

Input Similarity from the Neural Network Perspective

Guillaume Charpiat¹

Nicolas Girard²

Loris Felardos¹

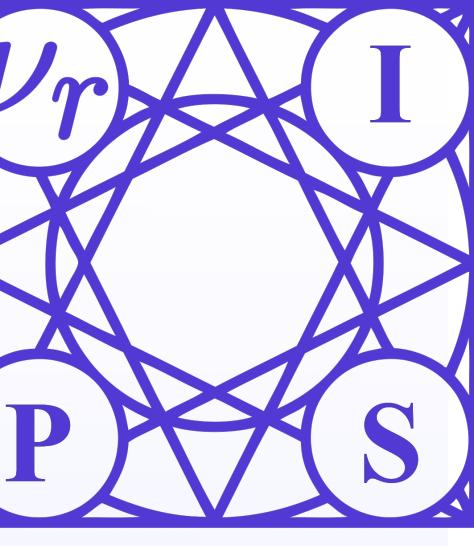
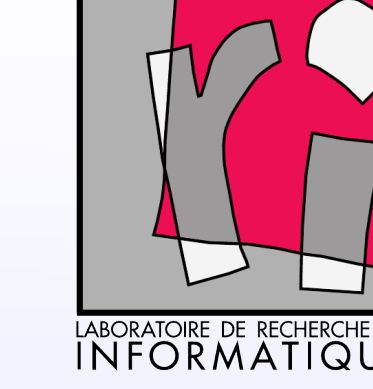
Yuliya Tarabalka^{2,3}

¹ TAU team, INRIA Saclay, LRI, Université Paris-Sud.

² Université Côte d’Azur, Inria.

³ LuxCarta Technology.

^{1,2,3} <firstname>.<lastname>@inria.fr



Introduction

- **Goal:** define and estimate the *similarity* of inputs, as perceived by the neural network
- **Motivation 1:** strong auto-denoising phenomenon in a multimodal image registration task (cf part IV)
 - ⇒ accuracy far better than label noise!
 - ⇒ analyze noise averaging effect over labels of *similar* examples
- **Motivation 2:** better understand neural network decisions
 - ⇒ display examples considered as similar *by the network*

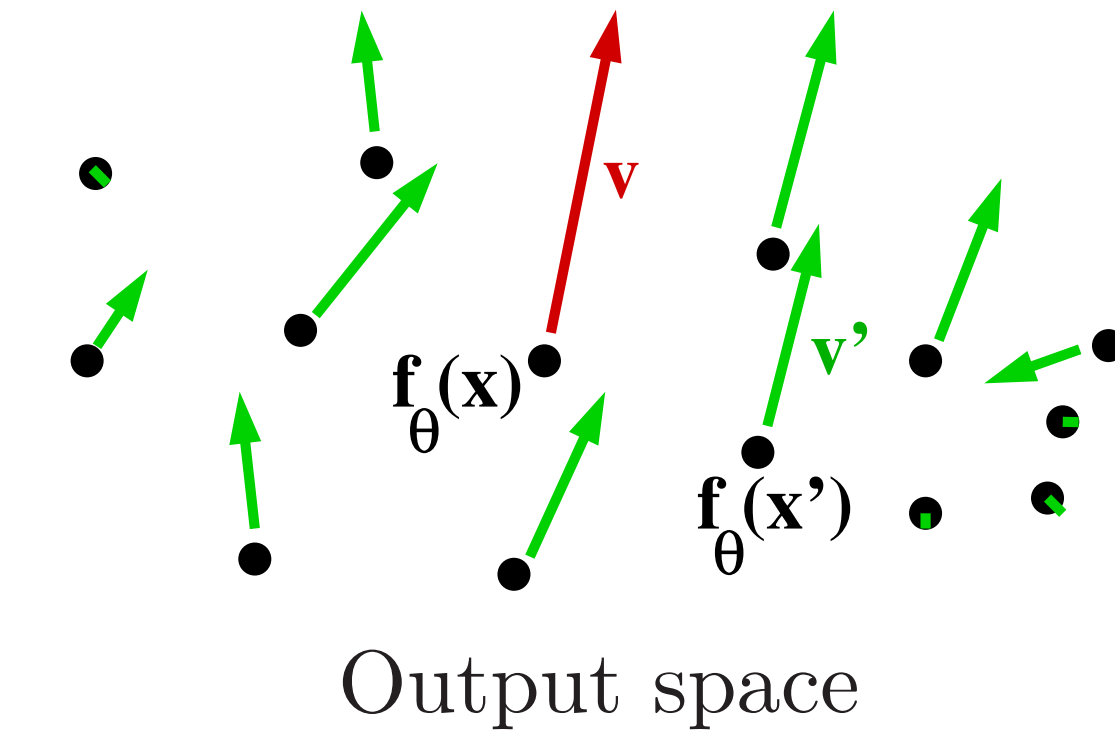
I - Building the similarity measure

Notations:

f_θ : trained neural network
 θ : network parameters
 \mathbf{x}, \mathbf{x}' : possible inputs

Similarity definition:

influence of \mathbf{x} over \mathbf{x}' , quantified as:
 how much an additional training step to push $f_\theta(\mathbf{x})$ in a certain direction would change the output for \mathbf{x}' as well.



- If \mathbf{x} and \mathbf{x}' are very different (for the NN): changing $f_\theta(\mathbf{x})$ will barely affect $f_\theta(\mathbf{x}')$
- If \mathbf{x} and \mathbf{x}' are very similar: changing $f_\theta(\mathbf{x})$ will greatly affect $f_\theta(\mathbf{x}')$

- Changing $f_\theta(\mathbf{x})$ by a small quantity ε means updating θ by $\delta\theta = \varepsilon \frac{\nabla_\theta f_\theta(\mathbf{x})}{\|\nabla_\theta f_\theta(\mathbf{x})\|^2}$
- After update, new values for \mathbf{x} and \mathbf{x}' :

$$f_{\theta+\delta\theta}(\mathbf{x}) = f_\theta(\mathbf{x}) + \nabla_\theta f_\theta(\mathbf{x}) \cdot \delta\theta + O(\|\delta\theta\|^2) = f_\theta(\mathbf{x}) + \varepsilon + O(\varepsilon^2)$$

$$f_{\theta+\delta\theta}(\mathbf{x}') = f_\theta(\mathbf{x}') + \nabla_\theta f_\theta(\mathbf{x}') \cdot \delta\theta + O(\|\delta\theta\|^2) = f_\theta(\mathbf{x}') + \varepsilon \underbrace{\frac{\nabla_\theta f_\theta(\mathbf{x}') \cdot \nabla_\theta f_\theta(\mathbf{x})}{\|\nabla_\theta f_\theta(\mathbf{x})\|^2}}_{\text{influence of } \mathbf{x} \text{ over } \mathbf{x}' } + O(\varepsilon^2)$$

- Symmetric kernel bounded in $[-1, 1]$: $k_\theta^C(\mathbf{x}, \mathbf{x}') = \frac{\nabla_\theta f_\theta(\mathbf{x})}{\|\nabla_\theta f_\theta(\mathbf{x})\|} \cdot \frac{\nabla_\theta f_\theta(\mathbf{x}')}{\|\nabla_\theta f_\theta(\mathbf{x}')\|}$

Properties for vanilla neural networks

Theorem. For any real-valued network f_θ without parameter sharing, if $\nabla_\theta f_\theta(\mathbf{x}) = \nabla_\theta f_\theta(\mathbf{x}')$ for two inputs \mathbf{x}, \mathbf{x}' , then all useful activities computed when processing \mathbf{x} are equal to the ones obtained when processing \mathbf{x}' .

Rewriting: $k_\theta(\mathbf{x}, \mathbf{x}') = \sum_{\text{activities } i} \lambda_i(\mathbf{x}, \mathbf{x}') a_i(\mathbf{x}) a_i(\mathbf{x}')$ where $\lambda_i(\mathbf{x}, \mathbf{x}') = \sum_{\text{neuron } j \text{ using } a_i} \frac{df_\theta(\mathbf{x})}{db_j} \frac{df_\theta(\mathbf{x}')}{db_j}$

⇒ data-dependent importance weights, vs. $\lambda_i(\mathbf{x}, \mathbf{x}') = \lambda_{\text{layer}(i)}$ for the *perceptual loss*

General case (parameter-sharing networks)

properties above do not hold anymore, as invariances (s.t. as translation) are incorporated in the similarity

Higher output dimension

- $f_\theta(\mathbf{x}) = (f_\theta^i(\mathbf{x}))_{i \in [1, d]} \in \mathbb{R}^d$ with $d > 1$
- $K_\theta(\mathbf{x}', \mathbf{x})$ the $d \times d$ kernel matrix defined by $K_\theta^{ij}(\mathbf{x}', \mathbf{x}) = \nabla_\theta f_\theta^i(\mathbf{x}') \cdot \nabla_\theta f_\theta^j(\mathbf{x})$
- Unitless symmetrized normalized kernel: $K_\theta^C(\mathbf{x}, \mathbf{x}') = K_\theta(\mathbf{x}, \mathbf{x})^{-1/2} K_\theta(\mathbf{x}, \mathbf{x}') K_\theta(\mathbf{x}', \mathbf{x}')^{-1/2}$
- Similarity in a single value: $k_\theta^C(\mathbf{x}, \mathbf{x}') = \frac{1}{d} \text{Tr } K_\theta^C(\mathbf{x}, \mathbf{x}')$.
- Other metrics possible in output space
 - e.g. for rotation invariance: $k_\theta^{C, \text{rot}}(\mathbf{x}, \mathbf{x}') = \frac{1}{2} \sqrt{\|K_\theta^C(\mathbf{x}, \mathbf{x}')\|_F^2 + 2 \det K_\theta^C(\mathbf{x}, \mathbf{x}')}$

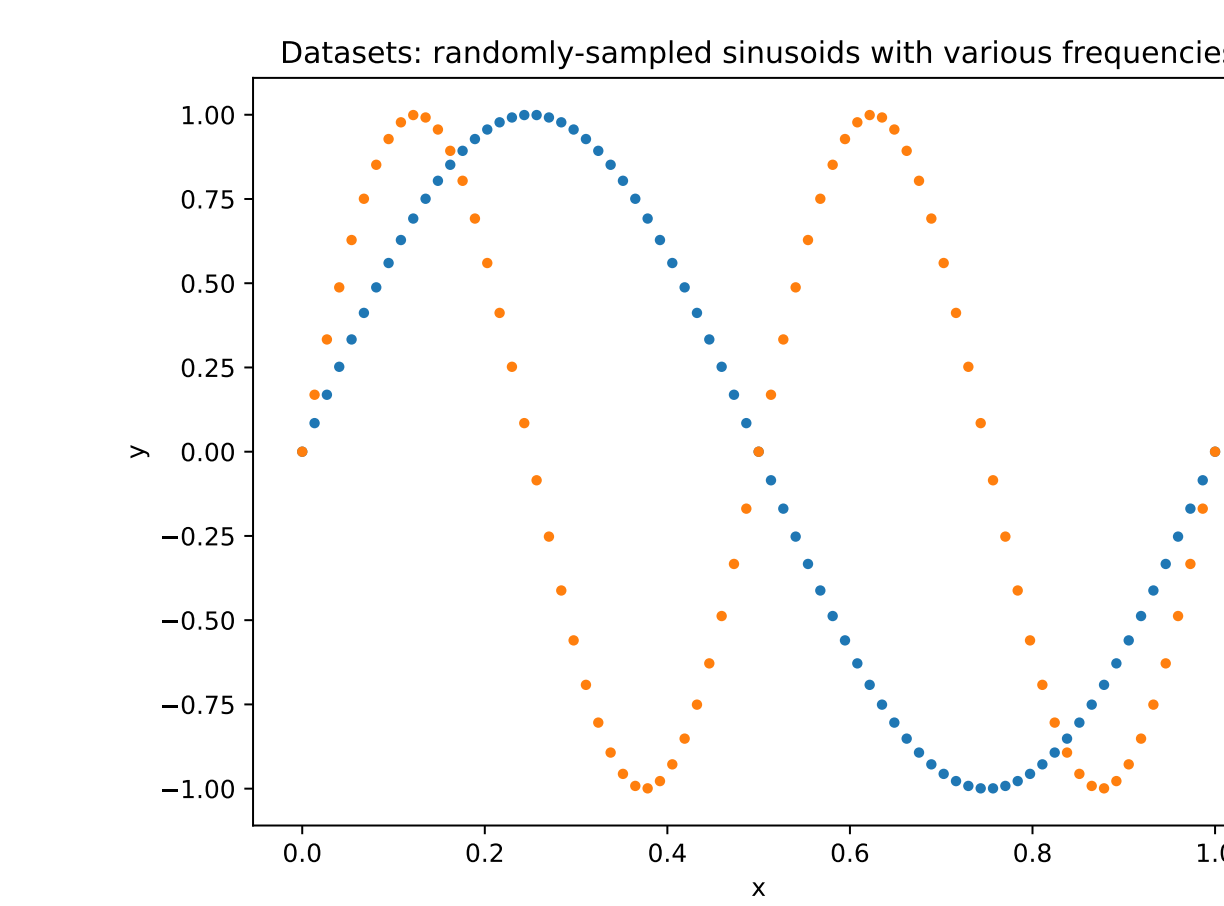
II - Estimating density

How many samples \mathbf{x}' are similar to \mathbf{x} according to the network? Many ways to count!

hard-thresholding with threshold $\tau \in [0, 1]$	soft estimate Fast! $O(NP)$	less-soft positive-only estimate ($\alpha > 0$)
$\sum_{\mathbf{x}'} \mathbb{1}_{k_\theta^C(\mathbf{x}, \mathbf{x}') \geq \tau}$	$\sum_{\mathbf{x}'} k_\theta^C(\mathbf{x}, \mathbf{x}') = \frac{\nabla_\theta f_\theta(\mathbf{x})}{\ \nabla_\theta f_\theta(\mathbf{x})\ } \cdot \frac{\nabla_\theta f_\theta}{\ \nabla_\theta f_\theta\ }$	$\sum_{\mathbf{x}'} \mathbb{1}_{k_\theta^C(\mathbf{x}, \mathbf{x}') > 0} k_\theta^C(\mathbf{x}, \mathbf{x}')^\alpha$

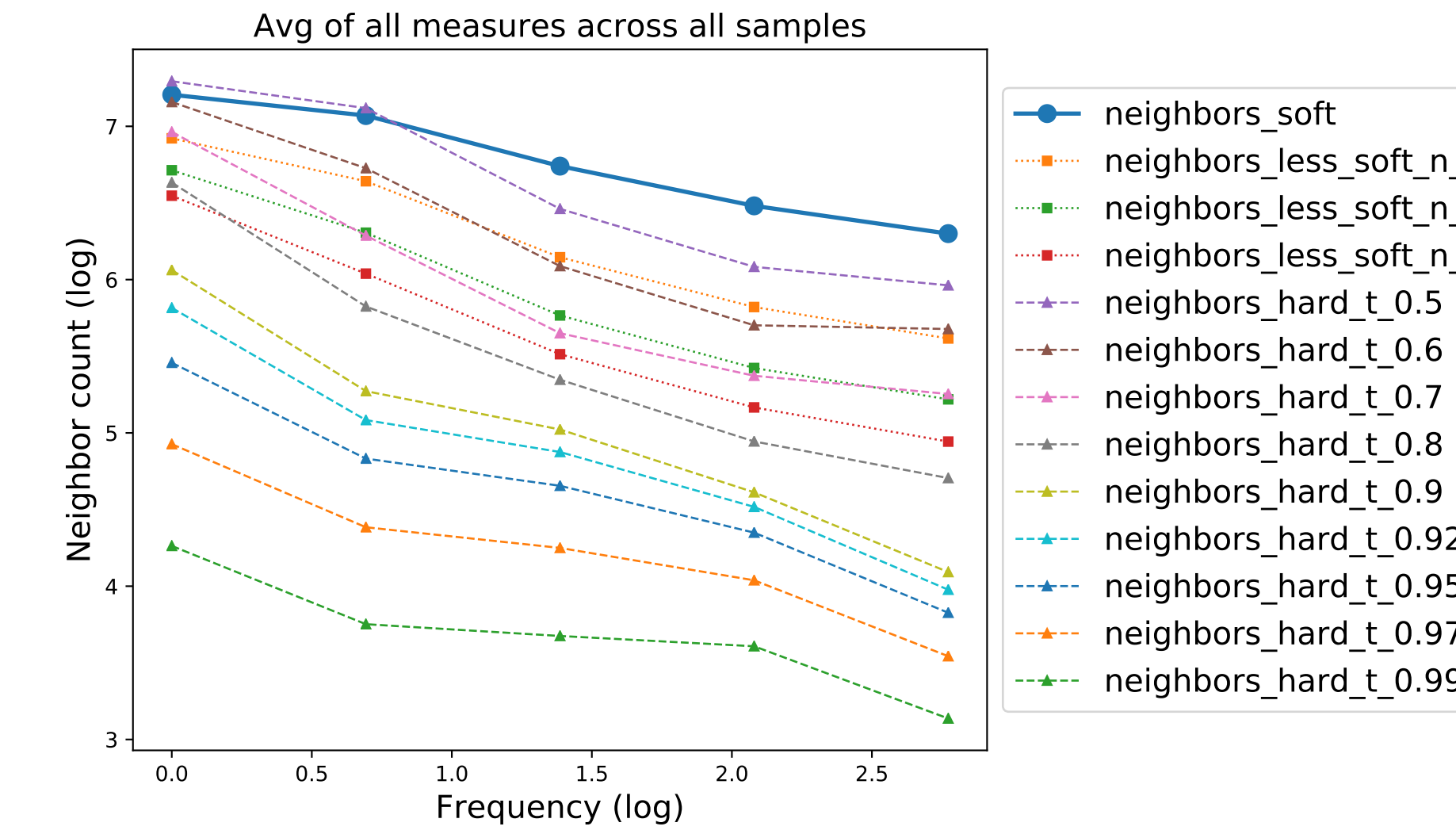
Toy dataset:

- Task f : learn to reproduce $\sin(2\pi f x)$, from $N = 2048$ random samples
- increase frequency $f \Rightarrow$ fewer neighbors

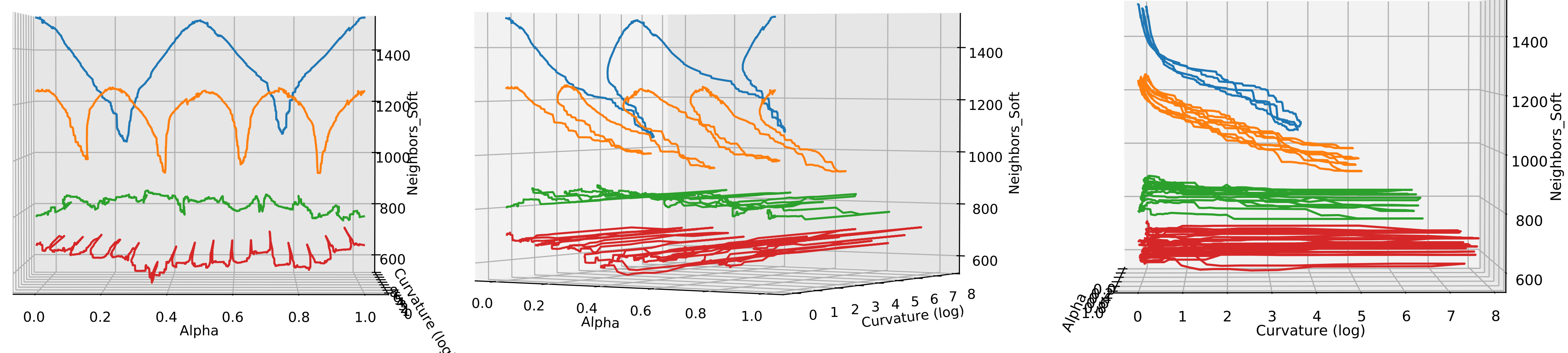


average
per task

⇒



no average



Soft estimate of the number of neighbors as a function of location and curvature

III - Futher possible uses of similarity

- similarity computation \Rightarrow density estimator \Rightarrow **active learning**
 - \Rightarrow + local variance (of predictions / of target labels): **overfit/underfit detection**
- **uncertainty** estimation: how much should the training set be changed to change the prediction $f_\theta(\mathbf{x})$ at point \mathbf{x} ?
 - \Rightarrow sum influence factors from all training points $\mathbf{x}_i \Rightarrow$ **reliability** $\propto \sum_i k_\theta(\mathbf{x}, \mathbf{x}_i) / \|\nabla_\theta f_\theta(\mathbf{x})\|^2$
- all quantities are **differentiable** \Rightarrow enforce similarity properties during training!
 - incite the network to consider given \mathbf{x} and \mathbf{x}' as (dis)similar : add $(\pm) k_\theta(\mathbf{x}, \mathbf{x}')$ to the loss
 - asking a subset \mathcal{S} of samples to be treated as similar :
 - $\sum_{i, j \in \mathcal{S}, i \neq j} k_\theta^C(\mathbf{x}_i, \mathbf{x}_j)$ rewrites as a function of $\mu = \frac{1}{n} \sum_{i \in \mathcal{S}} \frac{\nabla_\theta f_\theta(\mathbf{x}_i)}{\|\nabla_\theta f_\theta(\mathbf{x}_i)\|}$
 - **dynamics of learning** : speed-up effect noticed when training on MNIST with class-similarity enforcement

References

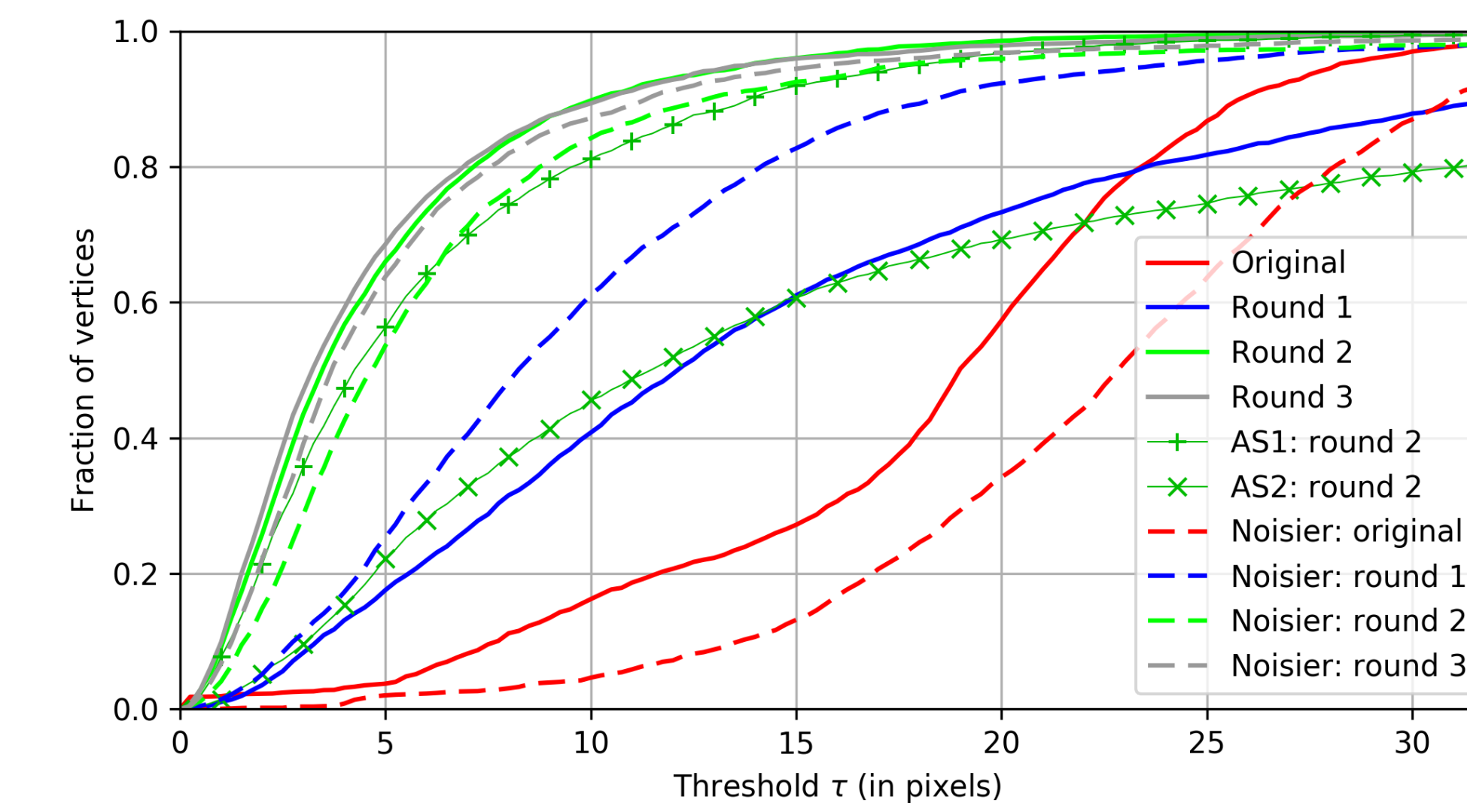
- [1] N. Girard et al. Noisy Supervision for Correcting Misaligned Cadaster Maps Without Perfect Ground Truth Data. In *IGARSS* 2019.
 [3] J. Lehtinen et al. Noise2noise: Learning image restoration without clean data. In *ICML* 2018.

IV - Analysis of self-denoising phenomena

Multimodal image registration task (RGB image to cadaster map) from noisy labels [1]



Red: initial dataset annotations (noisy)
Blue: prediction after learning from Red
Green: prediction after learning from Blue



Accuracy cumulative distributions (ground truth from [2]), i.e. fraction of pixels whose registration error \leq abscissa.

- ⇒ impressive self-denoising effect. Explanation from Noise2Noise [3]: same example x showed N times with noisy labels $\sim \mathcal{N}(y, \sigma_\varepsilon) \Rightarrow$ best prediction $\hat{y} = y \pm \frac{1}{\sqrt{N}} \sigma_\varepsilon$
- ⇒ noise averaging effect over labels of *similar* examples?

Formalism:

- input : \mathbf{x}_i
- true (unknown) label : y_i
- (unknown) noise : ε_i (iid, centered)
- noisy (available) label : $\tilde{y}_i = y_i + \varepsilon_i$
- predicted label : $\hat{y}_i = f_\theta(\mathbf{x}_i)$
- training loss : $L(\theta) = \sum_j \|\hat{y}_j - \tilde{y}_j\|^2$
- at convergence $\nabla_\theta E = 0 \Rightarrow \mathbb{E}_k[\hat{y}] = \mathbb{E}_k[\tilde{y}]$
- $\mathbb{E}_k[a] := \sum_j a_j k_\theta(\mathbf{x}_i, \mathbf{x}_j)$: mean value around \mathbf{x}_i
- $\hat{y}_i - \mathbb{E}_k[y] = \mathbb{E}_k[\varepsilon] + (\hat{y}_i - \mathbb{E}_k[\hat{y}])$
- $\hat{y}_i - \mathbb{E}_k[y]$: prediction error to smoothed true labels
- $\mathbb{E}_k[\varepsilon] \propto \sigma_\varepsilon \|k_\theta(\mathbf{x}_i, \cdot)\|_{L2} \Rightarrow$ denoising factor: 0.02
- Shift: $(\hat{y}_i - \mathbb{E}_k[\hat{y}])$: 4.4 px

Similar patch retrieval, and comparison with *perceptual loss*:



Discussion

- Fast similarity / density estimation opens the door to underfit/overfit/uncertainty analyses and control
- Extended Noise2Noise [3] to non-identical inputs: self-denoising effect as a function of inputs similarities
- Links with *Neural Tangent Kernel* [4]: same concept! used differently
- Future work: analyze and improve robustness to adversarial attacks
- Code available on GitHub: <http://github.com/Lydonr/netsimilarity> or scan the QR code:



- [2] E. Maggiori et al. Can semantic labeling methods generalize to any city? the Inria aerial image labeling benchmark. In *IGARSS* 2017.
 [4] A. Jacot et al. Neural tangent kernel: Convergence and generalization in neural networks. In *NIPS* 2018.