Understanding Deep Image Representation by Inverting Them

Aravindh Mahendran Andrea Vedaldi University of Oxford

Presented by Xinmeng Li and Beibei Du Mar 28, 2017

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo



Introduction







Framework to do visualizations

Image space (pre-images)

Representation space



Framework to invert representations



Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Contributions

- Propose a general method to invert representations, including SIFT, HOG and CNNS
- Inverting representations method performs better for DSIFT and HOG compared to recent alternatives.
- Apply inversion technique to the analysis of deep CNNs, showing that CNN gradually builds an increasing amount of invariance, layer by layer
- Study locality of information stored in the representations by reconstructing image from selected groups of neurons, either spatially or by channel.

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Related Work #1 SIFT and Dense-SIFT

Reconstruct an image given keypoint of SIFT based on a huge database







$$\rightarrow R_i = \{\mathbf{v}_i, \mathbf{x}_i, s_i, o_i, A_i\}$$

Appearance descriptor vi; Coordinates of region's center Xi;

Scale Si;

Dominant gradient orientation Oi; Region support information Ai.

Image from: P. Weinzaepfel, H. J'egou, and P. P'erez. Reconstructing an image from its local descriptors. In CVPR, 2011.

Related Work #1 Dense-SIFT



Figure 4. From left to right: the original picture and the reconstruction before and after completion of uncovered regions.

1. For each query appearance descriptor v_i , search its nearest neighbor in the descriptor database

$$j^* = \arg \max_{j \in \{1 \cdots m\}} \|\mathbf{v}_i - \mathbf{v}_j\|_2,$$
(1)

and recover the corresponding elliptic image patch

$$\mathbf{q}_j^* = \mathbf{I}_{k(j^*)}(S_{j^*}). \tag{2}$$

- 2. Seamlessly stitch all patches \mathbf{p}_i , $i = 1 \cdots n$, together to obtain a partial reconstruction with support $S = \bigcup_{i=1}^n S_i \subset \Omega$ (see details below).
- 3. Complete remaining empty zone $\overline{S} = \Omega \setminus S$ by smooth interpolation, as shown in Figure 4 (see details below).

Related Work #1 HOG and HOGgles



Related Works #2

- Visualizing and understanding convolutional networks(DeConvNet). In ECCV, 2014.
 - Backtrack the network computations to identify which image patches are responsible for certain neural activations

• Visualising image classification models and saliency maps.

ICLR, 2014

• Reconstruction via maximizing class-neuron scores

- Visualizing and understanding convolutional networks(DeConvNet).
 - 1. Feed image into net



2. Pick a layer, set the gradient there to be all zero except for one 1 for

some neuron of interest 3. Backprop to image:



"Guided backpropagation:" only propagate positive gradients



 Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n



[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]

[Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al., 2014] [Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]



comp150dl 😫 Tufts

13

Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

2. Visualize the Data gradient:

(note that the gradient on data has three channels. Here they visualize M, s.t.:

 $M_{ij} = \max_c |w_{h(i,j,c)}|$

(at each pixel take abs val, and max over channels)





26

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Inverting Method



Inverting Method

• To solve equation

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Loss: compare image representation and target one

R: regulariser to capture a natural image prior

- Given:
 - Representation function $\Phi : \mathbb{R}^{H \times W \times C} \to \mathbb{R}^{d}$
 - Representation: $\Phi_0 = \Phi(\mathbf{x}_0)$
- Find: x* that minimizes the objective

Inverting method - loss

Loss Function: Euclidean distance $\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$ Objective equation: $\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$ Code: # Let's suppose we are using AlexNet and the representation is obtained in its 3rd fully controls

h is the representation of image x given by the activations of the 3rd fc layer of the CNN.

 $h = cnn_fc3(x)$



Code: http://cesarsalgado.com/ipython_notes_on_understanding_deep_image_representations_by_inverting_them/

Inverting method - regularisers

Regularisers: *a*-norm image prior

• *a*-norm (*a* = 6) (X is vectorised and mean-subtracted image) $\mathcal{R}_{\alpha}(\mathbf{x})$

 $\mathcal{R}_{\alpha}(\mathbf{x}) = \|\mathbf{x}\|_{\alpha}^{\alpha}$



- Example:
 - Assume image x with shape(H,W,C) H = 2, W = 2, C = 3(RGB)

```
>>> x = np.random.randn(2,2,3)
>>> x
array([[[ 0.47451494,  0.29584752,  2.03439944],
        [-2.14093253,  0.68461236,  0.91091761]],
        [[ 1.06140708,  0.28991584, -0.30025526],
        [-0.34013064, -0.28461196,  0.84689234]]])
```

Then a-norm will be: [>>> np.sum(x**6) 169.68199780729023

Code: http://cesarsalgado.com/ipython_notes_on_understanding_deep_image_representations_by_inverting_them/

Inverting method - regularizes

Regularisers: TV-norm image prior

Ο

TV(total variation)
$$\mathcal{R}_{V^{eta}}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{rac{eta}{2}}$$



• Example: total_variation(x) will be with shape(3,)

```
>>> total(x)
array([ 7.18500788, 0.15117329, 6.71282397])
24
```

Inverting method - regularizes

Objective function: argmin $\|\Phi(\mathbf{x}) - \Phi_0\|^2 + \lambda_{\alpha} \mathcal{R}_{\alpha}(\mathbf{x}) + \lambda_{V^{\beta}} \mathcal{R}_{V^{\beta}}(\mathbf{x})$

Regulariser defined as sum of both subterms

Code:



Example: regularizer result will be with shape(3,)

```
[>>> regularizer(x)
(3,)
array([ 4.41173195e+10, 4.41173194e+10, 4.41173195e+10])
```

Inverting method - loss and regulariser

• Balancing loss and regularizer

$$\|\Phi(\sigma \mathbf{x}) - \Phi_0\|_2^2 / \|\Phi_0\|_2^2 + \lambda_\alpha \mathcal{R}_\alpha(\mathbf{x}) + \lambda_{V^\beta} \mathcal{R}_{V^\beta}(\mathbf{x})$$

- Normalized Loss
- σ : average Euclidean norm of natural images in a training set
- B = 128 *σ*x ∈[-B,B]
- $\mathcal{R}_{\alpha}(\mathbf{x}) \approx HWB^{\alpha}/\sigma^{\alpha} \quad \lambda_{\alpha} \approx \sigma^{\alpha}/(HWB^{\alpha})$
- $\lambda_{V^{\beta}} \approx \sigma^{\beta}/(HW(aB)^{\beta})$ (2.16 * 10^8)

Inverting method - Optimisation

- Gradient Descent(GD)
- GD extensions
 - Momentum



- Computing derivatives
 - CNNs : backpropagation
 - CNN-HOG and CNN-DSIFT: same as CNNs

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Representations - CNN-A(Caffe-Alex) deep networks

CNN-A structure



layer	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
name	conv1	relul	mpool1	norm1	conv2	relu2	mpool2	norm2	conv3	relu3	conv4	relu4	conv5	relu5	mpool5	fc6	relu6	fc7	relu7	fc8
type	cnv	relu	mpool	nrm	cnv	relu	mpool	nrm	cnv	relu	cnv	relu	cnv	relu	mpool	cnv	relu	cnv	relu	cnv
channels	96	96	96	96	256	256	256	256	384	384	384	384	256	256	256	4096	4096	4096	4096	1000
rec. field	11	11	19	19	51	51	67	67	99	99	131	131	163	163	195	355	355	355	355	355

Table 2. CNN-A structure. The table specifies the structure of CNN-A along with receptive field size of each neuron. The filters in layers from 16 to 20 operate as "fully connected": given the standard image input size of 227×227 pixels, their support covers the whole image. Note also that their receptive field is larger than 227 pixels, but can be contained in the image domain due to padding.

Representations - CNN-DSIFT and CNN-HOG

How DSIFT and HOG implemented in CNNs:

- Step1: Computing and binning image gradients
- Step2: Pooling binned gradients into cell histograms
- Step 3: Grouping cells into blocks (CNN layers)
- Step 4: Normalising the blocks (CNN layers)

Outline

- Introduction
- Contributions
- Related Works
- Inverting Method
- Representation
- Experiments
- Demo

Experiments with shallow representation (HOG & SIFT)

• Quantitative Analysis

descriptors	HOG	HOG	HOGb	DSIFT
method	HOGgle	our	our	our
error (%)	66.20	28.10	10.67	10.89
	± 13.7	± 7.9	± 5.2	±7.5

• Error:
$$\|\Phi(\mathbf{x}^*) - \Phi(\mathbf{x}_i)\|_2 / N_{\Phi}$$

Qualitative Analysis



Experiments with shallow representation (HOG & SIFT)

Comparison HOG and HOGle



descriptors	HOG	HOG	HOGb	DSIFT
method	HOGgle	our	our	our
error (%)	66.20	28.10	10.67	10.89

Table 1. Average reconstruction error of different representation inversion methods, applied to HOG and DSIFT. HOGb denotes HOG with bilinear orientation assignments. The standard deviation shown is the standard deviation of the error and not the standard deviation of the mean error.

Experiments with shallow representation (HOG & SIFT)

Comparison HOG, DSIFT(quantative)



descriptors	HOG	HOG	HOGb	DSIFT
method	HOGgle	our	our	our
error (%)	66.20	28.10	10.67	10.89
1 1 2 020	± 13.7	± 7.9	± 5.2	±7.5

Table 1. Average reconstruction error of different representation inversion methods, applied to HOG and DSIFT. HOGb denotes HOG with bilinear orientation assignments. The standard deviation shown is the standard deviation of the error and not the standard deviation of the mean error.

- Train Data: 1.2M images of the ImageNet ILSVRC 2012
- Validation Data: 100 ILSVRC validation images
- $\lambda_{\alpha} = 2.16 \times 10^8$
- Increasing $\lambda_{V^{\beta}}$ ten-folds, starting from 0.5

Quantitative Analysis

$\lambda_{V^{\beta}}$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	conv1	relu1	pool1	norm1	conv2	relu2	pool2	norm2	conv3	relu3	conv4	relu4	conv5	relu5	pool5	fc6	relu6	fc7	relu7	fc8
λ_1	10.0	11.3	21.9	20.3	12.4	12.9	15.5	15.9	14.5	16.5	14.9	13.8	12.6	15.6	16.6	12.4	15.8	12.8	10.5	5.3
	±5.0	± 5.5	± 9.2	± 5.0	± 3.1	±5.3	± 4.7	± 4.6	±4.7	± 5.3	± 3.8	± 3.8	± 2.8	±5.1	± 4.6	± 3.5	± 4.5	± 6.4	± 1.9	± 1.1
λ_2	20.2	22.4	30.3	28.2	20.0	17.4	18.2	18.4	14.4	15.1	13.3	14.0	15.4	13.9	15.5	14.2	13.7	15.4	10.8	5.9
1.000	±9.3	± 10.3	± 13.6	± 7.6	± 4.9	± 5.0	±5.5	±5.0	± 3.6	± 3.3	± 2.6	±2.8	± 2.7	± 3.2	± 3.5	± 3.7	± 3.1	± 10.3	± 1.6	± 0.9
λ_3	40.8 ±17.0	$45.2 \\ \pm 18.7$	54.1	48.1	39.7	32.8	32.7 ±8.0	32.4	25.6 ± 5.6	26.9 ±5.2	$23.3 \\ \pm 4.1$	23.9	25.7	20.1	19.0	18.6	18.7	17.1	15.5	8.5

Table 3. Inversion error for CNN-A. Average inversion percentage error (normalized) for all the layers of CNN-A and various amounts of V^{β} regularisation: $\lambda_1 = 0.5$, $\lambda_2 = 10\lambda_1$ and $\lambda_3 = 100\lambda_1$. In bold face are the error values corresponding to the regularizer that works best both qualitatively and quantitatively. The deviations specified in this table are the standard deviations of the errors and not the standard deviations of the mean error value.

- Qualitatively Analysis:
 - reconstruction for a image from each layer of CNN-A



Figure 6. CNN reconstruction. Reconstruction of the image of Fig. 5.a from each layer of CNN-A. To generate these results, the regularization coefficient for each layer is chosen to match the highlighted rows in table 3. This figure is best viewed in color/screen.

- Qualitatively Analysis:
 - reconstruction for a image from each layer of CNN-A



Figure 6. CNN reconstruction. Reconstruction of the image of Fig. 5.a from each layer of CNN-A. To generate these results, the regularization coefficient for each layer is chosen to match the highlighted rows in table 3. This figure is best viewed in color/screen.

• Reconstructions obtained from subset of representation in different CNN layers



Figure 9. CNN receptive field. Reconstructions of the image of Fig. 5.a from the central 5×5 neuron fields at different depths of CNN-A. The white box marks the field of view of the 5×5 neuron field. The field of view is the entire image for conv5 and relu5.

Weakness and Future Work

- Weakness
 - Generate an image with the feature representation that similar to the given representation, not the image similar to the given image
 - Error detect for comparing the original image representation and the inverse image representation
 - Optimization implements at test time, requires computed gradient of feature representation, makes it relatively slow
 - Implementation on only AlexNet, but not other state-of-art CNN
 - "Spikes" generated by using regularizer

• Future Work

- Implementation on other state-of-art CNN
- Experiment with more expressive natural image priors
- Extract subsets of neurons that encode object parts and try to establish sub-networks that capture different details of the image

Recent work

Inverting visual representations with convolutional networks

Dosovitskiy, Alexey, and Thomas Brox

Visualizing Deep Convolutional Neural Networks Using Natural Pre-images

Aravindh Mahendran, Andrea Vedald

http://link.springer.com/article/10.1007/s11263-016-0911-8





Original video:

http://techtalks.tv/talks/understanding-deep-image-representations-by-inverting-them/61629/