# A. Appendix

This section contains explanatory comments, clarifications and mathematical derivations to support the main text.

## A.1. Properties of the softmax function

Here, we discuss properties of the softmax function that are used below. We start with its stationarity:

$$\text{smax}(x + d, y + d) = \frac{\exp(x + d)}{\exp(x + d) + \exp(y + d)} \tag{26}$$

$$= \frac{\exp(d)\exp(x)}{\exp(d)(\exp(x) + \exp(y))} \tag{27}$$

$$= \frac{\exp(x)}{\exp(x) + \exp(y)} \tag{28}$$

$$= \text{smax}(x, y). \tag{29}$$

Setting $d = -x - y$, we get the following property:

$$\text{smax}(x, y) = \text{smax}(x - x - y, y - x - y) \tag{30}$$

$$= \text{smax}(-y, -x). \tag{31}$$

## A.2. Probabilistic derivation of $\mathcal{L}_{\text{log}}$ and $\mathcal{L}_{\text{sub}}$

The log loss has been previously discussed in the literature (e.g. Hoffer and Ailon [9]), but it was reported to perform quite badly. We believe that this is primarily because it was used at the wrong scale, not because it is unsuitable per se (quite the opposite, as we have empirically demonstrated in our experiments). In the following, we give a derivation of both the SSE loss (3) and the log loss (2) where we start from a probabilistic point of view, showing that the log loss in particular has sound theoretical foundations.

Our objective is to encourage that $d_p$ becomes smaller than $d_n$, where a larger gap is better. This can be expressed with the softmax function:

$$\text{smax}(d_p, d_n) = \frac{\exp(d_p)}{\exp(d_p) + \exp(d_n)} \tag{32}$$

The softmax function is in $[0, 1]$, that is, we can interpret it probabilistically as the uncertainty whether the two patches in the positive pair are indeed more similar than the ones in the negative pair. Our objective is to minimize this uncertainty:

$$\text{smax}(d_p, d_n) \longrightarrow 0, \text{ or, equivalently,} \tag{33}$$

$$1 - \text{smax}(d_p, d_n) \longrightarrow 1. \tag{34}$$

To measure the deviation from this objective, we can use the cross-entropy between actual and desired values:

$$\mathcal{L}_{\text{log}}(d_p, d_n) = -0 \cdot \log(\text{smax}(d_p, d_n)) \tag{35}$$

$$- 1 \cdot \log(1 - \text{smax}(d_p, d_n)) \tag{36}$$

$$= -\log(1 - \text{smax}(d_p, d_n)). \tag{37}$$

For the sake of a shorter notation, we rewrite the term under the logarithm:

$$1 - \text{smax}(d_p, d_n) = 1 - \frac{\exp(d_p)}{\exp(d_p) + \exp(d_n)} \tag{38}$$

$$= \frac{\exp(d_p) + \exp(d_n) - \exp(d_p)}{\exp(d_n) + \exp(d_p)} \tag{39}$$

$$= \frac{\exp(d_n)}{\exp(d_n) + \exp(d_p)} = \text{smax}(d_n, d_p). \tag{40}$$

Using (31), this can be further rewritten:

$$\text{smax}(d_n, d_p) = \text{smax}(-d_p, -d_n). \tag{41}$$

Hence, we get the final formulation of the log loss:

$$\mathcal{L}_{\text{log}}(d_p, d_n) = -\log \text{smax}(-d_p, -d_n). \tag{42}$$

As an alternative to the cross-entropy, one can simply calculate the squared sum of errors of the two cases (33) and (34):

$$\mathcal{L}_{\text{sse}}^*(d_p, d_n) = (0 - \text{smax}(d_p, d_n))^2 \tag{43}$$

$$+ (1 - (1 - \text{smax}(d_p, d_n))^2 \tag{44}$$

$$= 2\,\text{smax}(d_p, d_n)^2. \tag{45}$$

For simplicity, we leave away the constant factor 2 (this can be compensated by increasing the learning rate), yielding the SSE loss:

$$\mathcal{L}_{\text{sse}}(d_p, d_n) = \text{smax}(d_p, d_n)^2. \tag{46}$$

## A.3. Losses and performance functions

In this section, we show how the five losses discussed in Section 2 can be rewritten in terms of the performance functions given in Equations (7) to (9). The following three cases are straightforward:

$$\mathcal{L}_{\text{sub}} = [d_p - d_n + \alpha]_+ = [\alpha - \rho_{\text{sub}}]_+, \tag{47}$$

$$\mathcal{L}_{\text{sub2}} = [d_p^2 - d_n^2 + \alpha]_+ = [\alpha - \rho_{\text{sub2}}]_+, \tag{48}$$

$$\mathcal{L}_{\text{div}} = \left[1 - \frac{d_n}{d_p + \epsilon}\right]_+ = [1 - \rho_{\text{div}}]_+. \tag{49}$$

The two remaining cases, $\mathcal{L}_{\text{log}}$ and $\mathcal{L}_{\text{sse}}$ are less obvious, but using (29) we can write:

$$\mathcal{L}_{\text{log}} = -\log \text{smax}(-d_p, -d_n) \tag{50}$$

$$= -\log \text{smax}(-d_p + d_n, -d_n + d_n) \tag{51}$$

$$= -\log \text{smax}(d_n - d_p, 0) \tag{52}$$

$$= -\log \text{smax}(\rho_{\text{sub}}, 0), \tag{53}$$

$$\mathcal{L}_{\text{sse}} = \text{smax}(d_p, d_n)^2 \tag{54}$$

$$= \text{smax}(d_p - d_n, d_n - d_n)^2 \tag{55}$$

$$= \text{smax}(-(d_n - d_p), 0)^2 \tag{56}$$

$$= \text{smax}(-\rho_{\text{sub}}, 0)^2. \tag{57}$$

## A.4. Corner Cases of Generalized Scale Log Loss

In the following, we derive the two corner cases of $\tilde{\mathcal{L}}_{\log}$ from Section 3.3.

**Case $\delta \to \infty$:** Here, we show Equation (18). First, we rewrite the loss (setting $\rho := \rho_{\mathrm{sub}}$):

$$\tilde{\mathcal{L}}_{\log}(\rho) = \frac{1}{\delta}\mathcal{L}_{\log}(\delta(\rho - \alpha)) \tag{58}$$

$$= -\frac{1}{\delta}\log \mathrm{smax}(\delta(\rho - \alpha), 0) \tag{59}$$

$$= -\log \mathrm{smax}(\delta(\rho - \alpha), 0)^{1/\delta} \tag{60}$$

$$= -\log\left[\frac{\exp(\delta(\rho - \alpha))}{\exp(\delta(\rho - \alpha)) + \exp(0)}\right]^{1/\delta} \tag{61}$$

$$= -\log\left[\frac{\exp(\rho - \alpha)}{\sqrt[\delta]{\exp(\rho - \alpha)^{\delta} + \exp(0)}}\right]. \tag{62}$$

To see what happens as $\delta \to \infty$, we make a case distinction:

- $\rho - \alpha < 0$: in this case, the root in the denominator is dominated by $\exp(0) = 1$ because $\exp(\rho - \alpha)^{\delta}$ disappears. Hence, we have $\lim_{\delta \to \infty} \tilde{\mathcal{L}}_{\log}(\rho) = -\log[\exp(\rho - \alpha)] = \alpha - \rho$.

- $\rho - \alpha \geq 0$: in this case, $\exp(\rho - \alpha)$ dominates the root in the denominator. We have $\lim_{\delta \to \infty} \tilde{\mathcal{L}}_{\log}(\rho) = -\log[\exp(\rho - \alpha)/\exp(\rho - \alpha)] = 0$.

Concluding:

$$\lim_{\delta \to \infty} \tilde{\mathcal{L}}_{\log}(\rho) = \begin{cases} \alpha - \rho & \text{if } \alpha - \rho > 0 \\ 0 & \text{if } \alpha - \rho \leq 0 \end{cases} \tag{63}$$

$$= \max\{\alpha - \rho, 0\} \tag{64}$$

$$= [\alpha - \rho]_{+} = \tilde{\mathcal{L}}_{\mathrm{sub}}(\rho; \alpha). \tag{65}$$

**Case $\delta \to 0$:** Here, we show Equation (17). Actually, the loss $\mathcal{L}_{\log}$ diverges for $\delta \to 0$, but as we discussed in the text, we are not primarily interested in the actual value of a loss but in its derivative which tells us how strongly the loss acts. Therefore, we consider what happens for very small values of $\delta$ by investigating the derivative. It is given below in Equation (76). It is easy to see that $\lim_{\delta \to 0} \exp(\delta(\rho_{\mathrm{sub}} - \alpha)) = 1$, and consequently:

$$\lim_{\delta \to 0} \frac{\partial \mathcal{L}_{\log}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}} = -\frac{1}{1 + 1} - \frac{1}{2}. \tag{66}$$

By integrating, we get that the loss behaves asymptotically as the following loss:

$$\lim_{\delta \to 0} \tilde{\mathcal{L}}_{\log}(\rho_{\mathrm{sub}}; \alpha, \delta) \propto -0.5\rho_{\mathrm{sub}}. \tag{67}$$

## A.5. Antisymmetry

In Section 4.1, we mention that losses in $\rho_{\mathrm{sub}}$ are antisymmetric, i.e.,

$$\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}}(d_p, d_n))}{\partial d_p} = -\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}}(d_p, d_n))}{\partial d_n}. \tag{68}$$

Using the definition of the difference performance function in Equation (7) and the chain rule, we can derive:

$$\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial d_p} = \frac{\partial \mathcal{L}(d_n - d_p)}{\partial d_p} \tag{69}$$

$$= \frac{\partial(d_n - d_p)}{\partial d_p}\frac{\partial \mathcal{L}(d_n - d_p)}{\partial(d_n - d_p)} \tag{70}$$

$$= -\frac{\partial \mathcal{L}(d_n - d_p)}{\partial(d_n - d_p)} = -\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}}, \tag{71}$$

$$\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial d_n} = \frac{\partial \mathcal{L}(d_n - d_p)}{\partial d_n} \tag{72}$$

$$= \frac{\partial(d_n - d_p)}{\partial d_n}\frac{\partial \mathcal{L}(d_n - d_p)}{\partial(d_n - d_p)} \tag{73}$$

$$= +\frac{\partial \mathcal{L}(d_n - d_p)}{\partial(d_n - d_p)} = +\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}}. \tag{74}$$

It is immediately visible that

$$\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial d_p} = -\frac{\partial \mathcal{L}(\rho_{\mathrm{sub}})}{\partial d_n}, \tag{75}$$

i.e., any loss in $\rho_{\mathrm{sub}}$ is antisymmetric.

## A.6. Derivatives of the losses

For visualizing the losses, we need the derivatives of the losses w.r.t. their performance functions. For the three cases in Figure 4:

$$\frac{\partial \tilde{\mathcal{L}}_{\log}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}} = -\frac{1}{\exp(\delta(\rho_{\mathrm{sub}} - \alpha)) + 1} \tag{76}$$

$$\frac{\partial \tilde{\mathcal{L}}_{\mathrm{sse}}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}} = -2\frac{\mathrm{smax}(-\delta(\rho_{\mathrm{sub}} - \alpha), 0)^2}{\exp(\delta(\rho_{\mathrm{sub}} - \alpha))^2} \tag{77}$$

$$\frac{\partial \mathcal{L}_{\mathrm{sub}}(\rho_{\mathrm{sub}})}{\partial \rho_{\mathrm{sub}}} = \begin{cases} 0 & \text{if } \rho_{\mathrm{sub}} > \alpha \\ -1 & \text{otherwise} \end{cases} \tag{78}$$

For completeness, we also give the derivatives of the other two losses:

$$\frac{\partial \mathcal{L}_{\mathrm{sub2}}(\rho_{\mathrm{sub2}})}{\partial \rho_{\mathrm{sub2}}} = \begin{cases} 0 & \text{if } \rho_{\mathrm{sub2}} > \alpha \\ -1 & \text{otherwise} \end{cases} \tag{79}$$

$$\frac{\partial \mathcal{L}_{\mathrm{div}}(\rho_{\mathrm{div}})}{\partial \rho_{\mathrm{div}}} = \begin{cases} 0 & \text{if } \rho_{\mathrm{div}} > 1 \\ -1 & \text{otherwise} \end{cases} \tag{80}$$
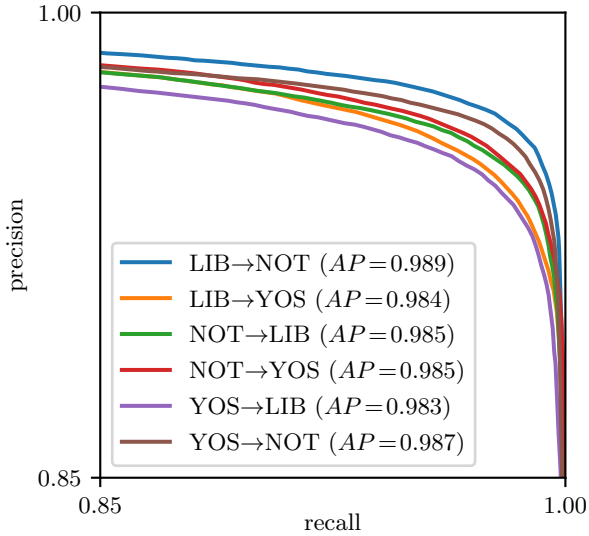
the loss in terms of the log loss $\mathcal{L}_{\log}$ (Equation (2)):

$$\mathcal{L}_{\text{L2}} = \sum_{i=1}^{N} -\log \text{smax}\left(-d_p^i, \log\sum_j \exp\left(-d_n^{i,j}\right)\right) \quad (82)$$

$$= \sum_{i=1}^{N} \mathcal{L}_{\log}\left(d_p^i, -\log\sum_j \exp\left(-d_n^{i,j}\right)\right). \quad (83)$$

In other words, we can interpret the loss of L2-Net as a kind triplet loss where many negative distances are amalgamated into one single value by means of a log-sum-exp expression (note that it therefore still suffers from the localized-context problem). This gives more weight to the smaller negative distances than to larger ones, that is, we can understand L2-Net as an approximation to our scale-aware sampling scheme where we simply pick the smallest negative distance $\min_j\{d_n^{i,j}\}$. The corresponding loss then is the one we presented in Equation (23), with $\mathcal{L}_{\text{triplet}} := \mathcal{L}_{\log}$:

$$\sum_{i=1}^{N} \mathcal{L}_{\log}\left(d_p^i, \min_j\{d_n^{i,j}\}\right). \quad (84)$$



Figure 8: PR curves and their APs on the mixed-context loss. Model is trained on each of the three scenes: `liberty`, `notredame`, and `yosemite` and then tested on the other datasets.

## A.7. Full PR curves and their AP

In this section, we report the full PR curves and their AP in Figure 8 on the UBC benchmark [7] by following the procedures described in [22]. The PR curve is computed from the default 100'000 evaluation patch pairs of each dataset. The values are analogous to the last row of Table 1

## A.8. Scale-aware sampling and L2-Net

Based on our discussion of the importance of scale in descriptor learning, we introduced scale-aware sampling in Section 4.3 and empirically demonstrated in Section 5 that it indeed brings a considerable boost in performance. In fact, if we look at the losses that were presented in key papers in the last two years, we can see that they followed a similar idea, only less consequently. One example is L2-Net [23] whose loss can be written as follows (it differs in details from the formulation in the original paper, but is true to the main idea):

$$\mathcal{L}_{\text{L2}} = -\sum_{i=1}^{N} \log\left[\frac{\exp(-d_p^i)}{\exp(-d_p^i) + \sum_j \exp(-d_n^{i,j})}\right]. \quad (81)$$

At first, glance, it is difficult to see the relationship with scale-aware sampling. However, if we use the fact that $\sum_j \exp(-d_n^{i,j}) = \exp\log\sum_j \exp(-d_n^{i,j})$, we can understand the fraction as a softmax function and hence rewrite