

Google DeepMind

Differentiable and Sparse Top-k: a Convex Analysis Perspective

Mathieu Blondel

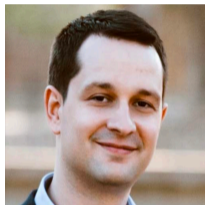
July 28th, ICML 2023



Michaël Sander



Joan Puigcerver



Josip Djolonga



Gabriel Peyré



Mathieu Blondel

Motivation for the research

- The top-k operator is increasingly used as a **building block** in neural networks (top-k classification, mixtures of expert, weight pruning)
- However, it is a **discontinuous** operation, making it difficult to use in end-to-end trainable networks
- A crucial property of the top-k is its **sparsity** but many existing differentiable top-k relaxations are **dense**
- **Smooth optimization** is known to enjoy **faster** convergence rates
- However, sparsity is crucial in certain applications as a **selection mechanism**: mixtures of experts, weight pruning

Related work

A large body of work on relaxations of sorting, ranking and top-k...

- Using **unimodal row-stochastic matrices** (Grover et al, 2019; Prillo and Eisenschlos, 2020)
- Using **optimal transport** (Cuturi et al, 2019)
- Using the **permutahedron** (Blondel et al, 2020)
- Using **perturbations** (Berthet et al, 2020)
- Using **sorting networks** (Petersen et al, 2021)

Contributions

- A **general** top-k framework, including top-k **in magnitude**
- Differentiable **and** sparse relaxations thanks to p -norm regularization
- Reduction to **isotonic optimization**, for computation and differentiation
- **GPU/TPU-friendly** algorithm based on Dykstra's algorithm
- Applications to top-k classification, mixtures of experts, weight pruning

1 Top-k mask

Top-k mask operator

Bit-encoding of the top-k indices (“k-hot encoding”)

$$[\mathbf{topkmask}(x)]_i := \begin{cases} 1, & \text{if } [\mathbf{rank}(x)]_i \leq k \\ 0, & \text{otherwise.} \end{cases} \in \{0, 1\}^n$$

$$x = (1.7, 3.2, -2.4)$$

$$\mathbf{top1mask}(x) = (0, 1, 0)$$

$$\mathbf{top2mask}(x) = (1, 1, 0)$$

$$\mathbf{top3mask}(x) = (1, 1, 1)$$

Discontinuous, piecewise constant with null derivatives

Top-k operator

Sparse vector containing the top-k values

$$\mathbf{topk}(x) := x \cdot \mathbf{topkmask}(x) \in \mathbb{R}^n$$

$$x = (1.7, 3.2, -2.4)$$

$$\mathbf{top1}(x) = (0.0, 3.2, 0)$$

$$\mathbf{top2}(x) = (1.7, 3.2, 0)$$

$$\mathbf{top3}(x) = (1.7, 3.2, -2.4)$$

Discontinuous, piecewise affine with constant derivatives

Regularized top-k mask: overview of the approach

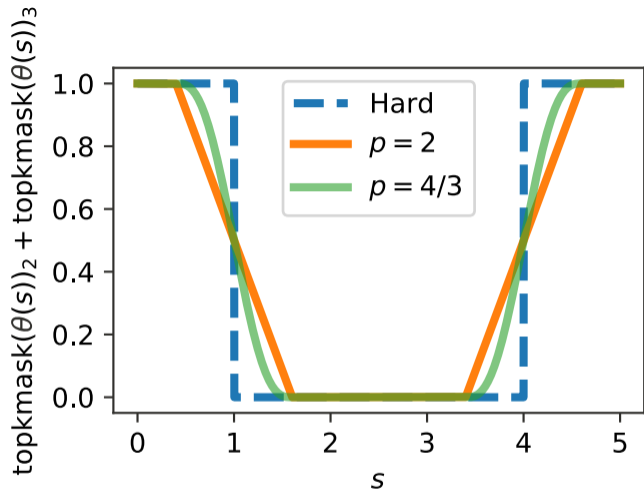
- Rewrite top-k mask as a linear program solution

$$\text{topkmask}(x) = y(x) := \underset{y \in \mathcal{C}}{\operatorname{argmax}} \langle x, y \rangle$$

- Add regularization R

$$\text{topkmask}_R(x) = y_R(x) := \underset{y \in \mathcal{C}}{\operatorname{argmax}} \langle x, y \rangle - R(y)$$

- Use a reduction to isotonic optimization to easily compute and differentiate $\text{topkmask}_R(x)$



$$\theta(s) = (3, 1, -1 + s, s) \in \mathbb{R}^4$$

$$k = 2$$

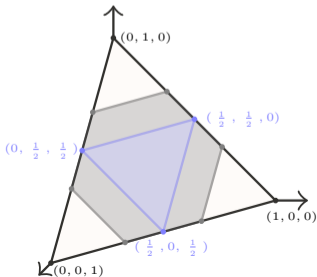
Top-k mask as a linear program

- With $\mathcal{C} = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y} \in [0, 1]^n, \mathbf{y}^\top \mathbf{1} = k\}$, we get

$$\text{topkmask}(\mathbf{x}) = \mathbf{y}(\mathbf{x}) = \underset{\mathbf{y} \in \mathcal{C}}{\text{argmax}} \langle \mathbf{x}, \mathbf{y} \rangle$$

- The vertices of \mathcal{C} are all possible bit encodings of cardinality k
- Relation with the **capped probability simplex**

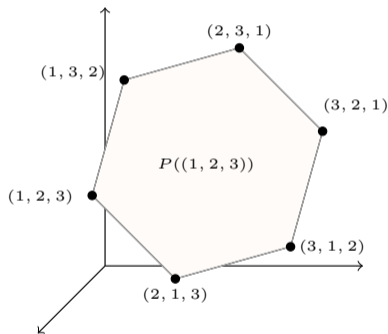
$$\mathcal{C}/k = \left\{ \mathbf{y} \in \mathbb{R}^n : \mathbf{y} \in [0, 1/k]^n, \mathbf{y}^\top \mathbf{1} = 1 \right\}$$



Relation with the permutahedron

- The convex hull of all permutations of w

$$P(w) := \text{conv}(\{(w_{\sigma_1}, \dots, w_{\sigma_n}) : \sigma \in \Sigma\})$$



- With $w = \mathbf{1}_k := (\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{n-k})$, we get

$$P(w) = \mathcal{C} = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y} \in [0, 1]^n, \mathbf{y}^\top \mathbf{1} = k\}$$

Top-k mask: value function and its conjugate

- Value function: support function of \mathcal{C}

$$f(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{C}} \langle \mathbf{x}, \mathbf{y} \rangle = \text{topksum}(\mathbf{x}) := \sum_{i=1}^k x_{\sigma_i} = \langle \mathbf{x}_{\sigma}, \mathbf{1}_k \rangle$$

where $\sigma = \text{argsort}(\mathbf{x}) \iff x_{\sigma_1} \geq \dots \geq x_{\sigma_n}$ and $\mathbf{x}_{\sigma} := (x_{\sigma_1}, \dots, x_{\sigma_n})$

- Conjugate: indicator function of \mathcal{C}

$$f^*(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbb{R}^n} \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}) = \delta_{\mathcal{C}}(\mathbf{y}) := \begin{cases} 0, & \text{if } \mathbf{y} \in \mathcal{C} \\ \infty, & \text{if } \mathbf{y} \notin \mathcal{C} \end{cases}$$

Regularized version

- The regularized version

$$\text{topmask}_R(x) = \mathbf{y}_R(x) := \mathbf{y}^*$$

is defined using the **dual** solution

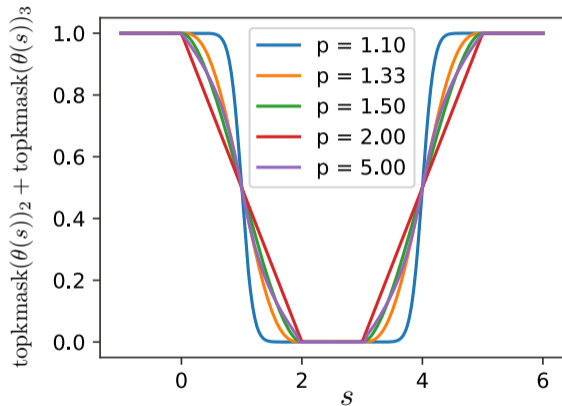
$$\begin{aligned}\mathbf{y}^* &= \underset{\mathbf{y} \in \mathcal{C}}{\operatorname{argmax}} \langle x, \mathbf{y} \rangle - R(\mathbf{y}) \\ &= \underset{\mathbf{y} \in \mathbb{R}^n}{\operatorname{argmax}} \langle x, \mathbf{y} \rangle - f^*(\mathbf{y}) - R(\mathbf{y})\end{aligned}$$

- Equivalently, if we define the **primal** solution (infimal convolution)

$$\mathbf{u}^* = \underset{\mathbf{u} \in \mathbb{R}^n}{\operatorname{argmin}} R^*(x - \mathbf{u}) + f(\mathbf{u})$$

then $\mathbf{y}^* = \nabla R^*(x - \mathbf{u}^*)$

Regularized version



$$R(\mathbf{y}) = \frac{1}{p} \|\mathbf{y}\|_p^p = \frac{1}{p} \sum_{i=1}^n |y_i|^p$$

Computing the regularized version

- Recall that the primal solution is

$$\begin{aligned} \mathbf{u}^* &= \underset{\mathbf{u} \in \mathbb{R}^n}{\operatorname{argmin}} R^*(\mathbf{x} - \mathbf{u}) + f(\mathbf{u}) \\ &= \underset{\mathbf{u} \in \mathbb{R}^n}{\operatorname{argmin}} R^*(\mathbf{x} - \mathbf{u}) + \langle \mathbf{u}_{\pi(\mathbf{u})}, \mathbf{1}_k \rangle \end{aligned}$$

where $\pi(\mathbf{u}) = \operatorname{argsort}(\mathbf{u})$

- Reduction to isotonic optimization

$$\begin{aligned} \mathbf{u}_\sigma^* &= \underset{v_1 \geq \dots \geq v_n}{\operatorname{argmin}} R^*(\mathbf{x}_\sigma - \mathbf{v}) + f(\mathbf{v}) \\ &= \underset{v_1 \geq \dots \geq v_n}{\operatorname{argmin}} R^*(\mathbf{x}_\sigma - \mathbf{v}) + \langle \mathbf{v}, \mathbf{1}_k \rangle \end{aligned}$$

where $\sigma = \operatorname{argsort}(\mathbf{x})$

- Differentiation available in closed form (implicit diff not needed) given \mathbf{v}^*

Pool Adjacent Violators (PAV)

$$\operatorname{argmin}_{v_1 \geq \dots \geq v_n} \sum_{i=1}^n h_i(v_i)$$

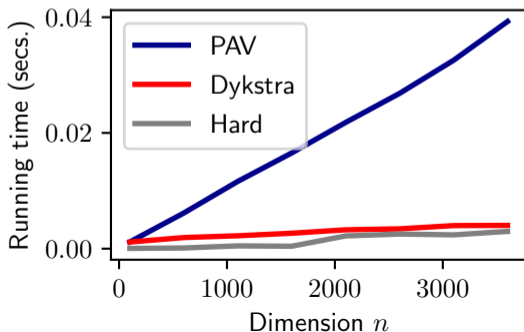
- Partitions the set $[n]$ into disjoint sets (B_1, \dots, B_m) , starting from $m = n$ and $B_i = \{i\}$
- Merges these sets until the isotonic condition is met
- Needs to be able to solve $\operatorname{argmin}_{\gamma \in \mathbb{R}} \sum_{i \in B_j} h_i(\gamma)$ in constant time to get $O(n)$ total complexity
 - p -norm regularization case: we need to find the root of a polynomial (easy when $p = 2$ or $p = 4/3$)

Using Dykstra's algorithm

- Key idea: alternate projections between C_1 and C_2

$$\{\mathbf{v} \in \mathbb{R}^n : v_1 \geq \dots \geq v_n\} = \underbrace{\{\mathbf{v} \in \mathbb{R}^n : v_1 \geq v_2, v_3 \geq v_4, \dots\}}_{C_1} \cap \underbrace{\{\mathbf{v} \in \mathbb{R}^n : v_2 \geq v_3, v_4 \geq v_5, \dots\}}_{C_2}$$

- Huge speedup on TPU



2 Top-k in magnitude

Top-k in magnitude operator

Same as top-k operator but selects elements with largest **absolute value**

$$\mathbf{topkmag}(x) := x \cdot \mathbf{topkmask}(|x|)$$

$$x = (1.7, 3.2, -2.4)$$

$$\mathbf{top1mag}(x) = (0.0, 3.2, 0)$$

$$\mathbf{top2mag}(x) = (0, 3.2, -2.4)$$

$$\mathbf{top3mag}(x) = (1.7, 3.2, -2.4)$$

Top-k in magnitude as a gradient

- We introduce a **nonlinearity** $\varphi(\mathbf{x}) = (\phi(x_1), \dots, \phi(x_n))$

$$f_\varphi(\mathbf{x}) := f(\varphi(\mathbf{x})) = \max_{\mathbf{y} \in \mathcal{C}} \langle \varphi(\mathbf{x}), \mathbf{y} \rangle$$

- With $\phi(x) = \frac{1}{2}x^2$, we have

$$\nabla f_\varphi(\mathbf{x}) = \text{topkmag}(\mathbf{x})$$

- With $\phi(x) = x$, we have

$$\nabla f_\varphi(\mathbf{x}) = \nabla f(\mathbf{x}) = \text{topkmask}(\mathbf{x})$$

Regularized version

$$\text{topkmag}_R(\mathbf{x}) := \mathbf{y}^* = \nabla R^*(\mathbf{x} - \mathbf{u}^*)$$

where we defined the **dual** solution

$$\mathbf{y}^* := \underset{\mathbf{y} \in \mathbb{R}^n}{\text{argmax}} \langle \mathbf{x}, \mathbf{y} \rangle - f_\varphi^*(\mathbf{y}) - R(\mathbf{y})$$

and the **primal** solution (inf-convolution)

$$\mathbf{u}^* = \underset{\mathbf{u} \in \mathbb{R}^n}{\text{argmin}} R^*(\mathbf{x} - \mathbf{u}) + f_\varphi(\mathbf{u})$$

Conjugate

- For $\varphi(x) = x$: **indicator function of \mathcal{C}**

$$f_{\varphi}^*(\mathbf{y}) = f^*(\mathbf{y}) = \delta_{\mathcal{C}}(\mathbf{y})$$

- For $\varphi(x) = \frac{1}{2}x^2$: **squared k-support norm**

$$f_{\varphi}^*(\mathbf{y}) = \frac{1}{2} \min_{z \in \mathcal{C}} \sum_{i=1}^n \frac{y_i^2}{z_i}$$

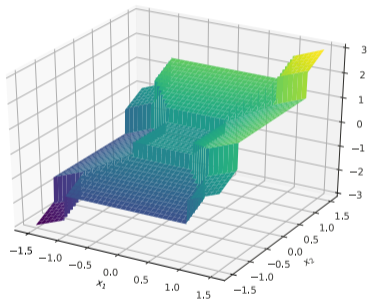
- For general φ : **minimum distance to \mathcal{C}**

$$f_{\varphi}^*(\mathbf{y}) = \min_{z \in \mathcal{C}} D_{\phi^*}(\mathbf{y}, z)$$

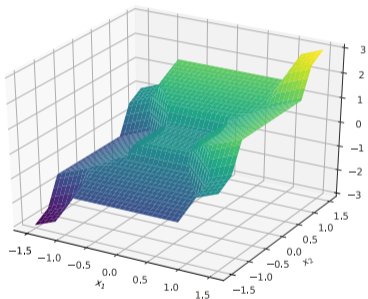
where we defined the **f-divergence**

$$D_f(\mathbf{y}, z) := \sum_{i=1}^n z_i f(y_i/z_i)$$

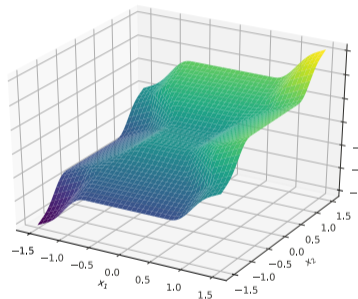
Regularized version



Hard



$p = 2$



$p = 4/3$

Computing the regularized version

- Primal solution

$$\mathbf{u}^* = \underset{\mathbf{u} \in \mathbb{R}^n}{\operatorname{argmin}} R^*(\mathbf{x} - \mathbf{u}) + f_\varphi(\mathbf{u})$$

- Reduction to isotonic optimization

$$\begin{aligned} \mathbf{u}_\sigma^* &= \underset{v_1 \geq \dots \geq v_n \geq 0}{\operatorname{argmin}} R^*(\mathbf{x}_\sigma - \mathbf{v}) + f_\varphi(\mathbf{v}) \\ &= \underset{v_1 \geq \dots \geq v_n \geq 0}{\operatorname{argmin}} R^*(\mathbf{x}_\sigma - \mathbf{v}) + f(\varphi(\mathbf{v})) \\ &= \underset{v_1 \geq \dots \geq v_n \geq 0}{\operatorname{argmin}} R^*(\mathbf{x}_\sigma - \mathbf{v}) + \langle \varphi(\mathbf{v}), \mathbf{1}_k \rangle \end{aligned}$$

where $\sigma = \operatorname{argsort}(|\mathbf{x}|)$ and assuming $\varphi(\mathbf{x}) = \varphi(-\mathbf{x})$

Nonconvex viewpoint: connection with the ℓ_0 pseudo-norm

- We have

$$f_\varphi(\mathbf{x}) = \max_{\mathbf{y} \in S_k} \langle \mathbf{x}, \mathbf{y} \rangle - \sum_{i=1}^n \phi(y_i)$$

where

$$\varphi(x) = (\phi(x_1), \dots, \phi(x_n)) \quad \text{and} \quad S_k := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y}\|_0 \leq k\}$$

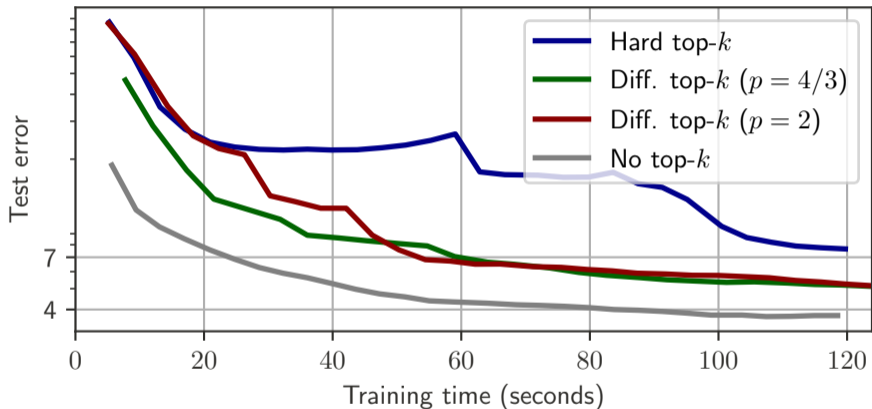
- $f_\varphi^*(\mathbf{y})$ is the **convex envelope** of

$$\mathbf{y} \mapsto \sum_{i=1}^n \phi^*(y_i) + \delta_{S_k}(\mathbf{y})$$

With $\phi(x) = \frac{1}{2}x^2$, $f_\varphi^*(\mathbf{y})$ is the **squared k-support norm**

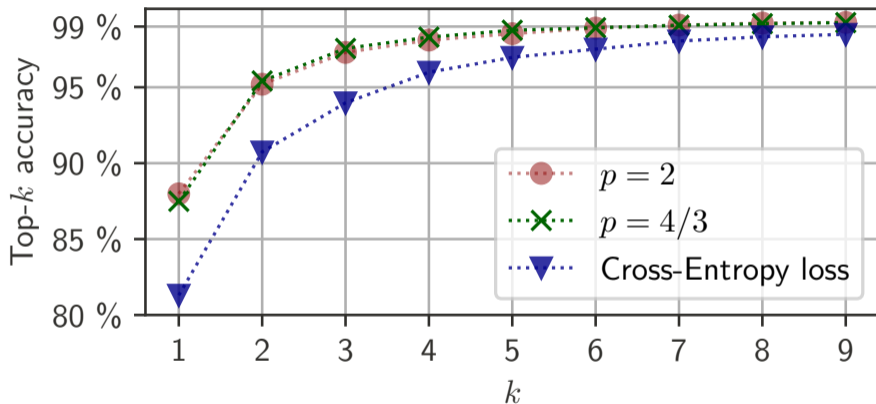
3 Applications

Weight pruning (multilayer perceptron, MNIST)



$$W_i \leftarrow \text{topkmag}(W_i), i \in \{2, 3\}$$
$$W_3 \sigma(W_2 \sigma(W_1 \mathbf{a} + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_3$$

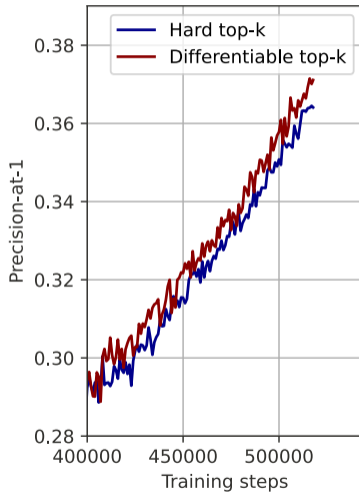
Top- k classification (vision transformer, CIFAR 100)



$$\ell_R(\theta, t) := [\max_{y \in \mathcal{C}} \langle \theta, y \rangle - R(y)] - \langle \theta, t \rangle$$

θ : logits, t : target

Sparse Mixture of Vision Transformers



JFT-300M dataset (305 million images)

Sparsity-constrained OT

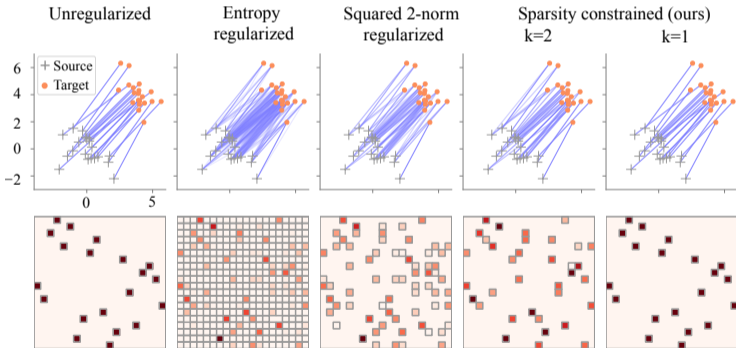


Sparsity-constrained OT.

Liu, Puigcerver, Blondel.

ICLR, 2023.

$$\min_{T \in \mathcal{U}(a,b)} \langle T, C \rangle + \sum_{j=1}^n f_{\varphi}^*(\mathbf{t}_j)$$



Google DeepMind

Thank you!