

STOCHASTIC SMOOTHING OF THE TOP-K CALIBRATED HINGE LOSS FOR DEEP IMBALANCED CLASSIFICATION

Camille Garcin^{1,2}

Maximilien Servajean^{3,4}

Alexis Joly²

Joseph Salmon^{1,5}

¹ IMAG, Univ Montpellier, CNRS

² Inria, LIRMM, Univ Montpellier, CNRS

³ LIRMM, Univ Montpellier, CNRS

⁴ AMIS, Univ Paul Valéry

⁵ Institut Universitaire de France (IUF)



UNIVERSITÉ DE
MONTPELLIER



Inria

INTER-CLASS AMBIGUITY

DIFFERENT LABELS BUT SIMILAR IMAGES



Cirsium rivulare



Chaerophyllum aromaticum



Conostomium kenysense



Adenostyles leucophylla



Sedum montanum



Cirsium tuberosum



Chaerophyllum temulum



Conostomium quadrangulare



Adenostyles alliariae



Sedum rupestre

Some species are visually similar

Necessary to return several classes

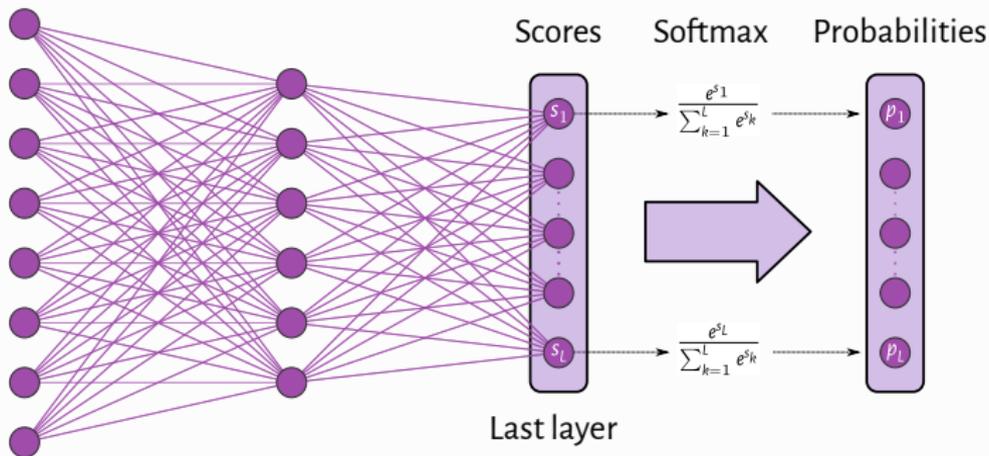


Possible solution: return the K "most likely" classes for all images

- ▶ Pros for a small K :
model output must remain informative
- ▶ Pros for a large K :
ensure the true class lies in the K returned classes

Choice of K :

- ▶ task-dependant, often $K = 3, 5, \dots$ or even larger for challenging tasks
- ▶ considered fixed by the user for the talk (not tuned)



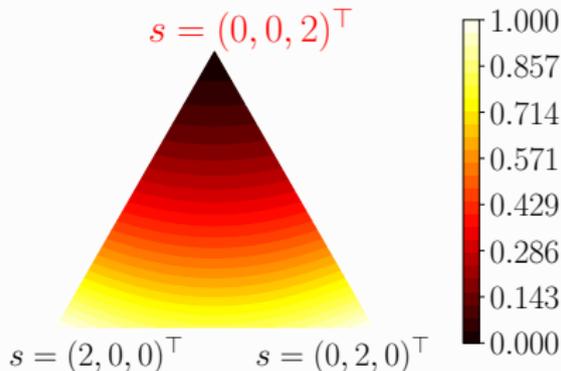
- ▶ L : number of **classes**, $[L] := \{1, \dots, L\}$, label space
- ▶ $K \in [L]$ is a fixed parameter used for top- K
- ▶ From an image, get a score vector $\mathbf{s} = (s_1, \dots, s_L)^T \in \mathbb{R}^L$ (aka logits)
- ▶ s_k : score for class k
- ▶ Prediction: output the K classes with the K highest scores

- ▶ Training: cross-entropy (CE) loss + Stochastic Gradient Descent (SGD)

- ▶ $l_{\text{CE}}(\mathbf{s}, y) = -\log\left(\frac{e^{s_y}}{\sum_{k \in [L]} e^{s_k}}\right)$

Example : $L = 3, K = 2, y = 3$

(Normalized) level set of $\mathbf{s} \mapsto l_{\text{CE}}(\mathbf{s}, y)$:



- ▶ Not designed to optimize top- K accuracy
- ▶ Can we do better than cross entropy ?



For a score vector $\mathbf{s} \in \mathbb{R}^L$:

Definition

$\text{top}_K : \mathbf{s} \mapsto s_{(K)}$ (K-th largest score)

Properties

- ▶ $\nabla \text{top}_K(\mathbf{s}) = \arg \text{top}_K(\mathbf{s}) \in \mathbb{R}^L$:
vector with a single 1 at the K-th largest coordinate of \mathbf{s} , 0 o.w.

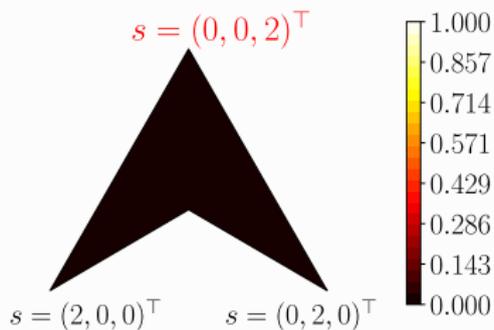
$$\text{For } \mathbf{s} = \begin{bmatrix} 4.0 \\ -1.5 \\ 2.5 \\ 1.0 \end{bmatrix}, \quad \text{top}_2(\mathbf{s}) = 2.5 \quad \nabla \text{top}_2(\mathbf{s}) := \arg \text{top}_2(\mathbf{s}) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

⁽¹⁾ F. Yang and S. Koyejo (2020). "On the consistency of top-k surrogate losses". In: *ICML*. vol. 119, pp. 10727–10735.

Objective: minimize top-K error (0/1 loss):

$$\ell^K(\mathbf{s}, \mathbf{y}) = \mathbb{1}_{\{\text{top}_K(\mathbf{s}) > s_y\}}$$

Problem: piecewise constant function w.r.t. \mathbf{s} , hard to optimize!!!

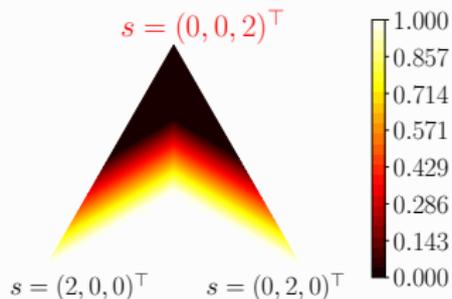


(Normalized) Level sets of $\mathbf{s} \mapsto \ell^K(\mathbf{s}, \mathbf{y})$, $L = 3$, $K = 2$, $y = 3$.



A top-K hinge-loss that is top-K calibrated:

$$\ell_{\text{Cal. Hinge}}^k(\mathbf{s}, y) = (1 + \text{top}_{K+1}(\mathbf{s}) - s_y)_+$$



Better theoretical properties, but still fails with deep learning (more later)

Problem: $\mathbf{s} \rightarrow \text{top}_K(\mathbf{s})$ non-smooth and sparse gradient



Motivation: top_K is a non-smooth, function, smooth it!

- ▶ smoothing parameter $\epsilon > 0$
- ▶ score $\mathbf{s} \in \mathbb{R}^L$

Definition

The ϵ -smoothed version of top_K :

$$\text{top}_{K,\epsilon}(\mathbf{s}) \triangleq \mathbb{E}_Z[\text{top}_K(\mathbf{s} + \epsilon Z)]$$

Z : standard normal random vector, $Z \sim \mathcal{N}(0, \text{Id}_L)$

⁽³⁾Q. Berthet et al. (2020). "Learning with differentiable perturbed optimizers". In: *NeurIPS*.



Proposition

For a smoothing parameter $\epsilon > 0$,

- ▶ $\text{top}_{K,\epsilon}$ is $\frac{4\sqrt{KL}}{\epsilon}$ -smooth.
- ▶ For any $\mathbf{s} \in \mathbb{R}^L$, $|\text{top}_{K,\epsilon}(\mathbf{s}) - \text{top}_K(\mathbf{s})| \leq \epsilon \cdot C_{K,L}$, where $C_{K,L} = K\sqrt{2 \log L}$.
- ▶ The gradient of $\text{top}_{K,\epsilon}$ reads:

$$\nabla_{\mathbf{s}} \text{top}_{K,\epsilon}(\mathbf{s}) = \mathbb{E}[\arg \text{top}_K(\mathbf{s} + \epsilon \mathbf{Z})]$$

- ▶ From non-smooth to smooth function with simple stochastic perturbation
- ▶ When $\epsilon \rightarrow 0$, recover the original function



Solution: Draw B noise vectors Z_1, \dots, Z_B , with $Z_b \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \text{Id}_L)$ for $b \in [B]$.

$$\text{top}_{K,\epsilon}(\mathbf{s}) = \mathbb{E}_Z[\text{top}_K(\mathbf{s} + \epsilon Z)]$$

Monte Carlo estimation :

$$\widehat{\text{top}}_{K,\epsilon,B}(\mathbf{s}) = \frac{1}{B} \sum_{b=1}^B \text{top}_K(\mathbf{s} + \epsilon Z_b)$$

Easy implementation with deep learning libraries *e.g.*, Pytorch, Tensorflow



Solution: Draw B noise vectors Z_1, \dots, Z_B , with $Z_b \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \text{Id}_L)$ for $b \in [B]$.

$$\nabla_{\mathbf{s}} \text{top}_{K, \epsilon}(\mathbf{s}) = \mathbb{E}[\arg \text{top}_K(\mathbf{s} + \epsilon Z)]$$

Monte Carlo estimation :

$$\widehat{\nabla}_{\text{top}_{K, \epsilon, B}}(\mathbf{s}) = \frac{1}{B} \sum_{b=1}^B \arg \text{top}_K(\mathbf{s} + \epsilon Z_b)$$

Pytorch implementation available at:

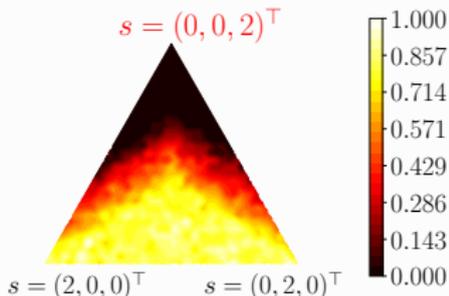
<https://github.com/garcinc/noised-topk>

Modification: use larger margins for classes with few examples⁽⁴⁾:

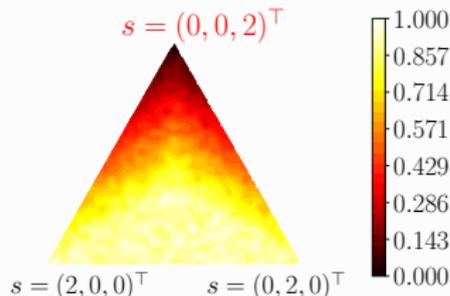
$$\ell_{\text{Noised Imbal.}}^{K, \epsilon, B, m_y}(\mathbf{s}, y) = (m_y + \widehat{\text{top}}_{K+1, \epsilon, B}(\mathbf{s}) - s_y)_+ \quad (1)$$

Set $m_y = C/n_y^{1/4}$, with n_y the number of samples in the training set with class y , and C a hyperparameter to be tuned on a validation set.

Intuition: add more emphasis on rarely seen examples



(a) $\ell_{\text{Noised bal.}}^{K, 0.3, 30}$



(b) $\ell_{\text{Noised bal.}}^{K, 1, 30}$

⁽⁴⁾ K. Cao et al. (2019). "Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss". In: *NeurIPS*. vol. 32, pp. 1565–1576.



K	ℓ_{CE}	$\ell_{\text{Smoothed Hinge}}^{K,0.1}$	$\ell_{\text{Noised bal.}}^{K,1.0,5}$	focal ($\gamma = 2.0$)	$\ell_{\text{LDAM}}^{\max m_y=0.2}$	$\ell_{\text{Noised imbal.}}^{K,0.01,5,\max m_y=0.2}$
1	36.3±0.3	35.7±0.2	35.8±0.3	37.6±0.3	40.6±0.1	42.4±0.3
3	58.8±0.4	50.3±0.2	58.7±0.4	60.4±0.3	63.3±0.3	64.9±0.4
5	68.7±0.2	50.9±0.3	66.4±0.5	69.7±0.2	71.9±0.3	73.2±0.5

Macro-average test top-K accuracy on Pl@ntNet-300K, ResNet-50.

- ▶ Imbalanced losses fare better than balanced losses
- ▶ Class-wise margin is effective compared to constant margin
- ▶ $\ell_{\text{Noised imbal.}}^{K,\epsilon,B,m_y}$ outperforms other losses on Pl@ntNet-300K
- ▶ Many more experiments in the paper !



Thank you for your attention !

Contact: camille.garcin@umontpellier.fr

-  Berthet, Q. et al. (2020). “Learning with differentiable perturbed optimizers”. In: *NeurIPS*.
-  Cao, K. et al. (2019). “Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss”. In: *NeurIPS*. Vol. 32, pp. 1565–1576.
-  Yang, F. and S. Koyejo (2020). “On the consistency of top-k surrogate losses”. In: *ICML*. Vol. 119, pp. 10727–10735.