

Invariant representations of images for better learning

Muthuvel Murugan K. V. Subrahmanyam

Chennai Mathematical Institute - Department of Computer Science
Chennai Mathematical Institute, Chennai, India
`{muthu,kv}@cmi.ac.in`

Abstract. We study the problem of obtaining representations of images which are invariant to transformation of the image under rotations, towards improving supervised learning. We show that using simple ideas from group representation theory we get invariant representations of images. Off the shelf learning algorithms perform much better on such representations. We develop on ideas by Cohen and Welling [1] to construct these invariant representations.

1 Introduction

This paper is primarily concerned with the issue of obtaining representations of images which are invariant to geometric transformations, with a view to using such representations to improve supervised learning. While our inspiration comes from SIFT [2] and the rotation invariant and scale invariant representation of images which SIFT outputs, our solution is best seen in the framework of representation learning, as articulated by [3, 4].

A key insight of deep learning seems to be that learning different *representations* of the data helps to improve the performance of machine learning algorithms. This has led to the development of RBM's [5] and autoencoders [6] modules, used to produce intermediate representations of data invariant to transformations. The problem we consider is much more modest, in a supervised set up and when the transformations are simple rotations of images. While this problem is significantly easier, we show that very simple ideas from group representation theory yields invariant representations of data and traditional supervised learning algorithms working on these representations give state of the art results or better.

Our methods are inspired largely by the ideas in Cohen and Welling [1]. In [1], the authors present a probabilistic model to learn compact commutative Lie groups and use that to produce invariant and disentangled representations of data. We follow their approach and produce a *more natural* invariant representation of the data. This representation is robust and gives 95%-96.5% accuracy on a number of off the shelf supervised learning algorithms.

The method of Cohen and Welling [1] is based on a simple but non-trivial idea from classical representation theory - when a reductive group acts on a vector space, the vector space splits into a direct sum of irreducible representations

i.e. there is a change of basis of the underlying vector space in which the matrix describing how a particular group element acts on the vector space is block diagonal. The block sizes are identical for all the group elements. The invariant representation we obtain is also based on a very simple idea, which reveals itself because one knows explicitly the block diagonal matrices which constitute building blocks of all representations.

So the problem of obtaining invariant representations of images splits into two subproblems - find the change of basis which reveals this block diagonal split, and then use the various invariant subspaces to obtain an invariant representation. Both these steps are also addressed in [1]. Our solution to the first step is different from theirs since we work in the simpler supervised set up. We use the learned conjugating matrix to obtain a natural invariant representation of the image. Due to constraints on the page limit of submissions we give very few details. The complete version of this paper with a relaxed introduction to representation theory is available with the authors.

2 Preliminaries on Representation theory and Cohen & Welling's framework

The problem we consider is classifying images under rotations i.e. given a collection of images and their rotated versions, we would like to classify test images into one of classes seen during the training phase. We will assume that images are rotated by an angle in $[0, 2\pi]$ and that an image rotated by 2π is itself. Assume that the vectorized image is D dimensional.

One regards the action of the rotation group G on images as a group homomorphism \mathcal{R} from G to $GL(D)$. Here $GL(D)$ is the group of invertible linear transformations of the D -dimensional real vector space.

As Cohen and Welling observe, the group of rotations is a one parameter compact, commutative subgroup of the group $SO(D)$, the special orthogonal group of invertible, orthogonal $D \times D$ matrices. It is well known that there is a $D \times D$ matrix W such that for all $g \in G$, conjugating $\mathcal{R}(g)$ by W yields a block diagonal matrix

$$R(\theta) = \begin{bmatrix} \cos(\omega_0\theta) & -\sin(\omega_0\theta) & \dots & 0 & 0 \\ \sin(\omega_0\theta) & \cos(\omega_0\theta) & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \cos(\omega_{D/2-1}\theta) & -\sin(\omega_{D/2-1}\theta) \\ 0 & 0 & \dots & \sin(\omega_{D/2-1}\theta) & \cos(\omega_{D/2-1}\theta) \end{bmatrix} \quad (1)$$

Here θ is the parameter describing the rotation the image undergoes and ω_j 's are integers, since we require that $R(\theta) = R(\theta + 2\pi)$.

For compact commutative subgroups of $SO(D)$, given a collection of images X and versions Y obtained by hitting X with an element of the subgroup, in [1] the authors derive an expression for the marginal likelihood of $Y|X$. Surprisingly they also give an elegant closed form solution to the gradient of the marginal

likelihood of $Y|X$ with respect to W , and use a gradient descent algorithm to compute W . After learning W they learn the parameters ω_j . This is then used for classification via manifold distance.

Our concern is more to do with representation learning i.e. to obtain representations of images which are invariant to rotations, more in the spirit of the papers [3, 4, 7]. The ideas in Cohen and Welling lend to this quite naturally.

Like in [1] we view the problem in two steps. In the first step we assume that we are given image pairs along with the angle of rotation, $(x^{(k)}, y^{(k)}, \theta^{(k)})$, where $y^{(k)}$ is the image $x^{(k)}$ rotated by some angle $\theta^{(k)}$.

$$y^{(k)} = WR(\theta^{(k)})W^T x^{(k)} + \eta \quad (2)$$

In this step we learn the conjugating W . This can be done exactly as Cohen and Welling suggest but we give a different algorithm in section 3. Then we address the problem of classifying images under rotations. We use the W identified in step 1 to project the image onto the irreducible subspaces of the given representation of the one parameter group of rotations, and show how this can be used to get an invariant representation of images. We then use this invariant representation to classify in a supervised set up.

3 Finding W and the invariant representation

ω_j being an integer, we may assume it takes values from $\{0, 1, \dots, 359\}$. Instead of working with Equation 2, we work with the following:

$$W^T y^{(k)} = R(\theta^{(k)})W^T x^{(k)} \quad (3)$$

We find W by finding pairs of columns at a time. We start this by picking a particular integer value in the range $\{0, 1, \dots, 359\}$ for ω_j . For a particular choice of ω_j , Equation 3 can be rewritten as

$$W_{(2j:2j+2)}^T y^{(k)} = R_j(\theta^{(k)})W_{(2j:2j+2)}^T x^{(k)} \quad (4)$$

where

$$R_j(\theta^{(k)}) = \begin{bmatrix} \cos(\omega_j \theta^{(k)}) & -\sin(\omega_j \theta^{(k)}) \\ \sin(\omega_j \theta^{(k)}) & \cos(\omega_j \theta^{(k)}) \end{bmatrix}$$

Expanding Equation 4 gives a set of simultaneous equations for every sample (k) . We can solve this using SVD, or by other well known techniques.

We then use the W and ω_j obtained above to get invariant representations of images. We experiment with 3 different representations all of which are invariant to rotation.

3.1 Invariant feature vector representation 1

Given an image x to get a feature vector we project x onto the 2-D subspaces naturally indexed by W , and take only the norms of the projection on each of the subspace. This is also reported in [1].

The feature vector given by x is $(\sqrt{(w_0^T x)^2 + (w_1^T x)^2}, \sqrt{(w_2^T x)^2 + (w_3^T x)^2}, \dots)$. From Equation 4 it is clear that this representation invariant to rotation.

3.2 Invariant Image representation

Given any image x consider $I(x)$ given by

$$I(x) = \sum_{j=1}^{\omega_{max}} \sum_{k=0}^{l_j} [W_{2k:2k+2}^{\omega=\omega_j}] \begin{bmatrix} r_k^{x,\omega_j} \cos(\alpha_k^{x,\omega_j} - \alpha_0^{x,\omega_j}) \\ r_k^{x,\omega_j} \sin(\alpha_k^{x,\omega_j} - \alpha_0^{x,\omega_j}) \end{bmatrix} + [W^{\omega=\omega_0}][W^{\omega=\omega_0}]^T x$$

where $[W_{2k:2k+2}^{\omega=\omega_j}]$ picks the columns of W corresponding to $\omega = \omega_j$ and then picks the columns $2k$ and $2k + 1$,

$$\alpha_k^{x,\omega_j} = \tan^{-1} \left(\frac{[W_{2k+1}^{\omega=\omega_j}]^T x}{[W_{2k}^{\omega=\omega_j}]^T x} \right)$$

and

$$r_k^{x,\omega_j} = \sqrt{([W_{2k}^{\omega=\omega_j}]^T x)^2 + ([W_{2k+1}^{\omega=\omega_j}]^T x)^2}$$

In other words fix a ω_j . Suppose the subspace has real dimension $2l_j$ for an integer l_j . Consider the projection of the image onto the l_j , 2-dimensional subspaces. Each such projection gives a point (x_k, y_k) or $(r_k \cos \alpha_k, r_k \sin \alpha_k)$ in polar coordinates. Let us view all the totality of these points T as sitting in a single two dimensional space, say the standard two dimensional plane. It is clear that under a rotation of the image, the point $(r_k \cos \alpha_k, r_k \sin \alpha_k)$ is transformed to $(r_k \cos(\alpha_k + \omega_j \theta), r_k \sin(\alpha_k + \omega_j \theta))$. Viewing all these points in the standard two dimensional plane, all that has happened is that each point in T has moved by an angle of $\omega_j \theta$. So the pairwise angles between these points (measured at the origin) has remained the same. Of course the norm of each point is also invariant. So the totality of norms plus pairwise angles is invariant. We exploit this invariance, by taking the point (x_0, y_0) and align the axes so that with respect to the new axes this point lies along the x -direction. This is akin to what is done in SIFT, where the direction of maximum gradient is chosen always as the reference direction.

This now describes an arrangement of points which is invariant with respect to rotations. We have

Theorem 1 *The representation $I(\cdot)$ is invariant to rotation.*

3.3 Invariant feature vector representation 2

In the previous section, we reconstructed a image from the projections. Here we will use the projections after realignment as is. We will not go back to the pixel space, instead we will work in the projected space. This gives a significant *dimension reduction*. In our experiments, we look at the average norm of each 2-D projections, and eliminated those subspaces, where the average norm of the projected vectors is less than a threshold. For the dataset we consider, the dimension reduced by a factor of 8 with almost no change in accuracy.

4 Experiments and results

To find W and ω_j , we generate samples $x^{(k)}$ by choosing each pixel from a standard normal distribution. The rotational angles $\theta^{(k)}$ are picked uniformly from $[0, 360)$ and the images are rotated to get $y^{(k)}$. We generate 5000 samples. Starting with $\omega_j = 0$ we find columns of W . In our experiments we stopped at $\omega_j = 47$ and obtained 294, 2-dimensional subspaces giving W of size 784×588 .

For classification we used MNIST_rot¹. This dataset consists of rotated handwritten digits of size 28×28 pixels with 12000 training and 50000 test samples.

We trained the different classifiers on the feature vectors and invariant image as explained in Section 3.1, 3.3 and 3.2 for the dataset MNIST-rot and tabulated (Table 1) the results. For comparison we also ran the classifiers on the pixel space, as is. For the feature vector described in 3.3, we set a threshold of 0.3 for the average norm of the projections onto each 2-D subspace and got feature vectors of dimension 108, giving us a factor 8 (almost) reduction in dimension. Our neural network has 2 layers with 4000 and 100 neurons and we use a dropout probability of 0.3. The CNN we use is the one given in tensorflow² tutorial.

Figure 1 shows images and their invariant representations as described in Section 3.2.

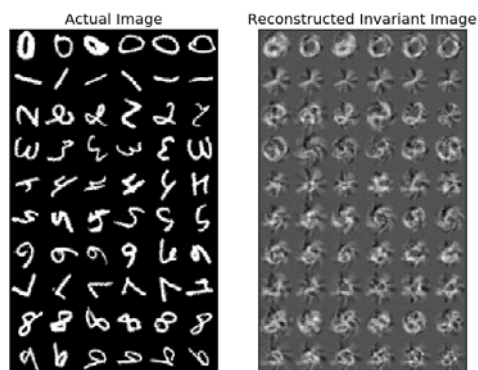


Fig. 1: Different handwritings and their invariant representations

In their 2016 paper [8], Cohen and Welling describe G-CNN an equivariant Convolutional neural network. For MNIST-rot they get an accuracy of 97.72%.

¹<http://www.iro.umontreal.ca/~lisa/twiki/bin/view.cgi/Public/MnistVariations>

²<https://www.tensorflow.org/versions/r0.9/tutorials/mnist/pros/index.html>

Method	Accuracy (%)			
	Actual Pixels	Invariant Vector 1	Invariant Vector 2	Reconstructed inv. image
kNN	83.89	87.97	94.29	94.29
SVD	73.56	81.88	91.05	91.05
SVM	82.00	83.37	95.26	95.72
Neural Networks	89.61	92.27	96.36	96.28
CNN	93.14	-NA-	-NA-	96.50

Table 1: MNIST-rot : Comparison of different methods

5 Conclusion

The results suggest that irreducible representations of images under group actions can yield interesting invariant representations of images, which could help in dimension reduction and better classification. As mentioned in all recent papers dealing with image transformations [8, 9] we need a way to understand and capture actions of groups that are not commutative. We are experimenting with extending the current work to the action of $SL(2)$ and $SU(2)$.

References

- [1] Taco S. Cohen and Max Welling. Learning the Irreducible Representations of Commutative Lie Groups. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pages 1755–1763, February 2014.
- [2] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110, November 2004.
- [3] Ian Goodfellow, Honglak Lee, Quoc V. Le, Andrew Saxe, and Andrew Y. Ng. Measuring invariances in deep networks. In *Advances in Neural Information Processing Systems 22*, pages 646–654. Curran Associates, Inc., 2009.
- [4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–828, August 2013.
- [5] Max Welling, Michal Rosen-Zvi, and Geoffrey E. Hinton. Exponential family harmoniums with an application to information retrieval. In *Advances in Neural Information Processing Systems 17 [NIPS]*, pages 1481–1488, 2004.
- [6] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, pages 1096–1103, 2008.
- [7] Fabio Anselmi, Joel Z. Leibo, Lorenzo Rosasco, Jim Mutch, Andrea Tacchetti, and Tomaso A. Poggio. Unsupervised learning of invariant representations in hierarchical architectures. *CoRR*, abs/1311.4158, 2013.
- [8] Taco S. Cohen and Max Welling. Group Equivariant Convolutional Networks. *Proceedings of The 33rd International Conference on Machine Learning*, 48, feb 2016.
- [9] Robert Gens and Pedro M Domingos. Deep symmetry networks. In *Advances in Neural Information Processing Systems 27*, pages 2537–2545. Curran Associates, Inc., 2014.