

1 First, we would like to thank all reviewers for their very positive comments on the theory presented in the paper. This is
 2 a theoretical paper indeed, introducing a new concept, mathematically principled and studied. In order to be applicable
 3 in practice, we show how to compute it and how to quickly approximate it, with code available (on the anonymous
 4 github link provided in the paper). We also check experimentally that our estimator behaves correctly. We believe this
 5 is already a nice set of contributions.

6 In addition, we also propose a certain number of possible uses and extensions of this concept, showing it could be
 7 useful in many different ways. We consider it is out of the scope of this paper to actually run such applications, which
 8 would be difficult to include in the paper anyway for space reasons without sacrificing ideas in the theoretical section.

9 However, we do understand the main criticisms about: (a) the lack of insights brought by the large-scale experiment on
 10 remote sensing image registration in Section 6, and (b) the lack of comparison to the *perceptual loss*. For (b), we propose
 11 to add such a comparison on nearest neighbor retrieval in Section 6. We notice that the *perceptual loss* sometimes
 12 performs reasonably well, but often not. For instance, we show below the closest neighbors to a structured residential area
 13 image, for the *perceptual loss* (first row: does not make sense) and for our similarity measure (second row: similar areas).



14
 15
 16 To tackle (a), we propose to show how the similarity experimental computations in Section 6 can be used to **solve the**
 17 **initial problem**, by explicitly turning similarity statistics into a **quantification of the self-denoising effect**, as follows.
 18 Let us denote by y_i the true (unknown) label for input \mathbf{x}_i , by \tilde{y}_i the noisy label given in the dataset, and by $\hat{y}_i = f_\theta(\mathbf{x}_i)$
 19 the label predicted by the network. We will denote the (unknown) noise by $\varepsilon_i = \tilde{y}_i - y_i$ and assume it is centered and
 20 i.i.d., with finite variance σ_ε . The training criterion is $E(\theta) = \sum_j \|\hat{y}_j - \tilde{y}_j\|^2$. At convergence, the training leads to
 21 a local optimum of the energy landscape: $\nabla_\theta E = 0$, that is, $\sum_j (\hat{y}_j - \tilde{y}_j) \nabla_\theta \hat{y}_j = 0$. Let's choose any sample i and
 22 multiply by $\nabla_\theta \hat{y}_i$: using $k_\theta^I(\mathbf{x}_i, \mathbf{x}_j) = \nabla_\theta \hat{y}_i \cdot \nabla_\theta \hat{y}_j$, we get: $\sum_j (\hat{y}_j - \tilde{y}_j) k_\theta^I(\mathbf{x}_j, \mathbf{x}_i) = 0$.

23 Let us denote by $k_\theta^{IN}(\mathbf{x}_j, \mathbf{x}_i) = k_\theta^I(\mathbf{x}_j, \mathbf{x}_i) (\sum_j k_\theta^I(\mathbf{x}_j, \mathbf{x}_i))^{-1}$ the normalized kernel, and by $\mathbb{E}_k[a] =$
 24 $\sum_j a_j k_\theta^{IN}(\mathbf{x}_j, \mathbf{x}_i)$ the mean of value a in the neighborhood of i , that is, the weighted average of the a_j with
 25 weights $k_\theta^I(\mathbf{x}_j, \mathbf{x}_i)$ normalized to sum up to 1. This is actually a Parzen window estimator. Then the previous property
 26 can be rewritten as $\mathbb{E}_k[\hat{y}] = \mathbb{E}_k[\tilde{y}]$. As $\mathbb{E}_k[\tilde{y}] = \mathbb{E}_k[y] + \mathbb{E}_k[\varepsilon]$, this yields: $\hat{y}_i - \mathbb{E}_k[y] = \mathbb{E}_k[\varepsilon] + (\hat{y}_i - \mathbb{E}_k[\hat{y}])$

27 *i.e.* the difference between the predicted \hat{y}_i and the average of the **true labels** in the neighborhood of i is equal to the
 28 average of the noise in the neighborhood of i , up to the deviation of the prediction \hat{y}_i from the average prediction in its
 29 neighborhood. **We want to bound the error** $\|\hat{y}_i - \mathbb{E}_k[y]\|$ **without knowing neither the true labels** y nor the noise ε .
 30 One can show that $\mathbb{E}_k[\varepsilon] \propto \text{var}_\varepsilon(\mathbb{E}_k[\varepsilon])^{1/2} = \sigma_\varepsilon \|k_\theta^{IN}\|_{L2}$. The **denoising factor** is thus $\|k_\theta^{IN}\|_{L2}$, which is between
 31 $1/\sqrt{N}$ and 1, depending on the neighborhood quality. It is $1/\sqrt{N}$ when all N data points are identical, *i.e.* all satisfying
 32 $k_\theta^C(\mathbf{x}_i, \mathbf{x}_j) = 1$. On the other extreme, this factor is 1 when all points are independent: $k_\theta^I(\mathbf{x}_i, \mathbf{x}_j) = 0 \quad \forall i \neq j$. This
 33 way we extend *noise2noise*[11] to real datasets with non-identical inputs. **In our remote sensing experiment**, we
 34 estimate this way a denoising factor of 0.02, consistent across all training rounds and inputs ($\pm 10\%$), implying that
 35 each training round contributed equally to denoising the labels. This is confirmed by Fig. 2, which shows the error
 36 steadily decreasing, on a control test where true labels are known. The shift $(\hat{y}_i - \mathbb{E}_k[\hat{y}])$ on the other hand can be
 37 directly estimated given the network prediction. In our case, it is 4.4px on average, which is close to the observed
 38 median error for the last round in Fig. 2. It is largely input-dependent, with variance 3.2px, which is reflected by the
 39 spread distribution of errors in Fig. 2. This input-dependent shift thus provides a hint about prediction reliability.

40 It is also possible to bound $(\hat{y}_i - \mathbb{E}_k[\hat{y}]) = \mathbb{E}_k[\hat{y}_i - \hat{y}]$ using only similarity information (without predictions \hat{y}).

41 **Theorem 1** implies that the application: $\frac{\nabla_\theta f_\theta(\mathbf{x})}{\|\nabla_\theta f_\theta(\mathbf{x})\|} \mapsto f_\theta(\mathbf{x})$ is well-defined, and it can actually be shown to be differ-
 42 entiable and Lipschitz with a network-dependent constant. Thus $\|f_\theta(\mathbf{x}) - f_\theta(\mathbf{x}')\| \leq C \left\| \frac{\nabla_\theta f_\theta(\mathbf{x})}{\|\nabla_\theta f_\theta(\mathbf{x})\|} - \frac{\nabla_\theta f_\theta(\mathbf{x}')}{\|\nabla_\theta f_\theta(\mathbf{x}')\|} \right\| =$
 43 $\sqrt{2}C \sqrt{1 - k_\theta^C(\mathbf{x}, \mathbf{x}')}$, yielding $\|\hat{y}_i - \hat{y}_j\| \leq \sqrt{2}C \sqrt{1 - k_\theta^C(\mathbf{x}_i, \mathbf{x}_j)}$ and thus $\mathbb{E}_k[\hat{y}_i - \hat{y}] \leq \sqrt{2}C \mathbb{E}_k \left[\sqrt{1 - k_\theta^C(\mathbf{x}_i, \cdot)} \right]$.

44 **Other:** thank you for the very relevant literature, and the nice application suggestion to GAN / cycle-consistency. We
 45 can postpone the paragraph *Dynamics of learning* to the appendix to make place for the section above if needed.