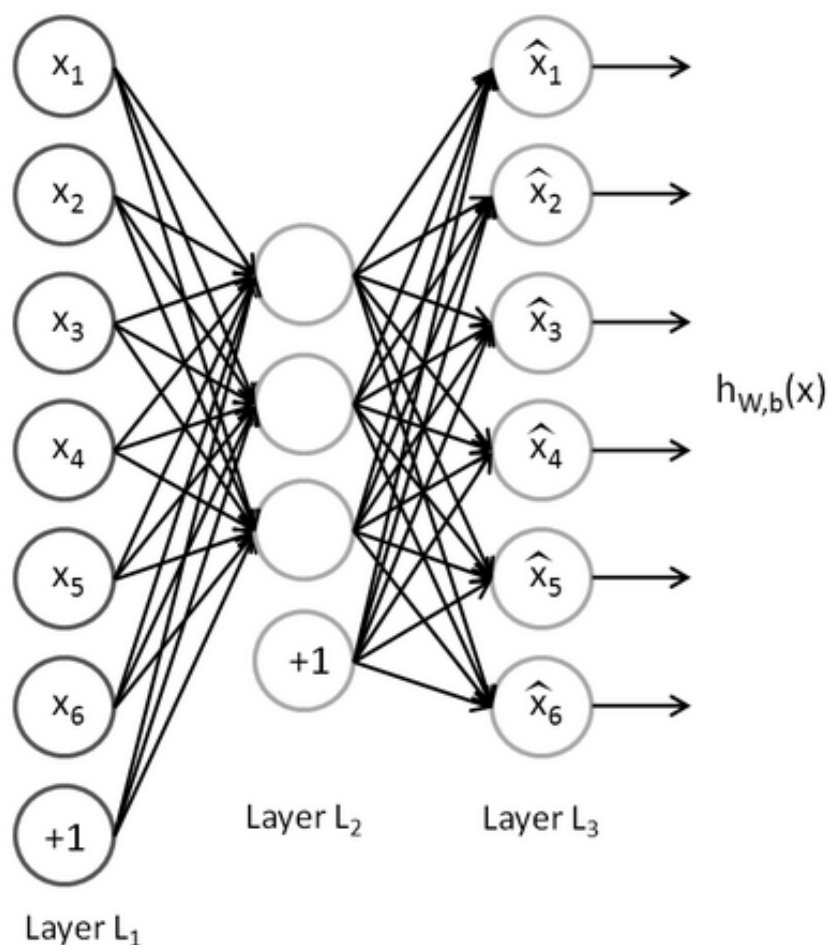


# Autoencoders and Sparsity

From Ufidl

So far, we have described the application of neural networks to supervised learning, in which we have labeled training examples. Now suppose we have only a set of unlabeled training examples  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots\}$ , where  $x^{(i)} \in \mathbb{R}^n$ . An autoencoder neural network is an unsupervised learning algorithm that applies backpropagation, setting the target values to be equal to the inputs. I.e., it uses  $y^{(i)} = x^{(i)}$ .

Here is an autoencoder:



The autoencoder tries to learn a function  $h_{W,b}(x) \approx x$ . In other words, it is trying to learn an approximation to the identity function, so as to output  $\hat{x}$  that is similar to  $x$ . The identity function seems a particularly trivial function to be trying to learn; but by placing constraints on the network, such as by limiting the number of hidden units, we can discover interesting structure about the data. As a concrete example, suppose the inputs  $x$  are the pixel intensity values from a  $10 \times 10$  image (100 pixels) so  $n = 100$ , and there are  $s_2 = 50$  hidden units in layer  $L_2$ . Note that we also have  $y \in \mathbb{R}^{100}$ . Since there are only 50 hidden units, the network is forced to learn a compressed representation of the input. I.e., given only the vector of hidden unit activations  $a^{(2)} \in \mathbb{R}^{50}$ , it must try to reconstruct the 100-pixel input  $x$ . If the input were completely random—say, each  $x_i$  comes from an IID Gaussian independent of the other features—then this compression task

would be very difficult. But if there is structure in the data, for example, if some of the input features are correlated, then this algorithm will be able to discover some of those correlations. In fact, this simple autoencoder often ends up learning a low-dimensional representation very similar to PCAs.

Our argument above relied on the number of hidden units  $s_2$  being small. But even when the number of hidden units is large (perhaps even greater than the number of input pixels), we can still discover interesting structure, by imposing other constraints on the network. In particular, if we impose a sparsity constraint on the hidden units, then the autoencoder will still discover interesting structure in the data, even if the number of hidden units is large.

Informally, we will think of a neuron as being "active" (or as "firing") if its output value is close to 1, or as being "inactive" if its output value is close to 0. We would like to constrain the neurons to be inactive most of the time. This discussion assumes a sigmoid activation function. If you are using a tanh activation function, then we think of a neuron as being inactive when it outputs values close to -1.

Recall that  $a_j^{(2)}$  denotes the activation of hidden unit  $j$  in the autoencoder. However, this notation doesn't make explicit what was the input  $\mathbf{x}$  that led to that activation. Thus, we will write  $a_j^{(2)}(\mathbf{x})$  to denote the activation of this hidden unit when the network is given a specific input  $\mathbf{x}$ . Further, let

$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m [a_j^{(2)}(x^{(i)})]$$

be the average activation of hidden unit  $j$  (averaged over the training set). We would like to (approximately) enforce the constraint

$$\hat{\rho}_j = \rho,$$

where  $\rho$  is a sparsity parameter, typically a small value close to zero (say  $\rho = 0.05$ ). In other words, we would like the average activation of each hidden neuron  $j$  to be close to 0.05 (say). To satisfy this constraint, the hidden unit's activations must mostly be near 0.

To achieve this, we will add an extra penalty term to our optimization objective that penalizes  $\hat{\rho}_j$  deviating significantly from  $\rho$ . Many choices of the penalty term will give reasonable results. We will choose the following:

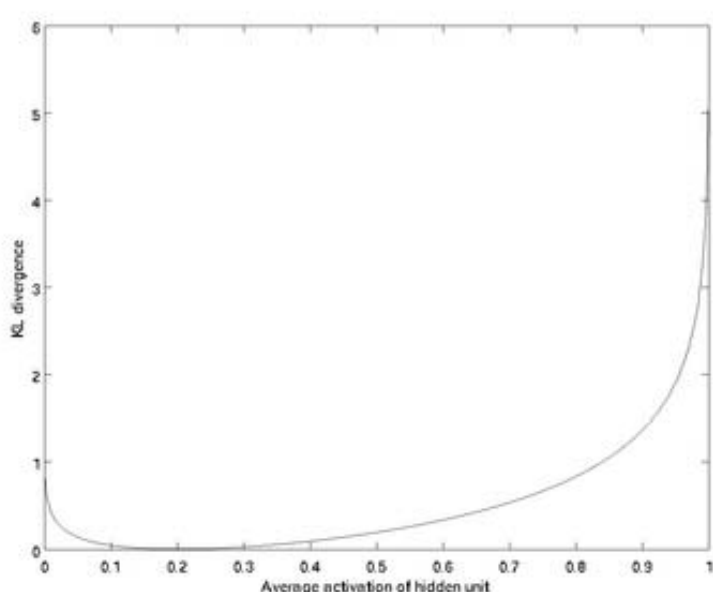
$$\sum_{j=1}^{s_2} \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}.$$

Here,  $s_2$  is the number of neurons in the hidden layer, and the index  $j$  is summing over the hidden units in our network. If you are familiar with the concept of KL divergence, this penalty term is based on it, and can also be written

$$\sum_{j=1}^{s_2} \text{KL}(\rho || \hat{\rho}_j),$$

where  $\text{KL}(\rho || \hat{\rho}_j) = \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1-\rho}{1-\hat{\rho}_j}$  is the Kullback-Leibler (KL) divergence between a Bernoulli random variable with mean  $\rho$  and a Bernoulli random variable with mean  $\hat{\rho}_j$ . KL-divergence is a standard function for measuring how different two different distributions are. (If you've not seen KL-divergence before, don't worry about it; everything you need to know about it is contained in these notes.)

This penalty function has the property that  $\text{KL}(\rho || \hat{\rho}_j) = 0$  if  $\hat{\rho}_j = \rho$ , and otherwise it increases monotonically as  $\hat{\rho}_j$  diverges from  $\rho$ . For example, in the figure below, we have set  $\rho = 0.2$ , and plotted  $\text{KL}(\rho || \hat{\rho}_j)$  for a range of values of  $\hat{\rho}_j$ :



We see that the KL-divergence reaches its minimum of 0 at  $\hat{\rho}_j = \rho$ , and blows up (it actually approaches  $\infty$ ) as  $\hat{\rho}_j$  approaches 0 or 1. Thus, minimizing this penalty term has the effect of causing  $\hat{\rho}_j$  to be close to  $\rho$ .

Our overall cost function is now

$$J_{\text{sparse}}(W, b) = J(W, b) + \beta \sum_{j=1}^{s_2} \text{KL}(\rho || \hat{\rho}_j),$$

where  $J(W, b)$  is as defined previously, and  $\beta$  controls the weight of the sparsity penalty term. The term  $\hat{\rho}_j$  (implicitly) depends on  $W, b$  also, because it is the average activation of hidden unit  $j$ , and the activation of a hidden unit depends on the parameters  $W, b$ .

To incorporate the KL-divergence term into your derivative calculation, there is a simple-to-implement trick involving only a small change to your code. Specifically, where previously for the second layer ( $l = 2$ ), during backpropagation you would have computed

$$\delta_i^{(2)} = \left( \sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) f'(z_i^{(2)}),$$

now instead compute

$$\delta_i^{(2)} = \left( \left( \sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) + \beta \left( -\frac{\rho}{\hat{\rho}_i} + \frac{1-\rho}{1-\hat{\rho}_i} \right) \right) f'(z_i^{(2)}).$$

One subtlety is that you'll need to know  $\hat{\rho}_i$  to compute this term. Thus, you'll need to compute a forward pass on all the training examples first to compute the average activations on the training set, before computing backpropagation on any example. If your training set is small enough to fit comfortably in computer memory (this will be the case for the programming assignment), you can compute forward passes on all your examples and keep the resulting activations in memory and compute the  $\hat{\rho}_i$ s. Then you can use your precomputed activations to perform backpropagation on all your examples. If your data is too large to fit in memory, you may have to scan through your examples computing a forward pass on each to accumulate (sum up) the activations and compute  $\hat{\rho}_i$  (discarding the result of each forward pass after you have taken its activations  $a_i^{(2)}$  into account for computing  $\hat{\rho}_i$ ). Then after having computed  $\hat{\rho}_i$ , you'd have to redo the forward pass for each example so that you can do backpropagation on that example. In this latter case, you would end up computing a forward pass twice on each example in your training set, making it computationally less efficient.

The full derivation showing that the algorithm above results in gradient descent is beyond the scope of these notes. But if you implement the autoencoder using backpropagation modified this way, you will be performing gradient descent exactly on the objective  $J_{\text{sparse}}(W, b)$ . Using the derivative checking method, you will be able to verify this for yourself as well.

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

Language : 中文

Retrieved from "[http://ufldl.stanford.edu/wiki/index.php/Autoencoders\\_and\\_Sparsity](http://ufldl.stanford.edu/wiki/index.php/Autoencoders_and_Sparsity)"

- This page was last modified on 7 April 2013, at 12:43.