

# Geometric losses for distributional learning

*Arthur Mensch*<sup>(1)</sup>, Mathieu Blondel<sup>(2)</sup>, Gabriel Peyré<sup>(1)</sup>

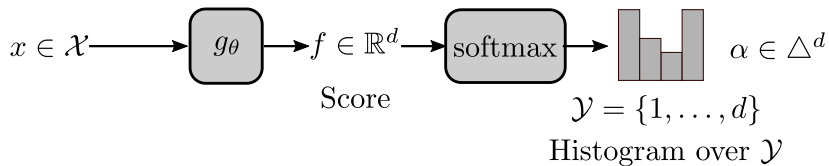
(1) École Normale Supérieure, DMA  
Centre National pour la Recherche Scientifique  
Paris, France

(2) NTT Communication Science Laboratories  
Kyoto, Japan

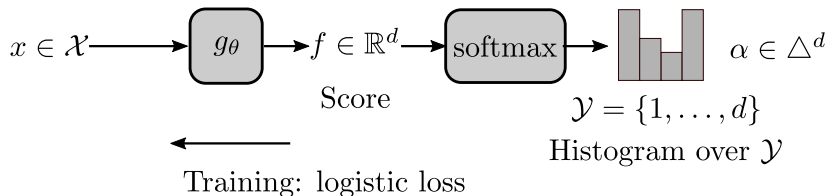


- 1 Introduction
- 2 Fenchel-Young losses for distribution spaces
- 3 Geometric softmax from Sinkhorn negentropies
- 4 Applications

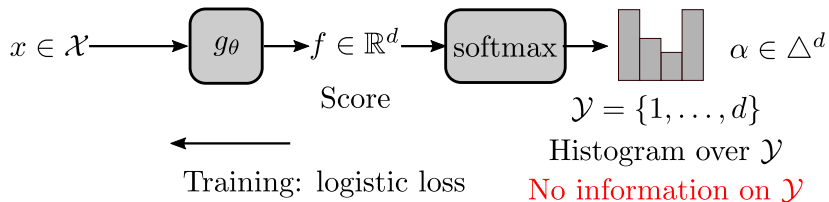
# Introduction: Predicting distributions



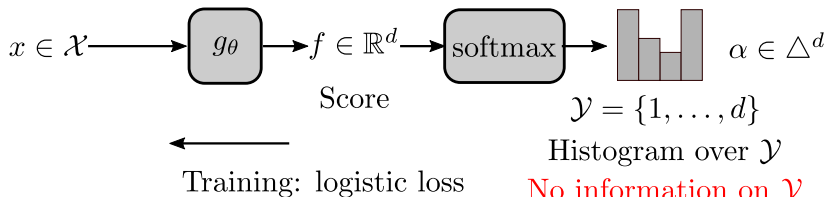
# Introduction: Predicting distributions



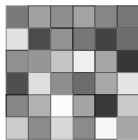
# Introduction: Predicting distributions



# Introduction: Predicting distributions



No information on  $\mathcal{Y}$

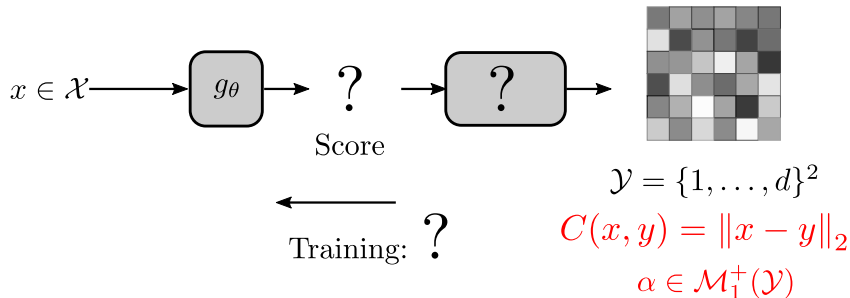
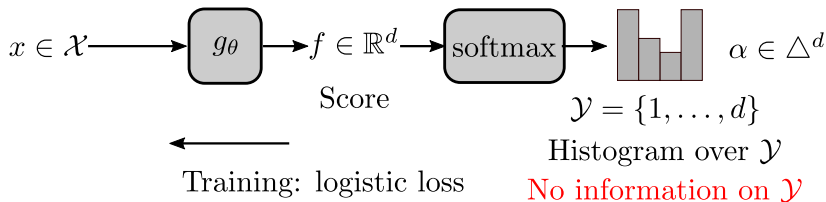


$$\mathcal{Y} = \{1, \dots, d\}^2$$

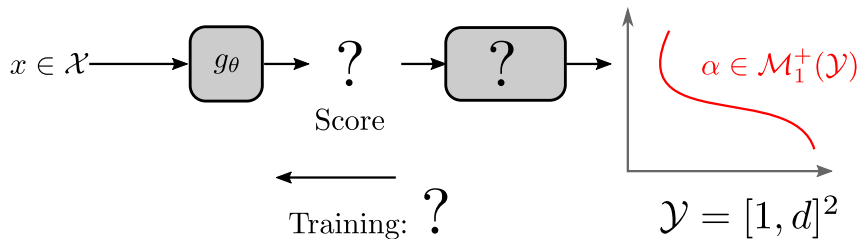
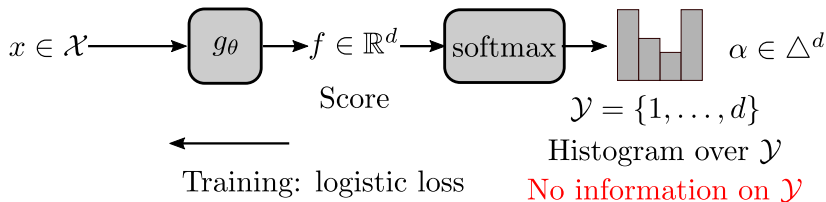
$$C(x, y) = \|x - y\|_2$$

$$\alpha \in \mathcal{M}_1^+(\mathcal{Y})$$

# Introduction: Predicting distributions



# Introduction: Predicting distributions





# Contribution: losses and links for continuous metrized output

## Handling output geometry

- Link and loss with cost between classes

$$C : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

- Output distribution over continuous space  $\mathcal{Y}$

# Contribution: losses and links for continuous metrized output

## Handling output geometry

- Link and loss with cost between classes

$$C : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

- Output distribution over continuous space  $\mathcal{Y}$

## New geometric losses and associated link functions:

- 1 Construction from duality between distributions and scores

# Contribution: losses and links for continuous metrized output

## Handling output geometry

- Link and loss with cost between classes

$$C : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

- Output distribution over continuous space  $\mathcal{Y}$


## New geometric losses and associated link functions:

- 1 Construction from duality between distributions and scores
- 2 **Need:** Convex functional on distribution space
  - Provided by **regularized optimal transport**

# Background: learning with a cost over outputs $\mathcal{Y}$

**Cost augmentation** of losses<sup>1,2</sup>:

- Convex cost-aware loss  $L_c : [1, d] \times \mathbb{R}^d \rightarrow \mathbb{R}$

 **Undefined link functions:**  $\mathbb{R}^d \rightarrow \Delta^d$ : what to predict at test time ?

---


<sup>1</sup>Ioannis Tsochantaridis et al. "Large margin methods for structured and interdependent output variables". In: *JMLR* (2005).

<sup>2</sup>Kevin Gimpel and Noah A Smith. "Softmax-margin CRFs: Training log-linear models with cost functions". In: *NAACL*. 2010.

# Background: learning with a cost over outputs $\mathcal{Y}$

**Cost augmentation** of losses<sup>1,2</sup>:


- Convex cost-aware loss  $L_c : [1, d] \times \mathbb{R}^d \rightarrow \mathbb{R}$

 **Undefined link functions:**  $\mathbb{R}^d \rightarrow \Delta^d$ : what to predict at test time ?

Use a Wasserstein **distance between output distributions**<sup>3</sup>:

- Ground metric  $C$  defines a distance  $W_C$  between distributions
- Prediction with a softmax link

$$\ell(\alpha, f) \triangleq W_C(\text{softmax}(f), \alpha)$$

 **Non-convex loss** and costly to compute

---

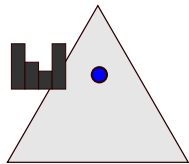
<sup>1</sup>Ioannis Tschantaris et al. "Large margin methods for structured and interdependent output variables". In: *JMLR* (2005).

<sup>2</sup>Kevin Gimpel and Noah A Smith. "Softmax-margin CRFs: Training log-linear models with cost functions". In: *NAACL*. 2010.

<sup>3</sup>Charlie Frogner et al. "Learning with a Wasserstein loss". In: *NIPS*. 2015.

- 1 Introduction
- 2 Fenchel-Young losses for distribution spaces**
- 3 Geometric softmax from Sinkhorn negentropies
- 4 Applications

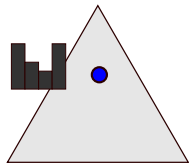
# Predicting distributions from topological duality



Primal histograms

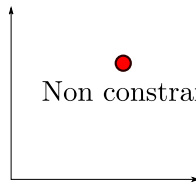
$$\Delta^3, \|\cdot\|_1$$

# Predicting distributions from topological duality



Primal histograms

$$\Delta^3, \|\cdot\|_1$$



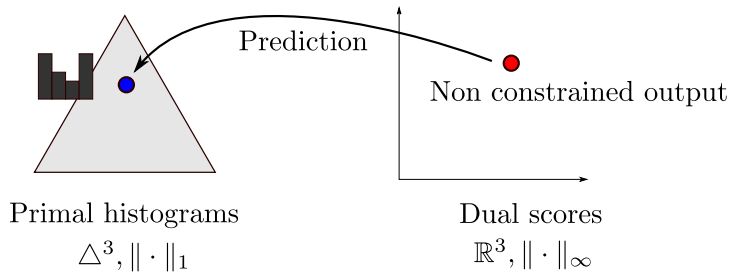
Non constrained output

Dual scores

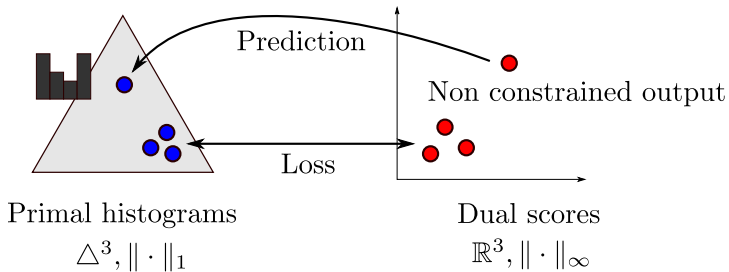
$$\mathbb{R}^3, \|\cdot\|_\infty$$



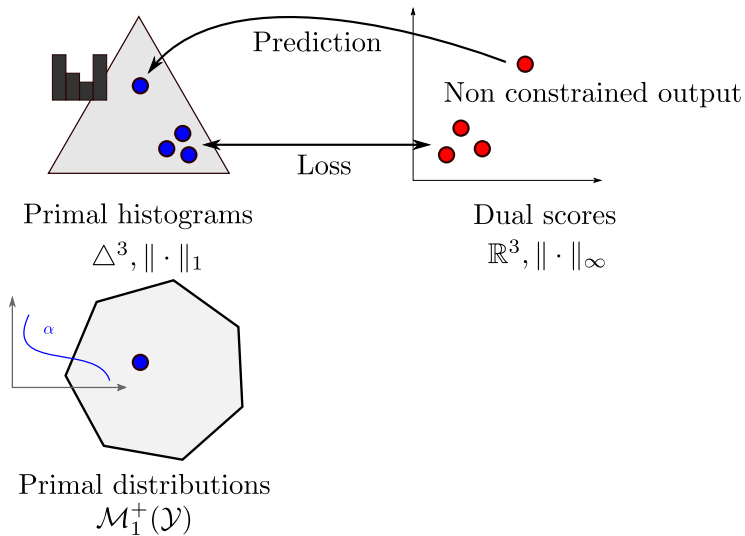
# Predicting distributions from topological duality



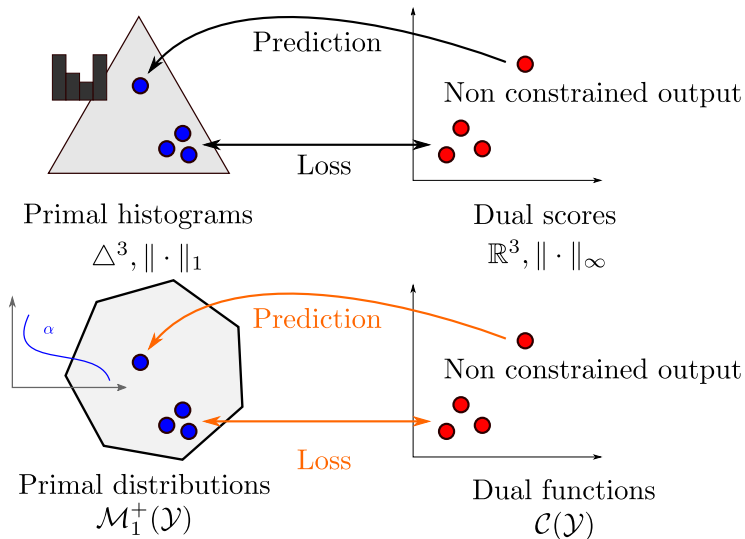
# Predicting distributions from topological duality



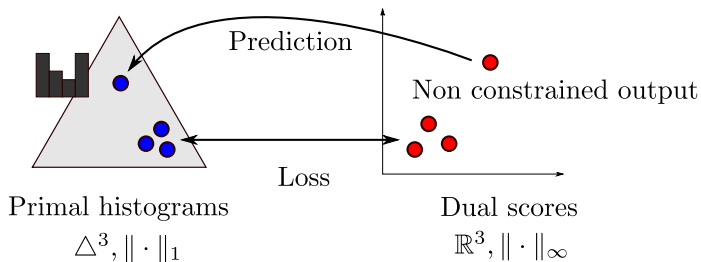
# Predicting distributions from topological duality



# Predicting distributions from topological duality



# All you need is a convex functional

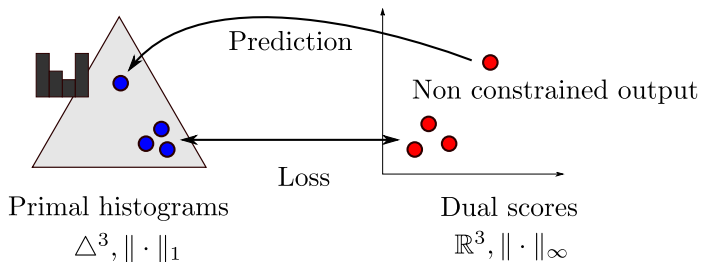


**Fenchel-Young losses**<sup>45</sup>: Convex function  $\Omega : \Delta^d \rightarrow \mathbb{R}$

<sup>4</sup>John C. Duchi et al. "Multiclass Classification, Information, Divergence, and Surrogate Risk". In: *Annals of Statistics* (2018).

<sup>5</sup>Mathieu Blondel et al. "Learning Classifiers with Fenchel-Young Losses: Generalized Entropies, Margins, and Algorithms". In: *AISTATS*. 2019.

# All you need is a convex functional



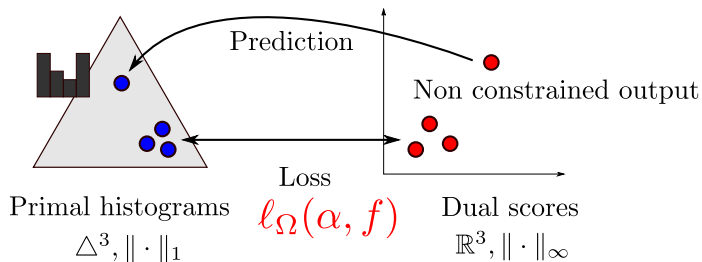
**Fenchel-Young losses**<sup>45</sup>: Convex function  $\Omega : \Delta^d \rightarrow \mathbb{R}$  and conjugate

$$\Omega^*(f) = \min_{\alpha \in \Delta^d} \Omega(\alpha) - \langle \alpha, f \rangle$$

<sup>4</sup>John C. Duchi et al. "Multiclass Classification, Information, Divergence, and Surrogate Risk". In: *Annals of Statistics* (2018).

<sup>5</sup>Mathieu Blondel et al. "Learning Classifiers with Fenchel-Young Losses: Generalized Entropies, Margins, and Algorithms". In: *AISTATS*. 2019.

# All you need is a convex functional



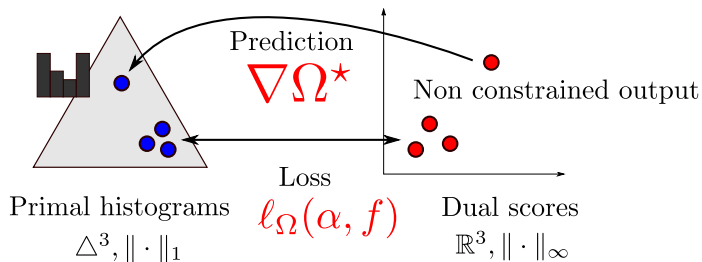
**Fenchel-Young losses**<sup>45</sup>: Convex function  $\Omega : \Delta^d \rightarrow \mathbb{R}$  and conjugate

$$\Omega^*(f) = \min_{\alpha \in \Delta^d} \Omega(\alpha) - \langle \alpha, f \rangle \quad \ell_\Omega(\alpha, f) = \Omega(\alpha) + \Omega^*(f) - \langle \alpha, f \rangle \geq 0$$

<sup>4</sup>John C. Duchi et al. "Multiclass Classification, Information, Divergence, and Surrogate Risk". In: *Annals of Statistics* (2018).

<sup>5</sup>Mathieu Blondel et al. "Learning Classifiers with Fenchel-Young Losses: Generalized Entropies, Margins, and Algorithms". In: *AISTATS*. 2019.

# All you need is a convex functional



**Fenchel-Young losses**<sup>45</sup>: Convex function  $\Omega : \Delta^d \rightarrow \mathbb{R}$  and conjugate

$$\Omega^*(f) = \min_{\alpha \in \Delta^d} \Omega(\alpha) - \langle \alpha, f \rangle \quad \ell_{\Omega}(\alpha, f) = \Omega(\alpha) + \Omega^*(f) - \langle \alpha, f \rangle \geq 0$$

Define link functions between dual and primal

$$\nabla \Omega(\alpha) = \operatorname{argmin}_{f \in \mathbb{R}^d} \ell_{\Omega}(\alpha, f)$$

$$\nabla \Omega^*(f) = \operatorname{argmin}_{\alpha \in \Delta^d} \ell_{\Omega}(\alpha, f)$$

<sup>4</sup>John C. Duchi et al. "Multiclass Classification, Information, Divergence, and Surrogate Risk". In: *Annals of Statistics* (2018).

<sup>5</sup>Mathieu Blondel et al. "Learning Classifiers with Fenchel-Young Losses: Generalized Entropies, Margins, and Algorithms". In: *AISTATS*. 2019.



## Discrete canonical example: Shannon entropy

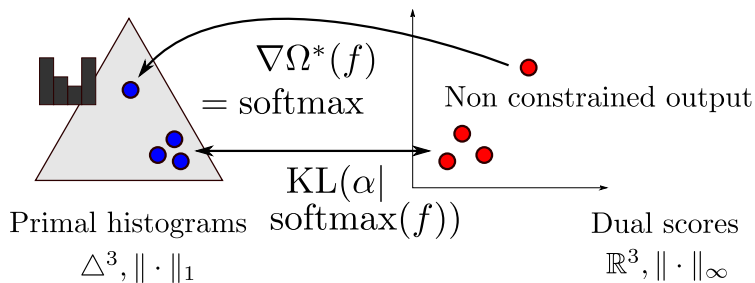
$$\Omega(\alpha) = -H(\alpha) = \sum_{i=1}^d \alpha_i \log \alpha_i$$

$$\Omega^*(f) = \text{logsumexp}(f)$$

# Discrete canonical example: Shannon entropy

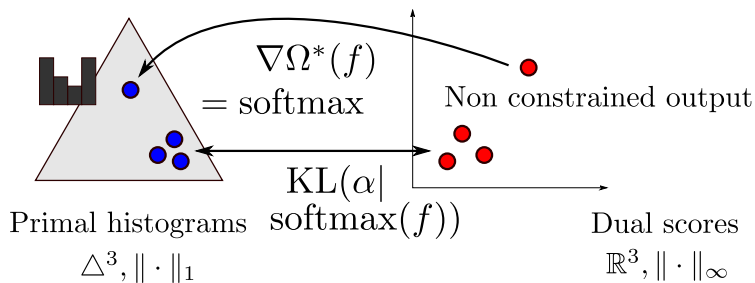
$$\Omega(\alpha) = -H(\alpha) = \sum_{i=1}^d \alpha_i \log \alpha_i$$

$$\Omega^*(f) = \text{logsumexp}(f)$$



# Discrete canonical example: Shannon entropy

$$\Omega(\alpha) = -H(\alpha) = \sum_{i=1}^d \alpha_i \log \alpha_i \quad \Omega^*(f) = \text{logsumexp}(f)$$



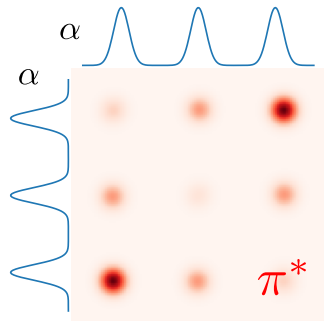
Not defined on continuous distributions, cost-agnostic

- 1 Introduction
- 2 Fenchel-Young losses for distribution spaces
- 3 Geometric softmax from Sinkhorn negentropies**
- 4 Applications

# Sinkhorn entropies from regularized optimal transport

**Self regularized optimal transportation distance:**

$$\Omega_C(\alpha) = -\frac{1}{2} \text{OT}_{C, \varepsilon=2}(\alpha, \alpha) = -\max_{f \in \mathcal{C}(\mathcal{Y})} \langle \alpha, f \rangle - \log \langle \alpha \otimes \alpha, \exp\left(\frac{f \oplus f - C}{2}\right) \rangle$$

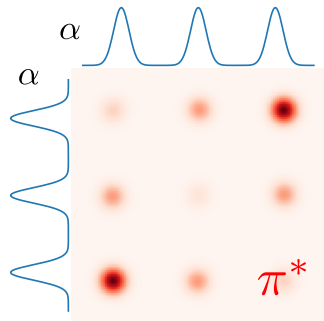


# Sinkhorn entropies from regularized optimal transport

**Self regularized optimal transportation distance:**

$$\Omega_C(\alpha) = -\frac{1}{2} \text{OT}_{C, \varepsilon=2}(\alpha, \alpha) = -\max_{f \in \mathcal{C}(\mathcal{Y})} \langle \alpha, f \rangle - \log \langle \alpha \otimes \alpha, \exp\left(\frac{f \oplus f - C}{2}\right) \rangle$$

**Continuous convex**



# Sinkhorn entropies from regularized optimal transport

**Self regularized optimal transportation distance:**

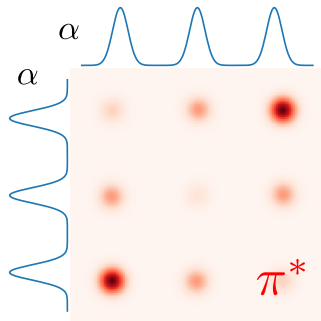
$$\Omega_C(\alpha) = -\frac{1}{2} \text{OT}_{C, \varepsilon=2}(\alpha, \alpha) = -\max_{f \in \mathcal{C}(\mathcal{Y})} \langle \alpha, f \rangle - \log \langle \alpha \otimes \alpha, \exp\left(\frac{f \oplus f - C}{2}\right) \rangle$$

**Continuous convex**

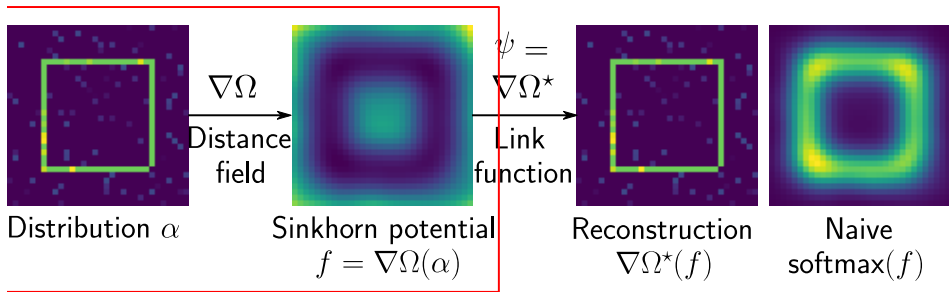
**Special cases**

$\varepsilon \rightarrow \infty$ : MMD autocorrelation

$$C = \begin{pmatrix} 0 & \infty & \dots \\ \infty & 0 & \ddots \\ \vdots & \ddots & \ddots \end{pmatrix} \quad \begin{array}{l} \text{Shannon entropy} \\ \text{Gini index} \end{array}$$



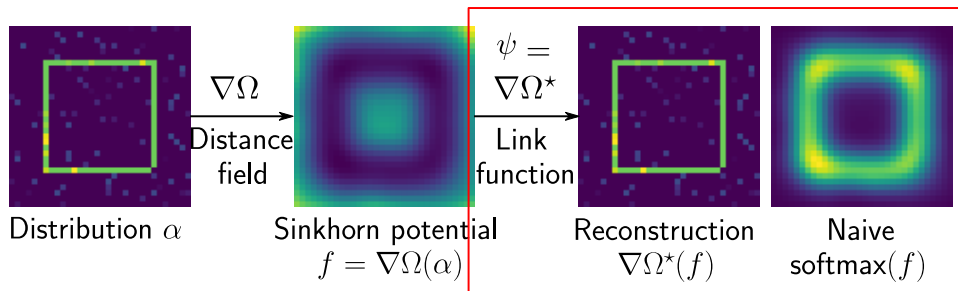
# Dual mapping from Sinkhorn negentropy



Sinkhorn entropy: 
$$\Omega(\alpha) = - \max_{f \in \mathcal{C}(\mathcal{Y})} \langle \alpha, f \rangle - \log \langle \alpha \otimes \alpha, e^{\frac{f \oplus f - C}{2}} \rangle$$



## Returning to primal: geometric softmax



$$\Omega^* = \text{g-logsumexp} : f \mapsto -\log \min_{\alpha \in \mathcal{M}_1^+(\mathcal{Y})} \langle \alpha \otimes \alpha, \exp\left(-\frac{f \oplus f + C}{2}\right) \rangle$$

$\nabla\Omega^*$  = geometric-softmax.

Minimizes a **simple quadratic**.

# Geometric loss construction and computation

**Training with the geometric logistic loss:**

$$\ell_C(\alpha, f) = \text{geometric-LSE}_C(f) + \text{sinkhorn-negentropy}_C(\alpha) - \langle \alpha, f \rangle$$

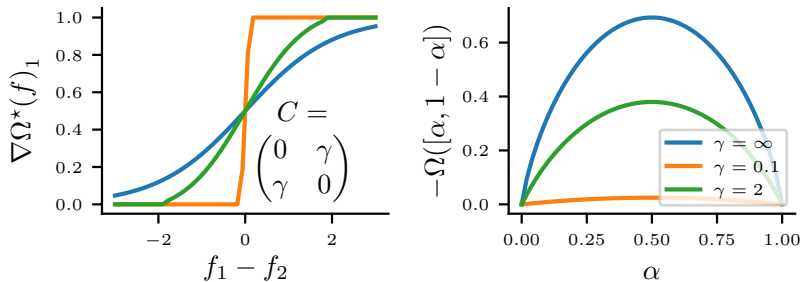
**Tractable discrete  $\mathcal{Y}$ :** use mirror descent/L-BFGS

- 10× as costly as a softmax
- Backpropagation:  $\nabla \Omega^* = \text{geometric-softmax}$

**Continuous case:**

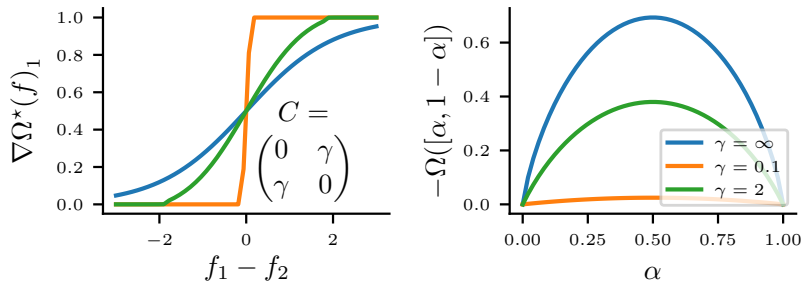
- Frank-Wolfe scheme, adding one Dirac at each iteration

# Properties of the geometric-softmax



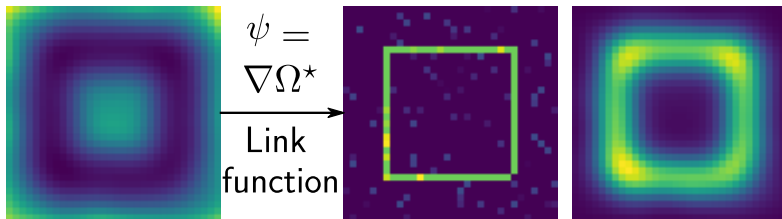
- $\nabla \Omega^*$  returns from Sinkhorn potentials:  $\nabla \Omega^* \circ \nabla \Omega = \text{Id}$

# Properties of the geometric-softmax



- $\nabla\Omega^*$  returns from Sinkhorn potentials:  $\nabla\Omega^* \circ \nabla\Omega = \text{Id}$
- $\nabla\Omega \circ \nabla\Omega^*$  projects  $f$  onto  $\mathcal{F}$  the set of symmetric Sinkhorn potentials
- **Sparse**  $\nabla\Omega^*(f) \leftarrow$  minimization on the simplex

# Properties of the geometric-softmax



$\varepsilon \rightarrow 0$  Mode finding

$\varepsilon \rightarrow \infty$  Positive deconvolution

**Bregman divergence** from Sinkhorn negentropy:

$$D_g(\alpha|\beta) = \Omega(\alpha) - \Omega(\beta) - \langle \nabla \Omega(\alpha), \alpha - \beta \rangle.$$

**Bregman divergence** from Sinkhorn negentropy:

$$D_g(\alpha|\beta) = \Omega(\alpha) - \Omega(\beta) - \langle \nabla \Omega(\alpha), \alpha - \beta \rangle.$$

Sample distribution  $(x_i, \alpha_i)_i \in \mathcal{X} \times \mathcal{M}_1^+(\mathcal{Y})$ .

# Consistent learning with the geometric logistic loss

**Bregman divergence** from Sinkhorn negentropy:

$$D_g(\alpha|\beta) = \Omega(\alpha) - \Omega(\beta) - \langle \nabla \Omega(\alpha), \alpha - \beta \rangle.$$

Sample distribution  $(x_i, \alpha_i)_i \in \mathcal{X} \times \mathcal{M}_1^+(\mathcal{Y})$ .

**Fisher consistency:**

$$\min_{\beta: \mathcal{X} \rightarrow \mathcal{M}_1^+(\mathcal{Y})} \mathbb{E} [D_g(\alpha, \beta(x))] = \min_{g: \mathcal{X} \rightarrow \mathcal{C}(\mathcal{Y})} \mathbb{E} [\ell_g(\alpha, \nabla \Omega^*(g(x)))]$$

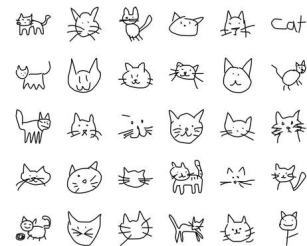


- 1 Introduction
- 2 Fenchel-Young losses for distribution spaces
- 3 Geometric softmax from Sinkhorn negentropies
- 4 Applications**

# Applications: variational auto-encoder

**Goal:** Generate nearly 1D distribution in 2D images

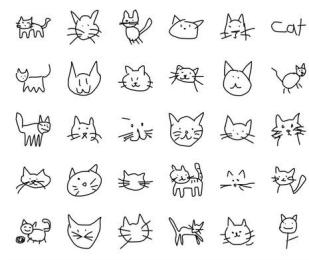
- Dataset: Google Quickdraw
  - Traditional sigmoid activation layer
- ← replaced by **geometric softmax**



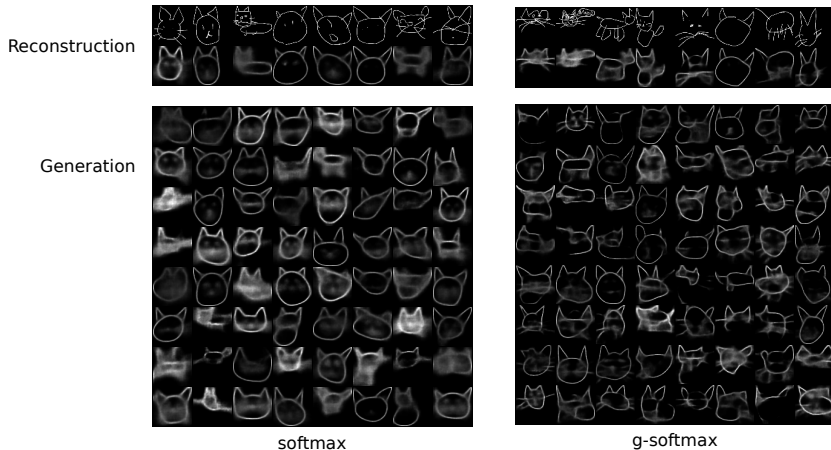
# Applications: variational auto-encoder

**Goal:** Generate nearly 1D distribution in 2D images

- Dataset: Google Quickdraw
- Traditional sigmoid activation layer
- ← replaced by **geometric softmax**
- Deconvolutional effect
- Cost-informed non-linearity



# Applications: variational auto-encoders



Better defined generated images

## Geometric softmax:

- New loss and projector onto output probabilities
- Discrete/continuous, aware of a cost between outputs
- Fenchel duality in Banach spaces + regularized optimal transport
- Application in VAE and ordinal regression

# Conclusion

## Geometric softmax:

- New loss and projector onto output probabilities
- Discrete/continuous, aware of a cost between outputs
- Fenchel duality in Banach spaces + regularized optimal transport
- Application in VAE and ordinal regression

## Future directions:

- How to improve computation methods (continuous FW)
- Geometric logistic loss in super resolution<sup>6</sup>

---

<sup>6</sup>Nicholas Boyd et al. "DeepLoco: Fast 3D localization microscopy using neural networks". In: *BioRxiv* (2018), p. 267096.



Mathieu Blondel



Gabriel Peyré

Poster # 179