathscr{C}}(\theta, y) \quad \forall \theta, \varphi(y) \in \mathscr{C}$$

# Upper bounds

$$S_{\mathscr{C}}(\theta, y) \leq S_{\mathscr{C}'}(\theta, y) \quad \forall \mathscr{C} \subseteq \mathscr{C}'$$

# Upper bounds

**Superset upper bound**

$$S_{\mathscr{C}}(\theta, y) \leq S_{\mathscr{C}'}(\theta, y) \quad \forall \mathscr{C} \subseteq \mathscr{C}'$$

# Link with Fenchel duality

Let $\quad \Omega(u) \triangleq \dfrac{1}{2}\|u\|^2$ **if** $u \in \mathscr{C}$, $\infty$ **otherwise**

"Fenchel-Young losses", Blondel, Martins, Niculae, 2019

# Link with Fenchel duality

Let $\quad \Omega(u) \triangleq \dfrac{1}{2}\|u\|^2$ **if** $u \in \mathscr{C}$, $\infty$ **otherwise**

primal space
$\mathbf{dom}(\Omega) = \mathscr{C}$

$\varphi(y)$

"Fenchel-Young losses", Blondel, Martins, Niculae, 2019

# Link with Fenchel duality

Let $\Omega(u) \triangleq \dfrac{1}{2}\|u\|^2$ **if** $u \in \mathscr{C}$, $\infty$ **otherwise**

primal space
$\mathbf{dom}(\Omega) = \mathscr{C}$

dual space
$\mathbf{dom}(\Omega^*) = \mathbb{R}^2$

$\varphi(y)$

$\theta$

"Fenchel-Young losses", Blondel, Martins, Niculae, 2019

# Link with Fenchel duality

Let $\Omega(u) \triangleq \dfrac{1}{2}\|u\|^2$ **if** $u \in \mathscr{C}$, $\infty$ **otherwise**



primal space
$\mathbf{dom}(\Omega) = \mathscr{C}$

dual space
$\mathbf{dom}(\Omega^*) = \mathbb{R}^2$

$\nabla\Omega^* = P_\mathscr{C}$

$u$

$\varphi(y)$

$\theta$

"Fenchel-Young losses", Blondel, Martins, Niculae, 2019

# Link with Fenchel duality

Let $\Omega(u) \triangleq \dfrac{1}{2}\|u\|^2$ **if** $u \in \mathscr{C}$, $\infty$ **otherwise**

primal space
$\mathbf{dom}(\Omega) = \mathscr{C}$

dual space
$\mathbf{dom}(\Omega^*) = \mathbb{R}^2$

$\nabla \Omega^* = P_{\mathscr{C}}$

$u$

$\theta$

$\varphi(y)$

$$S_{\mathscr{C}}(\theta, y) = \Omega^*(\theta) + \Omega(\varphi(y)) - \langle \varphi(y), \theta \rangle$$

"Fenchel-Young losses", Blondel, Martins, Niculae, 2019

# Kullback Leibler geometry

$$\Omega(u) \triangleq \langle u, \log u \rangle \text{ if } u \in \mathscr{C}, \infty \textbf{ otherwise}$$

# Kullback Leibler geometry

$$\Omega(u) \triangleq \langle u, \log u \rangle \text{ if } u \in \mathscr{C}, \infty \text{ otherwise}$$

$$\nabla \Omega^*(\theta) = \arg\min_{u \in \mathscr{C}} KL(u, e^{\theta - 1})$$

# Kullback Leibler geometry

$$\Omega(u) \triangleq \langle u, \log u \rangle \textbf{ if } u \in \mathscr{C}, \ \infty \textbf{ otherwise}$$

$$\nabla \Omega^*(\theta) = \arg \min_{u \in \mathscr{C}} KL(u, e^{\theta-1})$$

$$S_{\mathscr{C}}(\theta, y) = \Omega^*(\theta) + \Omega(\varphi(y)) - \langle \varphi(y), \theta \rangle$$

# Kullback Leibler geometry

$$\Omega(u) \triangleq \langle u, \log u \rangle \textbf{ if } u \in \mathscr{C}, \infty \textbf{ otherwise}$$

$$\nabla \Omega^*(\theta) = \underset{u \in \mathscr{C}}{\arg\min} \, KL(u, e^{\theta-1})$$

$$S_{\mathscr{C}}(\theta, y) = \Omega^*(\theta) + \Omega(\varphi(y)) - \langle \varphi(y), \theta \rangle$$

**Proposition**

Let $\beta = \underset{u \in \mathscr{C}}{\max} \|u\|_1$. Then,

$S_{\mathscr{C}}(\theta, y)$ is $\beta$-smooth with respect to $\| \cdot \|_\infty$.

# Kullback Leibler geometry

$$\Omega(u) \triangleq \langle u, \log u \rangle \text{ if } u \in \mathscr{C}, \infty \text{ otherwise}$$

$$\nabla \Omega^*(\theta) = \arg\min_{u \in \mathscr{C}} KL(u, e^{\theta-1})$$

$$S_{\mathscr{C}}(\theta, y) = \Omega^*(\theta) + \Omega(\varphi(y)) - \langle \varphi(y), \theta \rangle$$

**Smaller set → smoother loss!**

**Proposition**

Let $\beta = \max_{u \in \mathscr{C}} \|u\|_1$. Then,

$S_{\mathscr{C}}(\theta, y)$ is $\beta$-smooth with respect to $\| \cdot \|_\infty$.

# Calibrated decoding

# Calibrated decoding

$$L(\hat{y}, y) = \langle \varphi(\hat{y}), V\varphi(y) + b \rangle + c(y)$$

# Calibrated decoding

**Affine decomposition of the target loss**

$$L(\hat{y}, y) = \langle \varphi(\hat{y}), V\varphi(y) + b \rangle + c(y)$$

**Decoding calibrated for loss L**

$$\hat{y}_L(u) \triangleq \arg\min_{y' \in \mathcal{Y}} \langle \varphi(y'), Vu + b \rangle$$

# Calibrated decoding

$$L(\hat{y}, y) = \langle \varphi(\hat{y}), V\varphi(y) + b \rangle + c(y)$$

$$\hat{y}_L(u) \triangleq \arg\min_{y' \in \mathcal{Y}} \langle \varphi(y'), Vu + b \rangle$$

$$= MAP(-Vu - b)$$

# Calibrated decoding

$$L(\hat{y}, y) = \langle \varphi(\hat{y}), V\varphi(y) + b \rangle + c(y)$$

Decoding calibrated for loss L

$$\hat{y}_L(u) \triangleq \arg\min_{y' \in \mathcal{Y}} \langle \varphi(y'), Vu + b \rangle$$

$$= MAP(-Vu - b)$$

Decomposition important both for computational tractability
and theoretical analysis

# Consistency

$$\delta\mathscr{L}(f) \triangleq \mathscr{L}(f) - \inf_{f' \colon \mathcal{X} \to \mathcal{Y}} \mathscr{L}(f')$$

$$\delta\mathcal{S}_{\mathscr{C}}(g) \triangleq \mathcal{S}_{\mathscr{C}}(g) - \inf_{g' \colon \mathcal{X} \to \Theta} \mathcal{S}_{\mathscr{C}}(g')$$

# Consistency

$$\delta\mathscr{L}(f) \triangleq \mathscr{L}(f) - \inf_{f' : \mathscr{X} \to \mathscr{Y}} \mathscr{L}(f')$$

$$\delta\mathcal{S}_{\mathscr{C}}(g) \triangleq \mathcal{S}_{\mathscr{C}}(g) - \inf_{g' : \mathscr{X} \to \Theta} \mathcal{S}_{\mathscr{C}}(g')$$

Calibration between excess risks

$$\forall g : \mathscr{X} \to \Theta : \quad \frac{\delta\mathscr{L}(dec \circ g)^2}{8\beta\sigma^2} \leq \delta\mathcal{S}_{\mathscr{C}}(g)$$

$$dec \triangleq \hat{y}_L \circ P_{\mathscr{C}} \qquad \beta \triangleq \text{Lipschitz constant of } P_{\mathscr{C}} \text{ w.r.t. } \|\cdot\| \qquad \sigma \triangleq \sup_{y \in \mathscr{Y}} \|V^{\top}\varphi(y)\|$$

# Probability simplex

$$\mathscr{Y} = [k] \triangleq \{1,\ldots,k\}$$

# Probability simplex

**Output set**

$$\mathscr{Y} = [k] \triangleq \{1,\ldots,k\}$$

**Encoding**

$$\varphi(y) = e_y$$

**Marginal polytope**

$$\mathscr{M} = conv(\varphi(\mathscr{Y})) = \triangle^k$$



$\varphi(2) = [0, 1, 0]$

$\varphi(1) = [1, 0, 0]$

$\varphi(3) = [0, 0, 1]$

# Probability simplex

**Output set**

$$\mathscr{Y} = [k] \triangleq \{1,\ldots,k\}$$

**Encoding**

$$\varphi(y) = e_y$$

**Marginal polytope**

$$\mathscr{M} = conv(\varphi(\mathscr{Y})) = \triangle^k$$

$\varphi(2) = [0, 1, 0]$

$\varphi(1) = [1, 0, 0]$

$\varphi(3) = [0, 0, 1]$

**Oracles**

MAP: O(k)

Euclidean: sparsemax, O(k) or O(k log k)

KL: softmax, O(k)

# Unit cube

$$\mathscr{Y} = 2^{[k]}$$

# Unit cube

**Output set**

$$\mathscr{Y} = 2^{[k]}$$

**Encoding**

$$\varphi(y) = \sum_{i=1}^{|y|} e_{y_i}$$

**Marginal polytope**

$$\mathscr{M} = [0,1]^k$$

# Unit cube

**Output set**

$$\mathcal{Y} = 2^{[k]}$$

**Encoding**

$$\varphi(y) = \sum_{i=1}^{|y|} e_{y_i}$$

**Marginal polytope**

$$\mathcal{M} = [0,1]^k$$



**Oracles**

MAP: O(k)

Euclidean: clipping to [0,1], O(k)

KL: O(k)

# Budget polytope

$$\mathscr{Y} = \{y \in 2^{[k]} : l \le |y| \le u\}$$

# Budget polytope

$$\mathcal{Y} = \{y \in 2^{[k]} : l \leq |y| \leq u\}$$

$$\varphi(y) = \sum_{i=1}^{|y|} e_{y_i}$$

$$\mathcal{M} = \{y \in [0,1]^k : l \leq y^\top \mathbf{1} \leq m\}$$

# Budget polytope

**Output set**

$$\mathcal{Y} = \{y \in 2^{[k]} : l \le |y| \le u\}$$

**Encoding**

$$\varphi(y) = \sum_{i=1}^{|y|} e_{y_i}$$

**Marginal polytope**

$$\mathcal{M} = \{y \in [0,1]^k : l \le y^\top \mathbf{1} \le m\}$$

**Oracles**

MAP: O(k log k)

Euclidean: O(k)

KL: O(k log k)

# Order simplex

$$\mathscr{Y} = [k] \qquad 1 \prec \dots \prec k$$

# Order simplex

$$\mathcal{Y} = [k] \qquad 1 < \ldots < k$$

Encoding

$$\varphi(y) = \sum_{1 \le i < y \le k} e_i \in \mathbb{R}^{k-1}$$

Marginal polytope

$$\mathcal{M} = \{\mu \in \mathbb{R}^{k-1} : 1 \ge \mu_1 \ge \mu_2 \ge \ldots \ge \mu_{k-1} \ge 0\}$$



$\varphi(3) = [1, 1, 0]$

$\varphi(1) = [0, 0, 0]$

$\varphi(4) = [1, 1, 1]$

$\varphi(2) = [1, 0, 0]$

# Order simplex

**Output set**

$$\mathcal{Y} = [k] \quad 1 < \ldots < k$$

**Encoding**

$$\varphi(y) = \sum_{1 \le i < y \le k} e_i \in \mathbb{R}^{k-1}$$

$\varphi(3) = [1, 1, 0]$

$\varphi(1) = [0, 0, 0]$

$\varphi(4) = [1, 1, 1]$

$\varphi(2) = [1, 0, 0]$

**Marginal polytope**

$$\mathcal{M} = \{\mu \in \mathbb{R}^{k-1} : 1 \ge \mu_1 \ge \mu_2 \ge \ldots \ge \mu_{k-1} \ge 0\}$$

**Oracles**

MAP: O(k)

Eucl: isotonic reg, O(k)

KL: isotonic optimization

# Birkhoff polytope

$$\mathscr{Y} = \textbf{Permutations}([m])$$

# Birkhoff polytope

$$\mathscr{Y} = \textbf{Permutations}([m])$$

**Encoding**

$$\varphi(y) = \quad \text{permutation matrix associated with y}$$

**Marginal polytope**

$$\mathscr{M} = \{P \in \mathbb{R}^{m \times m} : P^\top \mathbf{1}_m = 1, P\mathbf{1}_m = 1, 0 \leq P \leq 1\}$$



$\varphi((2,1,3))$

$\varphi((2,3,1))$

$\varphi((3,2,1))$

$\varphi((1,2,3))$

$\varphi((1,3,2))$

$\varphi((3,1,2))$

# Birkhoff polytope

$$\mathscr{Y} = \textbf{Permutations}([m])$$

## Encoding

$$\varphi(y) = \quad \text{permutation matrix associated with y}$$

## Marginal polytope

$$\mathscr{M} = \{P \in \mathbb{R}^{m \times m} : P^\top \mathbf{1}_m = 1, P\mathbf{1}_m = 1, 0 \le P \le 1\}$$



$\varphi((2,1,3))$    $\varphi((2,3,1))$

$\varphi((3,2,1))$

$\varphi((1,2,3))$

$\varphi((1,3,2))$    $\varphi((3,1,2))$

## Oracles

MAP: Hungarian, O(m³)

Eucl: LBFGS dual, O(m²/ε)

KL: Sinkhorn, O(m²/ε)

Marginal: **intractable**

# Birkhoff polytope

$$\mathcal{Y} = \textbf{Permutations}([m])$$

### Encoding

$$\varphi(y) = \quad \text{permutation matrix associated with y}$$

### Marginal polytope

$$\mathcal{M} = \{P \in \mathbb{R}^{m \times m} : P^\top \mathbf{1}_m = 1, P\mathbf{1}_m = 1, 0 \leq P \leq 1\}$$

$$\triangle^{m \times m} \triangleq \{P \in \mathbb{R}^{m \times m} : P^\top \mathbf{1}_m = 1, 0 \leq P \leq 1\} \supset \mathcal{M}$$

Row-stochastic matrices



$\varphi((2,1,3))$        $\varphi((2,3,1))$

$\varphi((3,2,1))$

$\varphi((1,2,3))$

$\varphi((1,3,2))$        $\varphi((3,1,2))$

### Oracles

MAP: Hungarian, O(m³)

Eucl: LBFGS dual, O(m²/ε)

KL: Sinkhorn, O(m²/ε)

Marginal: **intractable**

# Permutahedron

$$\mathscr{Y} = \textbf{Permutations}([m])$$

# Permutahedron

**Output set**

$$\mathcal{Y} = \textbf{Permutations}([m])$$

**Encoding**

$$\varphi(y) = \quad \text{permutation of a vector w according to y}$$

w = [1, 0, 0]

w = [2, 1, 0] / 3

w = [1, 1, 0] / 2

w = [3, 2, 1] / 6

[0, 1, 0]

[1, 0, 0]

[0, 0, 1]

**Marginal polytope**

$$\mathcal{M} = \{\mu \in \mathbb{R}^m : \sum_{i \in S} \mu_i \leq \sum_{i=1}^{|S|} w_i \, \forall S \subset [m], \sum_{i=1}^{m} \mu_i = \sum_{i=1}^{m} w_i\}$$

# Permutahedron

## Output set

$$\mathscr{Y} = \textbf{Permutations}([m])$$

## Encoding

$$\varphi(y) = \quad \text{permutation of a vector w according to y}$$

w = [1, 0, 0]

w = [2, 1, 0] / 3

w = [1, 1, 0] / 2

w = [3, 2, 1] / 6

[0, 1, 0]

[1, 0, 0]

[0, 0, 1]

## Marginal polytope

$$\mathscr{M} = \{\mu \in \mathbb{R}^m : \sum_{i \in S} \mu_i \leq \sum_{i=1}^{|S|} w_i \, \forall S \subset [m], \sum_{i=1}^{m} \mu_i = \sum_{i=1}^{m} w_i\}$$

## Oracles

MAP: O(m log m)

Eucl: isotonic reg, O(m log m)

KL: isotonic optimization

# Outline

1. Background

2. Proposed framework

3. Experiments

# Experiments

$$\frac{1}{n} \sum_{i=1}^{n} S_{\mathscr{C}}(Wx_i, y_i) + \lambda \|W\|_F^2$$

- Label ranking

- Ordinal regression

- Multilabel classification

# Label ranking

Full-ranking supervision setting (no relevance scores)

e.g. $2 \succ 1 \succ 3 \succ 4$

# Label ranking

Full-ranking supervision setting (no relevance scores)

e.g. $\quad 2 \succ 1 \succ 3 \succ 4$

| Projection Decoding | $\mathbb{R}^{m \times m}$ $\mathcal{B}$ |
|---|---|
| Authorship | 5.70 |
| Glass | 7.11 |
| Iris | 19.26 |
| Vehicle | 9.04 |
| Vowel | 10.57 |
| Wine | **1.23** |

= squared loss

L = Hamming loss        Using Euclidean projections

# Label ranking

Full-ranking supervision setting (no relevance scores)

e.g. $2 > 1 > 3 > 4$

| Projection<br>Decoding | $\mathbb{R}^{m \times m}$<br>$\mathcal{B}$ | $[0, 1]^{m \times m}$<br>$\mathcal{B}$ |
|---|---|---|
| Authorship | 5.70 | 5.18 |
| Glass | 7.11 | 5.68 |
| Iris | 19.26 | 4.44 |
| Vehicle | 9.04 | 7.57 |
| Vowel | 10.57 | 9.56 |
| Wine | **1.23** | 1.85 |

= squared loss

L = Hamming loss        Using Euclidean projections

# Label ranking

Full-ranking supervision setting (no relevance scores)

e.g. $2 \succ 1 \succ 3 \succ 4$

| Projection<br>Decoding | $\mathbb{R}^{m \times m}$<br>$\mathcal{B}$ | $[0,1]^{m \times m}$<br>$\mathcal{B}$ | $\triangle^{m \times m}$<br>$\mathcal{B}$ |
|---|---|---|---|
| Authorship | 5.70 | 5.18 | 5.70 |
| Glass | 7.11 | 5.68 | 5.04 |
| Iris | 19.26 | 4.44 | **1.48** |
| Vehicle | 9.04 | 7.57 | 6.99 |
| Vowel | 10.57 | 9.56 | 9.18 |
| Wine | **1.23** | 1.85 | 1.85 |

$\mathcal{B}$ = squared loss

L = Hamming loss          Using Euclidean projections

# Label ranking

Full-ranking supervision setting (no relevance scores)

e.g. $2 > 1 > 3 > 4$

| Projection<br>Decoding | $\mathbb{R}^{m \times m}$<br>$\mathcal{B}$ | $[0,1]^{m \times m}$<br>$\mathcal{B}$ | $\triangle^{m \times m}$<br>$\mathcal{B}$ | $\mathcal{B}$<br>$\mathcal{B}$ |
|---|---|---|---|---|
| Authorship | 5.70 | 5.18 | 5.70 | **5.10** |
| Glass | 7.11 | 5.68 | 5.04 | **4.65** |
| Iris | 19.26 | 4.44 | **1.48** | 2.96 |
| Vehicle | 9.04 | 7.57 | 6.99 | **5.88** |
| Vowel | 10.57 | 9.56 | 9.18 | **8.76** |
| Wine | **1.23** | 1.85 | 1.85 | 1.85 |

= squared loss

L = Hamming loss        Using Euclidean projections

# Label ranking

| | Euclidean | vs. | KL |
|---|---|---|---|
| Projection | $\mathcal{B}$ | | $\mathcal{B}$ |
| Decoding | $\mathcal{B}$ | | $\mathcal{B}$ |
| Authorship | **5.10** | | **5.10** |
| Glass | **4.65** | | **4.65** |
| Iris | 2.96 | | 2.96 |
| Vehicle | **5.88** | | **6.25** |
| Vowel | **8.76** | | **9.17** |
| Wine | 1.85 | | **1.85** |

# Label ranking



$$P_{\mathscr{B}}(Wx)$$

| | Euclidean | vs. | KL |
|---|---|---|---|
| Projection | $\mathcal{B}$ | | $\mathcal{B}$ |
| Decoding | $\mathcal{B}$ | | $\mathcal{B}$ |
| Authorship | **5.10** | | **5.10** |
| Glass | **4.65** | | **4.65** |
| Iris | 2.96 | | 2.96 |
| Vehicle | **5.88** | | **6.25** |
| Vowel | **8.76** | | **9.17** |
| Wine | 1.85 | | **1.85** |

"soft permutation matrix"

# Label ranking

Birkhoff vs. permutahedron

|  |  | Linear | Poly 2 | Poly 3 |
|---|---|---|---|---|
| Projection | $\mathcal{B}$ | $\mathcal{P}$ | $\mathcal{P}$ | $\mathcal{P}$ |
| Decoding | $\mathcal{B}$ | $\mathcal{P}$ | $\mathcal{P}$ | $\mathcal{P}$ |
| Authorship | **5.10** | 10.06 | 10.50 | 8.59 |
| Glass | **4.65** | 7.49 | 7.10 | 8.14 |
| Iris | 2.96 | 27.41 | 20.00 | 5.93 |
| Vehicle | **5.88** | 11.62 | 8.30 | 9.26 |
| Vowel | **8.76** | 14.35 | 11.74 | 10.21 |
| Wine | 1.85 | 8.02 | 3.08 | 6.79 |
|  | $W \in \mathbb{R}^{p \times m^2}$ | $W \in \mathbb{R}^{p \times m}$ | $W \in \mathbb{R}^{n \times m}$ | $W \in \mathbb{R}^{n \times m}$ |

Using Euclidean projections

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \dots \prec k$$

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \ldots \prec k$$

|  | Projection Decoding Baseline |
| --- | --- |
| Average MAE | 0.78 |
| Average rank | 4.75 |

Averaged over 16 datasets

L = MAE = Mean Absolute Error

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \ldots \prec k$$

| Projection Decoding | Baseline | $\mathbb{R}$ Round |
|---|---|---|
| Average MAE | 0.78 | 0.72 |
| Average rank | 4.75 | 2.9 |

Averaged over 16 datasets

L = MAE = Mean Absolute Error

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \ldots \prec k$$

| Projection Decoding | Baseline | $\mathbb{R}$ Round |
|---|---|---|
| Average MAE | 0.78 | 0.72 |
| Average rank | 4.75 | 2.9 |

Averaged over 16 datasets

L = MAE = Mean Absolute Error          OS = Order Simplex

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \dots \prec k$$

| Projection Decoding | Baseline | $\mathbb{R}$ Round | $\mathbb{R}^{k-1}$ $\mathcal{OS}$ |
|---|---|---|---|
| Average MAE | 0.78 | 0.72 | 0.47 |
| Average rank | 4.75 | 2.9 | 2.1 |

Averaged over 16 datasets

L = MAE = Mean Absolute Error          OS = Order Simplex

# Ordinal regression

$$\mathcal{Y} = [k] \qquad 1 \prec \ldots \prec k$$

| Projection Decoding | Baseline | $\mathbb{R}$ Round | $\mathbb{R}^{k-1}$ $\mathcal{OS}$ | $[0, 1]^{k-1}$ $\mathcal{OS}$ |
|---|---|---|---|---|
| Average MAE | 0.78 | 0.72 | 0.47 | 0.45 |
| Average rank | 4.75 | 2.9 | 2.1 | 1.6 |

Averaged over 16 datasets

L = MAE = Mean Absolute Error        OS = Order Simplex

# Ordinal regression

$$\mathscr{Y} = [k] \qquad 1 \prec \ldots \prec k$$

| Projection Decoding | Baseline | $\mathbb{R}$ Round | $\mathbb{R}^{k-1}$ $\mathcal{OS}$ | $[0,1]^{k-1}$ $\mathcal{OS}$ | $\mathcal{OS}$ $\mathcal{OS}$ |
|---|---|---|---|---|---|
| Average MAE | 0.78 | 0.72 | 0.47 | 0.45 | 0.43 |
| Average rank | 4.75 | 2.9 | 2.1 | 1.6 | **1.5** |

Averaged over 16 datasets

L = MAE = Mean Absolute Error      OS = Order Simplex

# Multilabel classification

lower bound = 0          upper bound = $\left\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]} \right\rceil$

$F_1$ score

# Multilabel classification

lower bound = 0     upper bound = $\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]} \rceil$

| | |
|---|---|
| Projection | $[0, 1]^k$ |
| Decoding | $[0, 1]^k$ |
| Birds | 38.87 |
| Emotions | 56.60 |
| Scene | 61.06 |

$\mathcal{K}$: budget polytope

$F_1$ score

# Multilabel classification

lower bound = 0 $\qquad$ upper bound = $\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]}\,\rceil$

|  | $[0,1]^k$ | $\mathbb{R}^k$ |
|---|---|---|
| Projection | $[0,1]^k$ | $\mathbb{R}^k$ |
| Decoding | $[0,1]^k$ | $\mathcal{K}$ |
| Birds | 38.87 | 37.75 |
| Emotions | 56.60 | 51.73 |
| Scene | 61.06 | 50.33 |

$\mathcal{K}$: budget polytope

$F_1$ score

# Multilabel classification

lower bound = 0        upper bound = $\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]} \rceil$

| Projection | $[0,1]^k$ | $\mathbb{R}^k$ | $[0,1]^k$ |
|---|---|---|---|
| Decoding | $[0,1]^k$ | $\mathcal{K}$ | $\mathcal{K}$ |
| Birds | 38.87 | 37.75 | 39.21 |
| Emotions | 56.60 | 51.73 | 53.98 |
| Scene | 61.06 | 50.33 | 58.95 |

$\mathcal{K}$: budget polytope

$F_1$ score

# Multilabel classification

lower bound = 0     upper bound = $\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]} \rceil$

| Projection | $[0,1]^k$ | $\mathbb{R}^k$ | $[0,1]^k$ | $\mathcal{K}$ |
|------------|-----------|----------------|-----------|---------------|
| Decoding   | $[0,1]^k$ | $\mathcal{K}$  | $\mathcal{K}$ | $\mathcal{K}$ |
| Birds      | 38.87     | 37.75          | 39.21     | **39.43**     |
| Emotions   | 56.60     | 51.73          | 53.98     | **62.57**     |
| Scene      | 61.06     | 50.33          | 58.95     | **69.01**     |

$\mathcal{K}$: budget polytope

$F_1$ score

# Multilabel classification

lower bound = 0          upper bound = $\lceil \mathbb{E}[\,|Y|\,] + \sqrt{\mathbb{V}[\,|Y|\,]}\, \rceil$

| Projection | $[0,1]^k$ | $\mathbb{R}^k$ | $[0,1]^k$ | $\mathcal{K}$ |
|---|---|---|---|---|
| Decoding | $[0,1]^k$ | $\mathcal{K}$ | $\mathcal{K}$ | $\mathcal{K}$ |
| Birds | 38.87 | 37.75 | 39.21 | **39.43** |
| Cal500 | 34.62 | **35.86** | 34.63 | 34.61 |
| Emotions | 56.60 | 51.73 | 53.98 | **62.57** |
| Mediamill | **56.22** | 55.35 | **56.22** | 54.53 |
| Scene | 61.06 | 50.33 | 58.95 | **69.01** |
| TMC | **60.45** | 58.61 | 60.37 | 60.25 |
| Yeast | **60.24** | 60.20 | 60.23 | 60.06 |

$\mathcal{K}$: budget polytope

$F_1$ score

# Conclusion

# Conclusion

- We proposed a generic framework for deriving a **loss** from the **projection** onto a convex set

# Conclusion

- We proposed a generic framework for deriving a **loss** from the **projection** onto a convex set

- If its projection is affordable, the **marginal polytope** is the best convex set

# Conclusion

- We proposed a generic framework for deriving a **loss** from the **projection** onto a convex set

- If its projection is affordable, the **marginal polytope** is the best convex set

- If not, any convex **superset** with cheaper projection can be used (e.g., unit cube)