# Convolutional Neural Networks learn compact local image descriptors

**Christian Osendorfer**                                 OSENDORF@IN.TUM.DE
**Justin Bayer**                                              BAYERJ@IN.TUM.DE
**Patrick van der Smagt**                              SMAGT@IN.TUM.DE
Technische Universität München
85748 Garching, Germany

## Abstract

A standard deep convolutional neural network paired with a suitable loss function learns compact local image descriptors that perform comparably to state-of-the art approaches.

## 1. General Learning Architecture

Recently, several machine learning based approaches (Brown et al., 2010; Simonyan et al., 2012; Trzcinski et al., 2012) have shown impressive results for finding compact low-level image representations. These representations are considered good when corresponding image patches are described by representations that are close by.

DrLim (Hadsell et al., 2006) is a framework for energy based models that learns representation using only such correspondence relationships. We utilize DrLim to train a convolution neural network for learning low-dimensional mappings for low-level image patches.

The main idea behind DrLim is to map similar (i.e. corresponding) image patches to nearby points on the output manifold and dissimilar image patches to distant points. It is defined over pairs of image patches, $x_1, x_2$. The $i$-th pair $(x_1^i, x_2^i)$ is associated with a label $y^i$, with $y^i = 1$ if $x_1^i$ and $x_2^i$ are deemed similar and $y^i = 0$ otherwise. We denote by $d(x_1, x_2; \theta)$ the parameterized distance function between the representations of $x_1$ and $x_2$ that we want to learn. Based on $d(x_1, x_2; \theta)$ we define DrLim's loss function $\ell(\theta)$:

$$\ell(\theta) = \sum_i y^i \ell_{\text{pll}}(d(x_1^i, x_2^i; \theta)) + (1 - y^i)\ell_{\text{psh}}(d(x_1^i, x_2^i; \theta))$$

*Preliminary work.*

We denote with $\ell_{\text{pll}}(\cdot)$ the partial loss function for similar pairs (it *pulls* similar pairs together) and with $\ell_{\text{psh}}(\cdot)$ the partial loss function for dissimilar pairs (it *pushes* dissimilar pairs apart). $\ell_{\text{psh}}$ is defined as in (Hadsell et al., 2006):

$$\ell_{\text{psh}}(d(x_1, x_2; \theta)) = c_{\text{psh}}[\max(0, m_{\text{psh}} - d(x_1, x_2; \theta))]^2$$

$m_{\text{psh}}$ is the push *margin*: Dissimilar pairs are not pushed farther apart if they already are at a distance greater than the push margin. $c_{\text{psh}}$ is a scaling factor.

For $\ell_{\text{pll}}$ we use a loss similar to hinge loss, differently to the loss function proposed in the original DrLim formulation:

$$\ell_{\text{pll}}(d(x_1, x_2; \theta)) = c_{\text{pll}}[\max(0, d(x_1, x_2; \theta) - m_{\text{pll}})]$$

$c_{\text{pll}}$ is a scaling factor, $m_{\text{pll}}$ is a pull *margin*: Similar pairs are pulled together only if they are at a distance above $m_{\text{pll}}$.

$d(x_1, x_2; \theta)$ is defined as the Euclidean distance between the learned representations of $x_1$ and $x_2$:

$$d(x_1, x_2; \theta) = \|f(x_1; \theta) - f(x_2; \theta)\|_2$$

$f(\cdot)$ denotes the mapping from the (high-dimensional) input space to the low-dimensional space. In this paper, $f$ is a convolutional neural network(Jarrett et al., 2009). The layers of the convolutional network comprise a convolutional layer $C_1$ (kernel size $5 \times 5$) with 6 feature maps, a subsampling layer $S_1$, a second convolutional layer $C_2$ (kernel size $6 \times 6$) with 21 feature maps, a subsampling layer $S_2$, a third convolutional layer $C_3$ (kernel size $5 \times 5$) with 55 feature maps and a fully connected layer with 32 units.

## 2. Experiments

We evaluate our proposed model on the dataset from (Brown et al., 2010). The dataset is based on more

than 1.5 million image patches ($64 \times 64$ pixels) of three different scenes: the Statue of Liberty (about 450,000 patches), Notre Dame (about 450,000 patches) and Yosemites Half Dome (about 650,000 patches). We denote these scenes with LY, ND and HD respectively. There are 250000 corresponding image patch pairs and 250000 non-corresponding image patch pairs available for every scene. We train on one scene and evaluate the learned embedding function on the other two scenes. Evaluation is done on the same test sets (50000 matching and non-matching pairs) used also by other approaches.

Table 1 shows that convolutional networks (last entry) perform comparably to other state-of-the-art approaches. The appeal of a simple parameteric model like a convolutional neural network is that it does not require any complex paramter tuning or pipeline optimization and that it can be integrated into larger systems that can then be trained in an end-to-end fashion (Hadsell, 2008).

The architecture is trained with standard gradient descent. Training stops when a local minima of the Dr-Lim objective is reached. Notably, the hyperparameters ($c_\text{pll}$, $m_\text{pll}$, $c_\text{psh}$, $m_\text{psh}$) used in our evaluation are *not* scene dependent.

## 3. More data

Convolutional Neural Networks benefit from abundant data (Ciresan et al., 2012; Krizhevsky et al., 2012). Utilizing data from two scenes improves error rates noticebly: We get 15.1% on LY with combined training on ND and HD (in total 1M patch pairs). Similarly, we get 8.5% on ND and 14.3% on HD.

## References

M. Brown, G. Hua, and S. Winder. Discriminative learning of local image descriptors. *IEEE PAMI*, 2010.

Dan Ciresan, Ueli Meier, and Jürgen Schmidhuber. Multi-column deep neural networks for image classification. In *Proc. CVPR*, 2012.

R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *Proc. CVPR*, 2006.

R.T. Hadsell. *Learning long-range vision for an offroad robot*. PhD thesis, New York University, 2008.

K. Jarrett, K. Kavukcuoglu, M.A. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In *Proc. ICCV*, 2009.

| Method | Tr. set | Test set | | |
|---|---|---|---|---|
| | | LY | ND | HD |
| SIFT | – | 31.7 | 22.8 | 25.6 |
| L-BGM (64d) | LY | – | 14.1 | 19.6 |
| | ND | 18.0 | – | 15.8 |
| | HD | 21.0 | 13.7 | – |
| Brown et al. (29d) | LY | – | × | × |
| | ND | 16.8 | – | 13.5 |
| | HD | 18.2 | 11.9 | – |
| Simonyan et al. (29d) | LY | – | × | × |
| | ND | 14.5 | – | 12.5 |
| | HD | 17.4 | 9.6 | – |
| CNN (32d) | LY | – | $11.2_{\pm0.3}$ | $18.5_{\pm0.5}$ |
| | ND | $16.4_{\pm0.3}$ | – | $16.2_{\pm0.3}$ |
| | HD | $18.9_{\pm0.4}$ | $10.7_{\pm0.2}$ | – |

*Table 1.* Error rates, i.e. the percent of incorrect matches when 95% of the true matches are found. Every subtable, indicated by an entry in the *Method* column, denotes a descriptor algorithm. The line below every method denotes the size of the desciptor (e.g. 32d denotes a 32 dimensional descriptor). The 128 dimensional SIFT descriptor (Lowe, 2004) does not require learning (denoted by − in the column *Tr. set* (i.e. Training set)). The numbers in the columns labeled LY, ND and HD are the error rates of a method on the respective test set for this scene. (Brown et al., 2010; Simonyan et al., 2012) do not have results when trainend on the LY scene (indicated by ×). L-BGM is presented in (Trzcinski et al., 2012). The mean error rates for convolutional neural networks (CNN) are given with a standard deviation over 10 runs.

Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. NIPS*, 2012.

D.G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110, 2004.

K. Simonyan, A. Vedaldi, and A. Zisserman. Descriptor learning using convex optimisation. In *Computer Vision–ECCV 2012*, 2012.

T. Trzcinski, M. Christoudias, V. Lepetit, and P. Fua. Learning image descriptors with the boosting-trick. In *Proc. NIPS*, 2012.