

VISUAL SCENE REPRESENTATIONS: SCALING AND OCCLUSION IN CONVOLUTIONAL ARCHITECTURES*

Stefano Soatto, Jingming Dong & Nikolaos Karianakis

UCLA Vision Lab

University of California, Los Angeles

Los Angeles, CA 90095, USA

{soatto,dong}@cs.ucla.edu, nikarianakis@ucla.edu

ABSTRACT

We study the structure of representations, defined as approximations of minimal sufficient statistics that are maximal invariants to nuisance factors, for visual data subject to scaling and occlusion of line-of-sight. We derive analytical expressions for such representations and show that, under certain restrictive assumptions, they are related to features commonly in use in the computer vision community. This link highlights the conditions tacitly assumed by these descriptors, and also suggests ways to improve and generalize them.

1 INTRODUCTION

Soatto & Chiuso (2014) define an optimal representation as a minimal sufficient statistic (of past data for the scene) and a maximal invariant (of future data to nuisance factors), and propose a measure of how “useful” (informative) a representation is, via the uncertainty of the prediction density. What is a nuisance depends on the task, that includes decision and control actions about the surrounding environment, or *scene*, and its geometry (shape, pose), photometry (reflectance), dynamics (motion) and semantics (identities, relations of “objects” within).

We show that optimal management of nuisance variability due to occlusion is generally intractable, but can be approximated leading to a composite (correspondence) hypothesis test, which provides grounding for the use of “patches” or “receptive fields,” ubiquitous in practice. The analysis reveals that the size of the domain of the filters should be *decoupled* from spectral characteristics of the image, unlike traditionally taught in scale-space theory, an unintuitive consequence of the analysis. This idea has been exploited by Dong & Soatto (2015) to approximate the optimal descriptor of a *single image*, under an explicit model of image formation (the Lambert-Ambient, or LA, model) and nuisance variability, leading to DSP-SIFT. Extensions to multiple training images, leading to MV-HoG and R-HoG, have been championed by Dong et al. (2013). Here, we apply domain-size pooling to the scattering transform Bruna & Mallat (2011) leading to DSP-SC, to a convolutional neural network, leading to DSP-CNN, and to deformable part models Felzenszwalb et al. (2008), leading to DSP-DPM, in Sect. 2.2, 2.3 and 2.4 respectively.

We treat images as random vectors x, y and the scene θ as an (infinite-dimensional) parameter. An optimal representation is a function ϕ of past images $x^t \doteq \{x_1, \dots, x_t\}$ that maximally reduces uncertainty on questions about the scene Geman et al. (2015) given images from it and regardless of nuisance variables $g \in G$. In Soatto & Chiuso (2014) the sampled orbit anti-aliased (SOA) likelihood is introduced as:

$$\hat{L}_{G,\epsilon}(\theta; x) = \max_i \hat{L}(\theta, g_i; x), \quad i = 1, \dots, N(\epsilon) \quad (1)$$

where

$$\hat{L}(\theta, g_i; x) \doteq \int_G L(\theta, g; x) dP(g) \quad (2)$$

and $L(\theta, g; x) \doteq p_{\theta,g}(x)$ is the joint likelihood, understood as a function of the parameter θ and nuisance g for fixed data x , with $dP(g) = w(g^{-1})d\mu(g)$ an *anti-aliasing* measure with positive

*Also UCLA Technical Report CSD140024, November 12, 2014

weights w . The SOA likelihood is an optimal representation in the sense that, for any ϵ , it is possible to choose N and a finite number of samples $\{g_i\}_{i=1}^N$ so that $\phi_\theta(x^t) \doteq \hat{L}_{G,\epsilon}(\theta; x^t)$ approximates to within ϵ a minimal sufficient statistic (of x^t for θ) that is maximally invariant to group transformations in G . This result is valid under the assumptions of the Lambert-Ambient (LA) model [Dong & Soatto \(2014\)](#), which is the simplest known to capture the phenomenology of image formation including scaling, occlusion, and rudimentary illumination.

2 CONSTRUCTING VISUAL REPRESENTATIONS

Theorem 1 (Contrast invariant). *Given a training image x and a test image y , assuming that the latter is affected by noise that is independent in the gradient direction and magnitude, then the maximal invariant of y to the group G of contrast transformations is given by*

$$p_{x,G}(y) = p(\angle \nabla y | x) \|\nabla x\|. \quad (3)$$

Note that, other than for the gradient, the computations above can be performed point-wise under the assumption of LA model, so we could write (3) at each pixel y_i : if $\alpha \doteq \angle \nabla y_i$,

$$\phi_x(\alpha) = \prod_i \mathcal{N}_{\mathbb{S}^1}(\alpha_i - \angle \nabla x_i; \epsilon_\alpha) \|\nabla x_i\| \quad (4)$$

Note that (4) is invariant to contrast transformations of y , but *not* of x . Invariance to contrast transformations in the (single) *training* image can be performed by normalizing the likelihood, which in turn can be done by simply dividing by the integral over α , which is the ℓ^1 norm of the histogram across the entire image/patch

$$\frac{\phi_x(\alpha)}{\|\phi_x(\alpha)\|_{\ell^1}} = \frac{p(\alpha|x)\|\nabla x\|}{\int p(\alpha|x)d\alpha\|\nabla x\|} = p(\alpha|x) \quad (5)$$

that should be used instead of the customary ℓ^2 [Lowe \(2004\)](#). Once invariance to contrast transformations is achieved, which can be done on a single image x , we are left with nuisances G that include general viewpoint changes, including the occlusions they induce. This can be handled by computing the SOA likelihood with respect to G of $SE(3)$ (Sect. 2.1) from a training sample x^t , leading to

$$\hat{L}(\theta, g_i; x^t) = \left\{ \int_G \phi_{x^t}(\alpha | g_i \circ g) dP(g) \right\}_{i=1}^N \quad (6)$$

Occlusion, or visibility, is arguably the single most critical aspect of visual representations. It enforces *locality*, as dealing with occlusion nuisances entails searching through, or marginalizing, all possible (multiply-connected) subsets of the test image. This power set is clearly intractable even for very small images. *Missed detections* (treating a co-visible pixel as occluded) and *false alarms* (treating an occluded pixel as visible) have different costs: Omitting a co-visible pixel from Ω decreases the likelihood by a factor corresponding to multiplication by a Gaussian for samples drawn from the same distribution; vice-versa, including a pixel from Ω^c (false alarm) decreases the log-likelihood by a factor equal to multiplying by a Gaussian evaluated at points drawn from another distribution, such as uniform. So, testing for correspondence on *subsets of the co-visible regions*, assuming the region is sufficiently large, reduces the power, but not the validity, of the test. This observation can be used to *fix the shape* of the regions, *leaving only their size to be marginalized, or searched over*. This reasoning justifies the use of “patches” or “receptive fields” to seed image matching, but emphasizes that a search over different *sizes* [Dong & Soatto \(2015\)](#) is needed.

Together with the SOA likelihood, this also justifies the local marginalization of *domain sizes*, along with translation, as recently championed in [Dong & Soatto \(2015\)](#).

Corollary 1 (DSP-SIFT). *The DSP-SIFT descriptor [Dong & Soatto \(2015\)](#) approximates an optimal representation (6) for G the group of planar similarities and local contrast transformations, when the scene is a single training image, and the test image is restricted to a subset of its domain.*

The assumptions underlying all local representations built using a single image break down when the scene is not flat and not moving parallel to the image plane. In this case, multiple views are necessary to manage nuisance due to general viewpoint changes.

2.1 GENERAL VIEWPOINT CHANGES

If a co-variant translation-scale *and size* sampling/anti-aliasing mechanism is employed, then around each sample the only residual variability to viewpoint $SE(3) = \mathbb{R}^3 \times SO(3)$ is reduced.

In some cases, a consistent reference (canonical element) for both training and test images is available when scenes or objects are geo-referenced: The projection of the gravity vector onto the image plane Jones & Soatto (2011). In this case, $\hat{\alpha}$ is the angle of the projection of gravity onto the image plane (well defined unless they are orthogonal). Alternatively, multiple (principal) orientation references can be selected based on the norm of the directional derivative Lowe (2004):

$$p_{\theta}(\alpha|G) = p_{\theta}(\alpha|\hat{\alpha}). \quad (7)$$

This leaves out-of-plane rotations to be managed. Dong et al. (2013) have proposed extensions of local descriptors to multiple views, based on a sampling approximation of the likelihood function, \hat{p}_{θ} , or on a point estimate of the scene $p_{\hat{\theta}}$, MV-HoG and R-HoG respectively. The estimated scene has a geometric component (shape) \hat{S} and a photometric component (radiance) $\hat{\rho}$, inferred from the LA model as described in Dong & Soatto (2014). Once the effects of occlusions are considered (which force the representation to be local), and the effects of general viewpoint changes are accounted for (which creates the necessity for multiple training images of the same scene), a maximal contrast/viewpoint/occlusion invariant can be approximated: the SOA likelihood (6) becomes:

$$\hat{L}_{SE(3),\epsilon(N)}(\alpha_i) = \max_k \left\{ \int_{SO(3)} \mathcal{N}_{\mathbb{S}^1}(\alpha_i - \hat{\rho} \circ g_k g \circ \pi_{\hat{S}}^{-1}(x_j); \epsilon_{\alpha}) \kappa_{\sigma}(i-j) d\mu(j) dP(\sigma) dP_{SO(3)}(g) \right\}_{k=1}^N \quad (8)$$

in addition to domain-size pooling. The assumption that all existing multiple-view extensions of SIFT do *not* overcome is the conditional independence of the intensity of different pixels. This is discussed in Soatto & Chiuso (2014) for the case of convolutional deep architectures, and in the next section for Scattering Networks. Capturing the joint statistics of different components of the SOA likelihood is key to modeling intra-class variability of object or scene categories.

2.2 DSP-SCATTERING NETWORKS

The scattering transform Bruna & Mallat (2011) convolves an image (or patch) with a Gabor filter bank at different rotations and dilations, takes the modulus of the responses, and applies an averaging operator to yield the scattering coefficients. This is repeated to produce coefficients at different layers in a scattering network. The first layer is equivalent to SIFT Bruna & Mallat (2011), in the sense that (3) can be implemented via convolution with a Gabor element with orientation α then taking the modulus of the response. One could conjecture that domain-size pooling (DSP) applied to a scattering network would improve performance in tasks that involve changes of scale and visibility. We call the resulting method DSP Scattering Transform (DSP-SC). Indeed, this is the case, as we show in the Appendix of Soatto et al. (2014), where we compare DSP-SC to the single-scale scattering transform (SC) to the datasets of Mikolajczyk & Schmid (2003) (Oxford) and Fischer et al. (2014).

2.3 DSP-CNN

Deep convolutional architectures can be understood as implementing successive approximations of an optimal representation by stacking layers of (conditionally) independent local representations of the form (8), which have been shown by Soatto & Chiuso (2014) to increasingly achieve invariance to large deformations, despite locally marginalizing only affine (or similarity) transformations. As Dong & Soatto (2015) did for SIFT, and as we did for the Scattering Transform above, we conjectured that pooling over domain size would improve the performance of a convolutional network. In the Appendix of Soatto et al. (2014), we report experiments to test the conjecture using a pre-trained network which is fine-tuned with domain-size pooling on benchmark datasets.

2.4 DSP-DPM

We have also developed domain-size pooling extensions of deformable part models (DPMs) Felzenszwalb et al. (2008), small trees of local HOG descriptors (“parts”), whereby local photometry is

encoded in the latter (nodes), and geometry is encoded in their position on the image relative to the root node (edges). Intra-class shape variability is captured by the posterior density of edge values, learned from samples. Photometry is captured by a “HOG pyramid” where the *size* of each part is pre-determined and fixed relative to the root. One could therefore conjecture that performing anti-aliasing with respect to the size of the parts would improve performance. Experimental results, reported in the Appendix of Soatto et al. (2014), validate the conjecture.

ACKNOWLEDGMENTS

We acknowledge discussions with Alessandro Chiuso, Joshua Hernandez, Arash Amini, Ying-Nian Wu, Taco Cohen, Virginia Estellers, Jonathan Balzer. Research supported by ONR N000141110863, NSF RI-1422669, and FA8650-11-1-7154.

REFERENCES

- Bruna, J. and Mallat, S. Classification with scattering operators. In *Proc. IEEE Conf. on Comp. Vision and Pattern Recogn.*, 2011.
- Dong, J. and Soatto, S. The Lambert-Ambient Shape Space and the Systematic Design of Feature Descriptors. R. Cipolla, S. Battiato, G.-M. Farinella (Eds), Springer Verlag, 2014.
- Dong, J. and Soatto, S. Domain-size pooling in local descriptors: DSP-SIFT. In *Proc. IEEE Conf. on Comp. Vision and Pattern Recogn.*, 2015.
- Dong, J., Karianakis, N., Davis, D., Hernandez, J., Balzer, J., and Soatto, S. Multi-view feature engineering and learning. In *Proc. IEEE Conf. on Comp. Vision and Pattern Recogn.*, 2015. (also *ArXiv: 1311.6048*, 2013).
- Felzenszwalb, P., McAllester, D., and Ramanan, D. A discriminatively trained, multiscale, deformable part model. In *CVPR*, pp. 1–8, 2008.
- Fischer, P., Dosovitskiy, A., and Brox, T. Descriptor matching with convolutional neural networks: a comparison to sift. *ArXiv:1405.5769*, 2014.
- Geman, D., Geman, S., Hallonquist, N., and Younes, L. Visual turing test for computer vision systems. *Proceedings of the National Academy of Sciences*, 112(12):3618–3623, 2015.
- Jones, E. and Soatto, S. Visual-inertial navigation, localization and mapping: A scalable real-time large-scale approach. *Intl. J. of Robotics Res.*, 2011.
- Lowe, D. G. Distinctive image features from scale-invariant keypoints. *IJCV*, 2(60):91–110, 2004.
- Mikolajczyk, K. and Schmid, C. A performance evaluation of local descriptors. 2003.
- Soatto, S. and Chiuso, A. Visual scene representations: sufficiency, minimality, invariance and deep approximation. Proc. of the ICLR Workshop, 2015 (also *ArXiv: 1411.7676*, 2014).
- Soatto, S., Dong, J., and Karianakis, N. Visual scene representation: scaling and occlusion in convolutional architectures. (Extended version of this manuscript) Technical report UCLA CSD140024, 2014.