# Visualizing a Trained Autoencoder

## From Ufldl

Having trained a (sparse) autoencoder, we would now like to visualize the function learned by the algorithm, to try to understand what it has learned. Consider the case of training an autoencoder on $10 \times 10$ images, so that $n = 100$. Each hidden unit $i$ computes a function of the input:

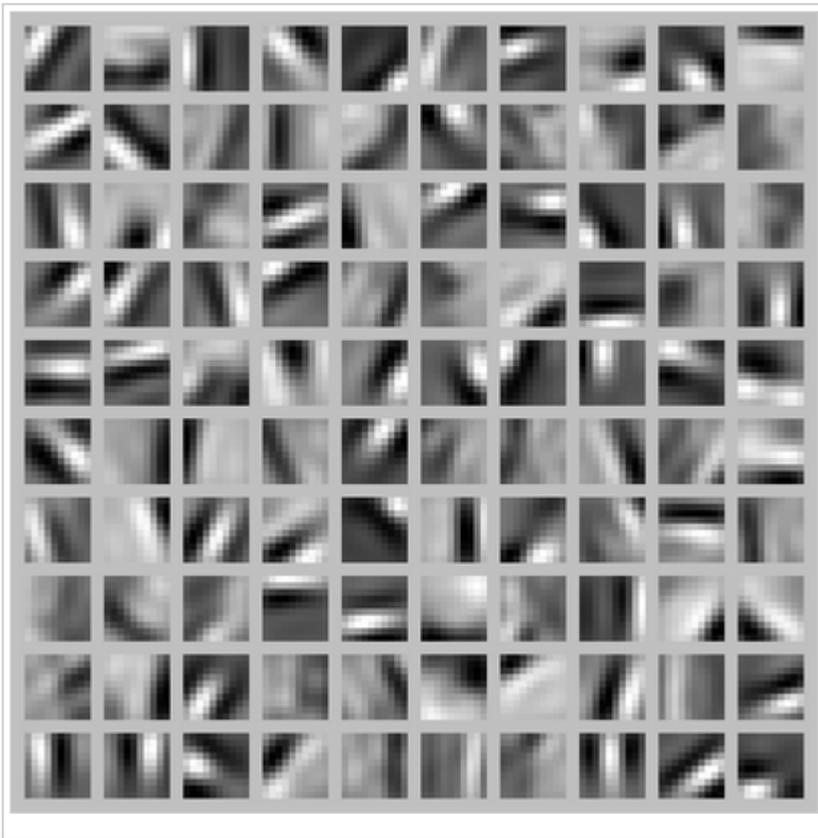$$a_i^{(2)} = f\left(\sum_{j=1}^{100} W_{ij}^{(1)} x_j + b_i^{(1)}\right).$$

We will visualize the function computed by hidden unit $i$---which depends on the parameters $W_{ij}^{(1)}$ (ignoring the bias term for now)---using a 2D image. In particular, we think of $a_i^{(2)}$ as some non-linear feature of the input $x$. We ask: What input image $x$ would cause $a_i^{(2)}$ to be maximally activated? (Less formally, what is the feature that hidden unit $i$ is looking for?) For this question to have a non-trivial answer, we must impose some constraints on $x$. If we suppose that the input is norm constrained by $||x||^2 = \sum_{i=1}^{100} x_i^2 \leq 1$, then one can show (try doing this yourself) that the input which maximally activates hidden unit $i$ is given by setting pixel $x_j$ (for all 100 pixels, $j = 1, \ldots, 100$) to

$$x_j = \frac{W_{ij}^{(1)}}{\sqrt{\sum_{j=1}^{100} (W_{ij}^{(1)})^2}}.$$

By displaying the image formed by these pixel intensity values, we can begin to understand what feature hidden unit $i$ is looking for.

If we have an autoencoder with 100 hidden units (say), then we our visualization will have 100 such images---one per hidden unit. By examining these 100 images, we can try to understand what the ensemble of hidden units is learning.

When we do this for a sparse autoencoder (trained with 100 hidden units on 10x10 pixel inputs[1] we get the following result:

Each square in the figure above shows the (norm bounded) input image $x$ that maximally actives one of 100 hidden units. We see that the different hidden units have learned to detect edges at different positions and orientations in the image.

These features are, not surprisingly, useful for such tasks as object recognition and other vision tasks. When applied to other input domains (such as audio), this algorithm also learns useful representations/features for those domains too.

[1] *The learned features were obtained by training on **whitened** natural images. Whitening is a preprocessing step which removes redundancy in the input, by causing adjacent pixels to become less correlated.*

Neural Networks | Backpropagation Algorithm | Gradient checking and advanced optimization | Autoencoders and Sparsity | **Visualizing a Trained Autoencoder** | Sparse Autoencoder Notation Summary | Exercise:Sparse Autoencoder

- This page was last modified on 7 April 2013, at 12:49.