

SCATTERING CONVOLUTIONAL HIDDEN MARKOV TREE

J.B. REGLI, J. D. B. NELSON

UCL, Department of Statistical Science

ABSTRACT

A Scattering convolutional hidden Markov tree proposes a new inference mechanism for high-dimensional signals by combining the interesting signal representation created by the scattering transform to a powerful probabilistic graphical model. A wavelet scattering network computes a translation invariant and stable to deformations representation of a signal while still preserving its informative content. Such properties are acquired by cascading wavelet transform convolutions with non-linear modulus and averaging operators. The network's structure and its distributions are described using a Hidden Markov Tree. This yield a generative model for high-dimensional inference. It offers a mean for performing several inference tasks among which are predictions. The scattering convolutional hidden Markov tree displays promising results on classification tasks of complex images in the difficult case where the number of training examples is extremely limited.

Index Terms— Scattering network, Hidden Markov Model, Classification, Deep network

1 Introduction

The standard approach to classify high dimensional signals can be expressed as a two steps procedure. First the data are projected in a feature space where the task at hand is simplified. Then prediction is done using a simple predictor in this new representational space. The mapping can either be hand-built—e.g. Fourier transform, wavelet transform—or learned. In the last decade methods for learning the projection have drastically improved under the impulsion of the so called deep learning. Deep neural networks (sometime enriched by convolutional architecture) have been able to learn very effective representations for a given dataset and a given task [1–3]. Such method have achieved state of the art on many standard problems [4, 5] as well as real world applications [6].

However deep learning methods are only efficient when we have access to a vast quantity of training examples [7]. But in some cases, such as in medical or defence applications for example, datapoints are rare or using an expert for hand-labelling them is either time-consuming, costly or subjective.

Hence in situations where training examples are expensive to collect, learning has to be performed on smaller datasets. In that case using a fix hand crafted set of filters seems to be one of the best solution [8]. Recently Mallat described the scattering transform [9] a fixed bank of wavelet filters used to generate data representation in a convolutional neural networks like architecture. This representational method associated to a support vector machine classifier achieved close to state of the art performances of some standard datasets [10]. Furthermore this method can be accurately applied to relatively smaller datasets [11].

When an extremely low number of training examples are available—one-shot learning [12]—generative classification methods provide better performances than their discriminative counterparts [13]. Model a signal representation using generative probabilistic graphical models have been successfully done for various wavelet transform [14, 15]. In those work Hidden Markov trees are used to model the wavelet coefficients distribution.

In a similar fashion, we propose to model Mallat's scattering convolutional network [10] using hidden Markov trees. This combines a recently proposed deterministic analytically tractable transformation inspired by deep convolutional to a probabilistic graphical model in order to create a powerful probabilistic tool to handle high dimensional prediction problems. As opposed to previous works modeling wavelet trees, the use of scattering transforms allow us to leverage their interesting representational properties but also force us to adapt the HMT model to non-homogeneous, non regular trees. In contrast to passing the raw scattering coefficients into a classifier, our proposed method captures dependencies between different layers in a generative probabilistic model. Moreover as opposed to the commonly used simple classification method, once trained our model can tackle prediction problems but also other inference tasks—e.g. generation, sensitivity analysis...—and can also outperform SVMs when only a very low number of training examples are available.

The remainder of this paper introduces the scattering convolutional hidden Markov tree and is organised as follow. In section 2 we review the scattering transform and some of its properties. Section 3 introduces the proposed Scattering Hidden Markov Tree (SCHMT). Section 4 we perform classification on a selection of standard datasets restricted to only a few training samples. We draw conclusions in Section 5.

J.-B. Regli is funded by a Dstl/UCL Impact studentship. J. D. B. Nelson is partially supported by grants from the Dstl and Innovate UK/EPSRC.

2 Scattering networks

Scattering convolutional networks (SCNs) [16] are Convolutional Neural Networks (CNNs) [3] using a fixed filter bank of wavelets. Those filters can be hand-crafted to yield descriptors with the desired invariances [9]. For image classification tasks, one is interested in descriptors that are—at least—stable to deformations and invariant to translations. Note that SCNs producing more complexes set of invariances exist but on the remainder of this paper we consider only on descriptors with the previously mentioned properties.

2.1 Scattering transform

Wavelets are localized functions stable to deformations. They are thus well adapted to construct descriptor that would also be translation invariant. A two-dimensional spatial wavelet transform W is obtained by scaling by 2^j and rotating by r_θ a mother wavelet ψ ,

$$\psi_\lambda(u) = \psi_{j,\theta}(u) = 2^{-2j} \psi(2^{-j} r_\theta p) \quad (1)$$

In the remainder of this paper we restrict to Morlet wavelet transforms defined on $\Lambda = G \times \llbracket 0, J \rrbracket$ where G is a finite group of rotations of cardinal L and where the wavelet is taken at scale J ,

$$W_J \mathbf{x} = \{\mathbf{x} * \phi_J(u); \mathbf{x} * \psi_\lambda(u)\}_{p \in \mathbb{R}^2, \lambda \in \Lambda} \quad (2)$$

While the averaging part ϕ_J of the wavelet transform is invariant to translations, the high frequency part ψ_λ is covariant to them [9]. Invariance within a limited range inferior 2^J can be achieved by averaging the positive envelope with a smooth window,

$$S_J[\lambda] \mathbf{x}(u) = |\mathbf{x} * \psi_\lambda| * \phi_J(u) \quad (3)$$

Such non-linearised averaged wavelet coefficients are used in various form in computer vision (SIFT [17], DAISY [18]), but the scattering transform proposes a new non-linearity as well as a layer based architecture.

2.2 Scattering convolutional network

While providing local translation invariance, the averaging convolution introduced in 3 also removes the spatial variability of the wavelet transform. SCNs cascade this wavelet modulus operator to recover the lost information and compute progressively more invariant descriptors. Let combine the wavelet transform and modulus operations into a single wavelet modulus operator,

$$\mathcal{U}_J \mathbf{x} = \{S[\emptyset] \mathbf{x}; U[\lambda] \mathbf{x}\}_{\lambda \in \Lambda_J} = \{\mathbf{x} * \phi_J; |\mathbf{x} * \psi_\lambda|\}_{\lambda \in \Lambda_J}, \quad (4)$$

A scattering transform can be interpreted as a CNN [19] illustrated in Figure 1 which propagates a signal \mathbf{x} across multiple layers of the network and which outputs at each layer m the scattering invariant coefficients $S[p_m] \mathbf{x}$ where $p_m = (\lambda_1 \dots \lambda_m)$ is a path of m orientations and scales.

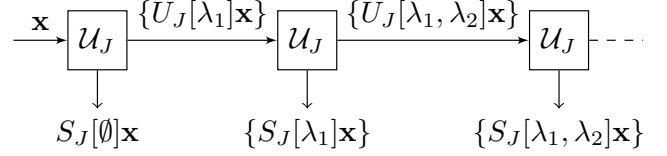


Fig. 1. Scattering networks can be seen as neural networks iterating over wavelet modulus operators \mathcal{U}_j . Each layer m outputs the averaged invariants $S[p_m] \mathbf{x}$ and covariant coefficients $U[p_{m+1}] \mathbf{x}$.

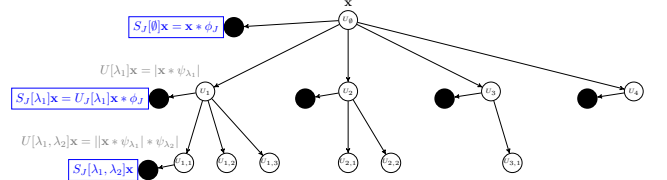


Fig. 2. Frequency decreasing scattering convolution network with $J = 4$, $L = 1$ and $M = 2$. A node i at scale j_i generates $(j_i - 1) \times L$ nodes.

The scattering energy is mainly concentrated along frequency decreasing paths, i.e. for which $|\lambda_{k+1}| \leq |\lambda_k|$ [9]. The energy contained in the other paths is negligible and thus for applications only frequency decreasing paths are considered. Moreover there exist a path length $M > 0$ after which all longer paths can be neglected. For signal processing applications, this decay appears to be exponential. And for classification applications, paths of length $M = 3$, i.e. two convolutions, provides the most interesting results [10, 20].

This restrictions yield an easier parametrization of a scattering network. Indeed its now completely defined by the mother wavelet ϕ , the maximum path length considered M , the finest scale level considered J and the number of orientation considered L .

Hence for a given set of parameter (ψ, M, J, L) , let $ST_{(\psi, M, J, L)}(\mathbf{x})$ denotes the unique frequency decreasing windowed scattering convolutional network with those parameters evaluated for signal \mathbf{x} . Each node i of this network generates a -possibly empty- set of nodes of size $(j_i - 1) \times L$ where j_i is the scale of node i and L is the number of orientations considered and it has the architecture displayed by Figure 2.

2.3 Scattering convolutional classifier

In the original framework [16], the scattering network $ST_{(\psi, M, J, L)}(\cdot)$ is used for classification task using a SVM classifier on the outputs of the network. Performance can be slightly improved by adding a feature selection step performing PCA on the scattering coefficients and keeping only the most informative ones. This classification framework provides results comparable with the state of the art on several datasets [10].

Training samples	MNIST		KTH-TIPS	
	SCHMT	SCN+SVM	SCHMT	SCN+SVM
2	30.0%	40.0%	40.0%	40.0%
	30.0%	40.0%	40.0%	40.0%
5	45.0%	56.5%	40.0%	40.0%
	45.5%	43.2%	40.0%	40.0%
10	30.0%	40.0%	40.0%	40.0%
	30.0%	40.0%	40.0%	40.0%
20	30.0%	40.0%	40.0%	40.0%
	30.0%	40.0%	40.0%	40.0%

Fig. 4. Classification score on 200 samples (top row) and on the full testset (bottom row) of MNIST and KTH-TIPS trained using only a limited number of training points per class

In this context the MAP algorithm [22] aims at finding the optimal hidden tree $\hat{h}_0^{new} = (\hat{h}_0^{new} \dots \hat{h}_{I-1}^{new})$ maximizing the probability of this sequence given the model’s parameters $P(\hat{h}_0 = \hat{h}_0^{new} | \mathcal{T}^{new}, \Theta_c)$. The MAP framework also provides \hat{P} the value of this maximum.

The MAP algorithm can be used in a multi-class classification problem by training an SCHMT model per class and then when presented with a new realization \mathbf{x}^{new} comparing the probability of the MAP hidden tree provided by each model

4 Classification results

We compare the performance of SCHMT to those of a SCN combined to an SVM (SCN+SVM) on restrictions to small number of training examples of two standard datasets. For both datasets and classification methods we use a scattering transform with $M = 3$ orders, $J = 5$ scales, $L = 3$ orientations and a Morlet mother wavelet. The hidden Markov tree has $K = 2$ states and is using a mixture of Gaussian to describe the relationship between the scattering coefficients and the hidden states. For the SVM, the best parameters are selected by cross-validation.

4.1 MNIST

We first test SCHMTs on the digit classification dataset MNIST [23]. SCHMT and SCN+SVM are both trained on a limited number of training examples per class and tested on 200 test samples (20 per class) and on the full testset (10000 images). The best results for each models are displayed in Figure 4.

Even when the number of training points is extremely limited —i.e. 2 or 5— SVM can reached good score on a small testset. However the model learned has poor generalization’s properties and highly decrease in performance when the testset size is scaled up. As opposed to this, SCHMT displays better generalization as its score remains constant while increasing the number of test samples and finally outperforms SVMs. —45% compared to 40% for 5 training examples. Finally, as

expected, when the number of training samples grow large enough —i.e. 10 and onward— for the SVMs, SCN+SVMs reach both better maximum and average score.

We should also notice that EM algorithm performances are undermined by convergence to local minima issues [24]. This sometime yields to very poor learning quality for the SCHMT. However when convergence occurs correctly SCHMTs reach better performances than the best SVMs.

This experiment confirm the superiority of generative model for limited number of training points. It also highly some weakness of the SCHMTs as sometime convergence is problematic and the problem seems to get stronger when the number of training samples increases. However when convergence occurs correctly SCHMT provides great classification score for low number of training examples.

4.2 KTH Texture

Further experiments are run on the KTH-TIPS texture dataset [25]. We perform the classification using 5 training examples per class and perform the testing on the rest of the dataset, i.e. 155 samples per class. SVM’s parameters are optimized using cross-validation. Results are displayed in Figure 4.

This texture dataset is more challenging than MNIST since the intra-class variability is higher. However the consequences are limited thanks to the scattering transform invariance properties. Again SCHMT perfoms well when given few training examples. The generalization capacity of SCHMT are also good since we see no change in performace when testing it on the 1550 remaining examples.

5 Conclusion

SCHMTs propose a new framework to process a high-dimensional signal from the raw data to the prediction task. The scattering transform projects the data into a representational space of even higher dimensionality but of reduced volume along the invariants in the data. Then a probabilistic graphical model —hidden Markov tree— is used to fit a generative model to the distribution of the representation of the data. This yield a powerful tool combining the interesting properties of the scattering transform for signal representation and the representational power of the hidden Markov models. The modeled distribution can be used to perform efficient classification tasks even given low information. Even though this document considers only classification, a generative model is much more versatile than a simple —yet efficient— discriminative counterpart. Because they model the full distribution of the data they can express more complex relationships between the observed and the unknown variables than simple discrimination.

To enhance SCHMT and especially the chance of converging toward a good minima during the EM learning, we plan on developing a variational to learn the model parameters [26].

6 References

- [1] Ruslan Salakhutdinov and Geoffrey E Hinton, “Deep boltzmann machines,” in *International Conference on Artificial Intelligence and Statistics*, 2009, pp. 448–455.
- [2] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol, “Stacked denoising autoencoders : Learning useful representations in a deep network with a local denoising criterion,” *The Journal of Machine Learning Research*, vol. 11, pp. 3371–3408, 2010.
- [3] Yann LeCun and Yoshua Bengio, “Convolutional networks for images, speech, and time series,” *The handbook of brain theory and neural networks*, vol. 3361, no. 10, 1995.
- [4] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [5] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors,” *arXiv preprint arXiv :1207.0580*, 2012.
- [6] Brody Huval, Tao Wang, Sameep Tandon, Jeff Kiske, Will Song, Joel Pazhayampallil, Mykhaylo Andriluka, Royce Cheng-Yue, Fernando Mujica, Adam Coates, et al., “An empirical evaluation of deep learning on highway driving,” *arXiv preprint arXiv :1504.01716*, 2015.
- [7] Xavier Glorot and Yoshua Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *International conference on artificial intelligence and statistics*, 2010, pp. 249–256.
- [8] Shih-Ta Hsiang and John W Woods, “Embedded video coding using invertible motion compensated 3-d sub-band/wavelet filter bank,” *Signal Processing : Image Communication*, vol. 16, no. 8, pp. 705–724, 2001.
- [9] S. Mallat, “Group invariant scattering,” *Communications in Pure and Applied Mathematics*, vol. 65, no. 10, pp. 1331–1398, Oct. 2012.
- [10] Joan Bruna and Stéphane Mallat, “Classification with scattering operators,” *arXiv preprint arXiv :1011.3023*, 2010.
- [11] Laurent Sifre and Stéphane Mallat, “Rotation, scaling and deformation invariant scattering for texture discrimination,” in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. 2013, pp. 1233–1240, IEEE.
- [12] Li Fei-Fei, Rob Fergus, and Pietro Perona, “One-shot learning of object categories,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 28, no. 4, pp. 594–611, 2006.
- [13] A Jordan, “On discriminative vs. generative classifiers : A comparison of logistic regression and naive bayes,” *Advances in neural information processing systems*, vol. 14, pp. 841, 2002.
- [14] Matthew S Crouse, Robert D Nowak, and Richard G Baraniuk, “Wavelet-based statistical signal processing using hidden markov models,” *Signal Processing, IEEE Transactions on*, vol. 46, no. 4, pp. 886–902, 1998.
- [15] Nick Kingsbury, “Complex wavelets for shift invariant analysis and filtering of signals,” *Applied and computational harmonic analysis*, vol. 10, no. 3, pp. 234–253, 2001.
- [16] Joan Bruna, *Scattering representations for recognition*, Ph.D. thesis, Ecole Polytechnique X, 2013.
- [17] Michael Grabner, Helmut Grabner, and Horst Bischof, “Fast approximated sift,” in *Computer Vision—ACCV 2006*, pp. 918–927. Springer, 2006.
- [18] Simon Winder, Gang Hua, and Michael Brown, “Picking the best daisy,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. 2009, pp. 178–185, IEEE.
- [19] Edouard Oyallon and Stéphane Mallat, “Deep roto-translation scattering for object classification,” *arXiv preprint arXiv :1412.8659*, 2014.
- [20] Joakim Andén and Stéphane Mallat, “Multiscale scattering for audio classification,” in *ISMIR*, 2011, pp. 657–662.
- [21] Jean-Baptiste Durand, Paulo Goncalves, and Yann Guédon, “Computational methods for hidden markov tree models—an application to wavelet trees,” *Signal Processing, IEEE Transactions on*, vol. 52, no. 9, pp. 2551–2560, 2004.
- [22] Jean-Baptiste Durand and Paulo Gonçalves, “Statistical inference for hidden markov tree models and application to wavelet trees,” 2001.
- [23] Yann LeCun, “Personal webpage : Mnist,” <http://yann.lecun.com/exdb/mnist/>, 2016, Accessed : 18-01-2016.
- [24] Toad K Moon, “The expectation-maximization algorithm,” *Signal processing magazine, IEEE*, vol. 13, no. 6, pp. 47–60, 1996.
- [25] Eric Hayman and Barbara Caputo, “Personal webpage : Kth-tips dataset,” <http://www.nada.kth.se/cvap/databases/kth-tips/index.html>, 2016, Accessed : 18-01-2016.
- [26] Martin J Wainwright and Michael I Jordan, “Graphical models, exponential families, and variational inference,” *Foundations and Trends® in Machine Learning*, vol. 1, no. 1-2, pp. 1–305, 2008.