

# CS294A Lecture notes

Andrew Ng

## Sparse autoencoder

### 1 Introduction

Supervised learning is one of the most powerful tools of AI, and has led to automatic zip code recognition, speech recognition, self-driving cars, and a continually improving understanding of the human genome. Despite its significant successes, supervised learning today is still severely limited. Specifically, most applications of it still require that we manually specify the input features  $x$  given to the algorithm. Once a good feature representation is given, a supervised learning algorithm can do well. But in such domains as computer vision, audio processing, and natural language processing, there're now hundreds or perhaps thousands of researchers who've spent years of their lives slowly and laboriously hand-engineering vision, audio or text features. While much of this feature-engineering work is extremely clever, one has to wonder if we can do better. Certainly this labor-intensive hand-engineering approach does not scale well to new problems; further, ideally we'd like to have algorithms that can automatically learn even better feature representations than the hand-engineered ones.

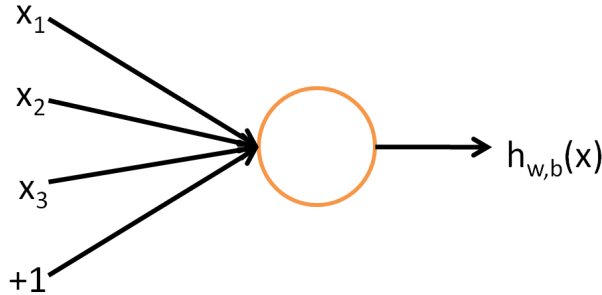
These notes describe the **sparse autoencoder** learning algorithm, which is one approach to automatically learn features from unlabeled data. In some domains, such as computer vision, this approach is not by itself competitive with the best hand-engineered features, but the features it can learn do turn out to be useful for a range of problems (including ones in audio, text, etc). Further, there're more sophisticated versions of the sparse autoencoder (not described in these notes, but that you'll hear more about later in the class) that do surprisingly well, and in many cases are competitive with or superior to even the best hand-engineered representations.

These notes are organized as follows. We will first describe feedforward neural networks and the backpropagation algorithm for supervised learning. Then, we show how this is used to construct an autoencoder, which is an unsupervised learning algorithm. Finally, we build on this to derive a sparse autoencoder. Because these notes are fairly notation-heavy, the last page also contains a summary of the symbols used.

## 2 Neural networks

Consider a supervised learning problem where we have access to labeled training examples  $(x^{(i)}, y^{(i)})$ . Neural networks give a way of defining a complex, non-linear form of hypotheses  $h_{W,b}(x)$ , with parameters  $W, b$  that we can fit to our data.

To describe neural networks, we will begin by describing the simplest possible neural network, one which comprises a single “neuron.” We will use the following diagram to denote a single neuron:



This “neuron” is a computational unit that takes as input  $x_1, x_2, x_3$  (and a  $+1$  intercept term), and outputs  $h_{W,b}(x) = f(W^T x) = f(\sum_{i=1}^3 W_i x_i + b)$ , where  $f : \mathbb{R} \mapsto \mathbb{R}$  is called the **activation function**. In these notes, we will choose  $f(\cdot)$  to be the sigmoid function:

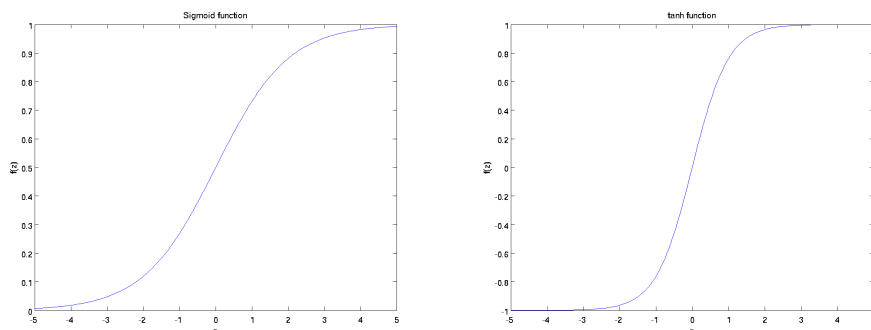
$$f(z) = \frac{1}{1 + \exp(-z)}.$$

Thus, our single neuron corresponds exactly to the input-output mapping defined by logistic regression.

Although these notes will use the sigmoid function, it is worth noting that another common choice for  $f$  is the hyperbolic tangent, or tanh, function:

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \quad (1)$$

Here are plots of the sigmoid and tanh functions:



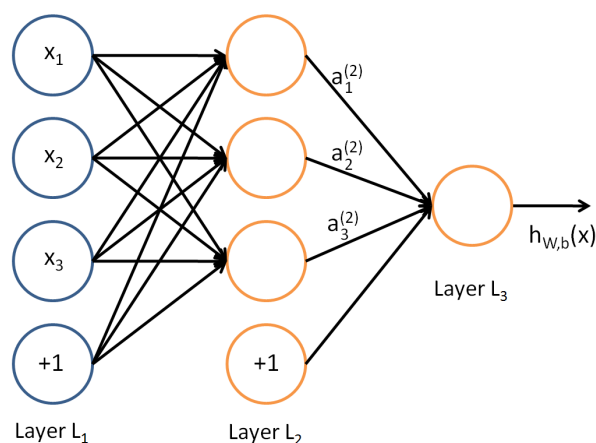
The  $\tanh(z)$  function is a rescaled version of the sigmoid, and its output range is  $[-1, 1]$  instead of  $[0, 1]$ .

Note that unlike CS221 and (parts of) CS229, we are not using the convention here of  $x_0 = 1$ . Instead, the intercept term is handled separately by the parameter  $b$ .

Finally, one identity that'll be useful later: If  $f(z) = 1/(1 + \exp(-z))$  is the sigmoid function, then its derivative is given by  $f'(z) = f(z)(1 - f(z))$ . (If  $f$  is the tanh function, then its derivative is given by  $f'(z) = 1 - (f(z))^2$ .) You can derive this yourself using the definition of the sigmoid (or tanh) function.

## 2.1 Neural network formulation

A neural network is put together by hooking together many of our simple “neurons,” so that the output of a neuron can be the input of another. For example, here is a small neural network:



In this figure, we have used circles to also denote the inputs to the network. The circles labeled “+1” are called **bias units**, and correspond to the intercept term. The leftmost layer of the network is called the **input layer**, and the rightmost layer the **output layer** (which, in this example, has only one node). The middle layer of nodes is called the **hidden layer**, because its values are not observed in the training set. We also say that our example neural network has 3 **input units** (not counting the bias unit), 3 **hidden units**, and 1 **output unit**.

We will let  $n_l$  denote the number of layers in our network; thus  $n_l = 3$  in our example. We label layer  $l$  as  $L_l$ , so layer  $L_1$  is the input layer, and layer  $L_{n_l}$  the output layer. Our neural network has parameters  $(W, b) = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$ , where we write  $W_{ij}^{(l)}$  to denote the parameter (or weight) associated with the connection between unit  $j$  in layer  $l$ , and unit  $i$  in layer  $l+1$ . (Note the order of the indices.) Also,  $b_i^{(l)}$  is the bias associated with unit  $i$  in layer  $l+1$ . Thus, in our example, we have  $W^{(1)} \in \mathbb{R}^{3 \times 3}$ , and  $W^{(2)} \in \mathbb{R}^{1 \times 3}$ . Note that bias units don’t have inputs or connections going into them, since they always output the value +1. We also let  $s_l$  denote the number of nodes in layer  $l$  (not counting the bias unit).

We will write  $a_i^{(l)}$  to denote the **activation** (meaning output value) of unit  $i$  in layer  $l$ . For  $l = 1$ , we also use  $a_i^{(1)} = x_i$  to denote the  $i$ -th input. Given a fixed setting of the parameters  $W, b$ , our neural network defines a hypothesis  $h_{W,b}(x)$  that outputs a real number. Specifically, the computation that this neural network represents is given by:

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}) \quad (2)$$

$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)}) \quad (3)$$

$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)}) \quad (4)$$

$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)}) \quad (5)$$

In the sequel, we also let  $z_i^{(l)}$  denote the total weighted sum of inputs to unit  $i$  in layer  $l$ , including the bias term (e.g.,  $z_i^{(2)} = \sum_{j=1}^n W_{ij}^{(1)}x_j + b_i^{(1)}$ ), so that  $a_i^{(l)} = f(z_i^{(l)})$ .

Note that this easily lends itself to a more compact notation. Specifically, if we extend the activation function  $f(\cdot)$  to apply to vectors in an element-wise fashion (i.e.,  $f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$ ), then we can write

Equations (2-5) more compactly as:

$$\begin{aligned} z^{(2)} &= W^{(1)}x + b^{(1)} \\ a^{(2)} &= f(z^{(2)}) \\ z^{(3)} &= W^{(2)}a^{(2)} + b^{(2)} \\ h_{W,b}(x) &= a^{(3)} = f(z^{(3)}) \end{aligned}$$

More generally, recalling that we also use  $a^{(1)} = x$  to also denote the values from the input layer, then given layer  $l$ 's activations  $a^{(l)}$ , we can compute layer  $l + 1$ 's activations  $a^{(l+1)}$  as:

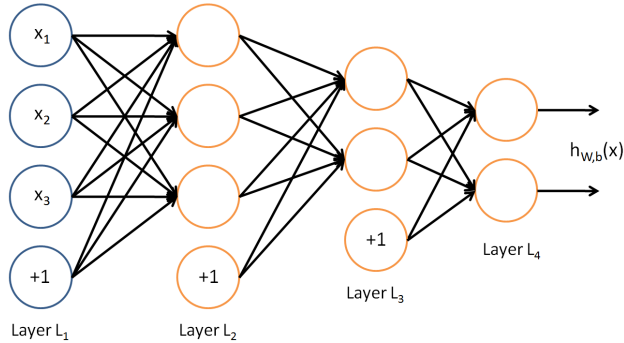
$$z^{(l+1)} = W^{(l)}a^{(l)} + b^{(l)} \quad (6)$$

$$a^{(l+1)} = f(z^{(l+1)}) \quad (7)$$

By organizing our parameters in matrices and using matrix-vector operations, we can take advantage of fast linear algebra routines to quickly perform calculations in our network.

We have so far focused on one example neural network, but one can also build neural networks with other **architectures** (meaning patterns of connectivity between neurons), including ones with multiple hidden layers. The most common choice is a  $n_l$ -layered network where layer 1 is the input layer, layer  $n_l$  is the output layer, and each layer  $l$  is densely connected to layer  $l + 1$ . In this setting, to compute the output of the network, we can successively compute all the activations in layer  $L_2$ , then layer  $L_3$ , and so on, up to layer  $L_{n_l}$ , using Equations (6-7). This is one example of a **feedforward** neural network, since the connectivity graph does not have any directed loops or cycles.

Neural networks can also have multiple output units. For example, here is a network with two hidden layers layers  $L_2$  and  $L_3$  and two output units in layer  $L_4$ :



To train this network, we would need training examples  $(x^{(i)}, y^{(i)})$  where  $y^{(i)} \in \mathbb{R}^2$ . This sort of network is useful if there're multiple outputs that you're interested in predicting. (For example, in a medical diagnosis application, the vector  $x$  might give the input features of a patient, and the different outputs  $y_i$ 's might indicate presence or absence of different diseases.)

## 2.2 Backpropagation algorithm

Suppose we have a fixed training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$  of  $m$  training examples. We can train our neural network using batch gradient descent. In detail, for a single training example  $(x, y)$ , we define the cost function with respect to that single example to be

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2.$$

This is a (one-half) squared-error cost function. Given a training set of  $m$  examples, we then define the overall cost function to be

$$\begin{aligned} J(W, b) &= \left[ \frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (W_{ji}^{(l)})^2 \\ &= \left[ \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{2} \|h_{W,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (W_{ji}^{(l)})^2 \end{aligned} \quad (8)$$

The first term in the definition of  $J(W, b)$  is an average sum-of-squares error term. The second term is a regularization term (also called a **weight decay** term) that tends to decrease the magnitude of the weights, and helps prevent overfitting.<sup>1</sup> The **weight decay parameter**  $\lambda$  controls the relative importance of the two terms. Note also the slightly overloaded notation:  $J(W, b; x, y)$  is the squared error cost with respect to a single example;  $J(W, b)$  is the overall cost function, which includes the weight decay term.

This cost function above is often used both for classification and for regression problems. For classification, we let  $y = 0$  or  $1$  represent the two class labels (recall that the sigmoid activation function outputs values in  $[0, 1]$ ; if

---

<sup>1</sup>Usually weight decay is not applied to the bias terms  $b_i^{(l)}$ , as reflected in our definition for  $J(W, b)$ . Applying weight decay to the bias units usually makes only a small difference to the final network, however. If you took CS229, you may also recognize weight decay this as essentially a variant of the Bayesian regularization method you saw there, where we placed a Gaussian prior on the parameters and did MAP (instead of maximum likelihood) estimation.

we were using a tanh activation function, we would instead use -1 and +1 to denote the labels). For regression problems, we first scale our outputs to ensure that they lie in the  $[0, 1]$  range (or if we were using a tanh activation function, then the  $[-1, 1]$  range).

Our goal is to minimize  $J(W, b)$  as a function of  $W$  and  $b$ . To train our neural network, we will initialize each parameter  $W_{ij}^{(l)}$  and each  $b_i^{(l)}$  to a small random value near zero (say according to a  $\mathcal{N}(0, \epsilon^2)$  distribution for some small  $\epsilon$ , say 0.01), and then apply an optimization algorithm such as batch gradient descent. Since  $J(W, b)$  is a non-convex function, gradient descent is susceptible to local optima; however, in practice gradient descent usually works fairly well. Finally, note that it is important to initialize the parameters randomly, rather than to all 0's. If all the parameters start off at identical values, then all the hidden layer units will end up learning the same function of the input (more formally,  $W_{ij}^{(1)}$  will be the same for all values of  $i$ , so that  $a_1^{(2)} = a_2^{(2)} = a_3^{(2)} = \dots$  for any input  $x$ ). The random initialization serves the purpose of **symmetry breaking**.

One iteration of gradient descent updates the parameters  $W, b$  as follows:

$$\begin{aligned} W_{ij}^{(l)} &:= W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) \\ b_i^{(l)} &:= b_i^{(l)} - \alpha \frac{\partial}{\partial b_i^{(l)}} J(W, b) \end{aligned}$$

where  $\alpha$  is the learning rate. The key step is computing the partial derivatives above. We will now describe the **backpropagation** algorithm, which gives an efficient way to compute these partial derivatives.

We will first describe how backpropagation can be used to compute  $\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x, y)$  and  $\frac{\partial}{\partial b_i^{(l)}} J(W, b; x, y)$ , the partial derivatives of the cost function  $J(W, b; x, y)$  defined with respect to a single example  $(x, y)$ . Once we can compute these, then by referring to Equation (8), we see that the derivative of the overall cost function  $J(W, b)$  can be computed as

$$\begin{aligned} \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) &= \left[ \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x^{(i)}, y^{(i)}) \right] + \lambda W_{ij}^{(l)}, \\ \frac{\partial}{\partial b_i^{(l)}} J(W, b) &= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b_i^{(l)}} J(W, b; x^{(i)}, y^{(i)}). \end{aligned}$$

The two lines above differ slightly because weight decay is applied to  $W$  but not  $b$ .

The intuition behind the backpropagation algorithm is as follows. Given a training example  $(x, y)$ , we will first run a “forward pass” to compute all the activations throughout the network, including the output value of the hypothesis  $h_{W,b}(x)$ . Then, for each node  $i$  in layer  $l$ , we would like to compute an “error term”  $\delta_i^{(l)}$  that measures how much that node was “responsible” for any errors in our output. For an output node, we can directly measure the difference between the network’s activation and the true target value, and use that to define  $\delta_i^{(n_l)}$  (where layer  $n_l$  is the output layer). How about hidden units? For those, we will compute  $\delta_i^{(l)}$  based on a weighted average of the error terms of the nodes that uses  $a_i^{(l)}$  as an input. In detail, here is the backpropagation algorithm:

1. Perform a feedforward pass, computing the activations for layers  $L_2$ ,  $L_3$ , and so on up to the output layer  $L_{n_l}$ .
2. For each output unit  $i$  in layer  $n_l$  (the output layer), set

$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{(n_l)}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)})$$

3. For  $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$

For each node  $i$  in layer  $l$ , set

$$\delta_i^{(l)} = \left( \sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

4. Compute the desired partial derivatives, which are given as:

$$\begin{aligned} \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x, y) &= a_j^{(l)} \delta_i^{(l+1)} \\ \frac{\partial}{\partial b_i^{(l)}} J(W, b; x, y) &= \delta_i^{(l+1)}. \end{aligned}$$

Finally, we can also re-write the algorithm using matrix-vectorial notation. We will use “ $\bullet$ ” to denote the element-wise product operator (denoted “ $\cdot$ ” in Matlab or Octave, and also called the Hadamard product), so that if  $a = b \bullet c$ , then  $a_i = b_i c_i$ . Similar to how we extended the definition of  $f(\cdot)$  to apply element-wise to vectors, we also do the same for  $f'(\cdot)$  (so that  $f'([z_1, z_2, z_3]) = [\frac{\partial}{\partial z_1} f(z_1), \frac{\partial}{\partial z_2} f(z_2), \frac{\partial}{\partial z_3} f(z_3)]$ ). The algorithm can then be written:



1. Perform a feedforward pass, computing the activations for layers  $L_2$ ,  $L_3$ , up to the output layer  $L_{n_l}$ , using Equations (6-7).
2. For the output layer (layer  $n_l$ ), set

$$\delta^{(n_l)} = -(y - a^{(n_l)}) \bullet f'(z^{(n_l)})$$

3. For  $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$

Set

$$\delta^{(l)} = ((W^{(l+1)})^T \delta^{(l+1)}) \bullet f'(z^{(l)})$$

4. Compute the desired partial derivatives:

$$\begin{aligned} \nabla_{W^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} (a^{(l)})^T, \\ \nabla_{b^{(l)}} J(W, b; x, y) &= \delta^{(l+1)}. \end{aligned}$$

**Implementation note:** In steps 2 and 3 above, we need to compute  $f'(z_i^{(l)})$  for each value of  $i$ . Assuming  $f(z)$  is the sigmoid activation function, we would already have  $a_i^{(l)}$  stored away from the forward pass through the network. Thus, using the expression that we worked out earlier for  $f'(z)$ , we can compute this as  $f'(z_i^{(l)}) = a_i^{(l)}(1 - a_i^{(l)})$ .

Finally, we are ready to describe the full gradient descent algorithm. In the pseudo-code below,  $\Delta W^{(l)}$  is a matrix (of the same dimension as  $W^{(l)}$ ), and  $\Delta b^{(l)}$  is a vector (of the same dimension as  $b^{(l)}$ ). Note that in this notation, “ $\Delta W^{(l)}$ ” is a matrix, and in particular it isn’t “ $\Delta$  times  $W^{(l)}$ .” We implement one iteration of batch gradient descent as follows:

1. Set  $\Delta W^{(l)} := 0$ ,  $\Delta b^{(l)} := 0$  (matrix/vector of zeros) for all  $l$ .
2. For  $i = 1$  to  $m$ ,
  - 2a. Use backpropagation to compute  $\nabla_{W^{(l)}} J(W, b; x, y)$  and  $\nabla_{b^{(l)}} J(W, b; x, y)$ .
  - 2b. Set  $\Delta W^{(l)} := \Delta W^{(l)} + \nabla_{W^{(l)}} J(W, b; x, y)$ .
  - 2c. Set  $\Delta b^{(l)} := \Delta b^{(l)} + \nabla_{b^{(l)}} J(W, b; x, y)$ .

3. Update the parameters:

$$\begin{aligned} W^{(l)} &:= W^{(l)} - \alpha \left[ \left( \frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \right] \\ b^{(l)} &:= b^{(l)} - \alpha \left[ \frac{1}{m} \Delta b^{(l)} \right] \end{aligned}$$

To train our neural network, we can now repeatedly take steps of gradient descent to reduce our cost function  $J(W, b)$ .

## 2.3 Gradient checking and advanced optimization

Backpropagation is a notoriously difficult algorithm to debug and get right, especially since many subtly buggy implementations of it—for example, one that has an off-by-one error in the indices and that thus only trains some of the layers of weights, or an implementation that omits the bias term—will manage to learn something that can look surprisingly reasonable (while performing less well than a correct implementation). Thus, even with a buggy implementation, it may not at all be apparent that anything is amiss. In this section, we describe a method for numerically checking the derivatives computed by your code to make sure that your implementation is correct. Carrying out the derivative checking procedure described here will significantly increase your confidence in the correctness of your code.

Suppose we want to minimize  $J(\theta)$  as a function of  $\theta$ . For this example, suppose  $J : \mathbb{R} \mapsto \mathbb{R}$ , so that  $\theta \in \mathbb{R}$ . In this 1-dimensional case, one iteration of gradient descent is given by

$$\theta := \theta - \alpha \frac{d}{d\theta} J(\theta).$$

Suppose also that we have implemented some function  $g(\theta)$  that purportedly computes  $\frac{d}{d\theta} J(\theta)$ , so that we implement gradient descent using the update  $\theta := \theta - \alpha g(\theta)$ . How can we check if our implementation of  $g$  is correct?

Recall the mathematical definition of the derivative as

$$\frac{d}{d\theta} J(\theta) = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}.$$

Thus, at any specific value of  $\theta$ , we can numerically approximate the derivative as follows:

$$\frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}$$

In practice, we set `EPSILON` to a small constant, say around  $10^{-4}$ . (There's a large range of values of `EPSILON` that should work well, but we don't set `EPSILON` to be "extremely" small, say  $10^{-20}$ , as that would lead to numerical roundoff errors.)

Thus, given a function  $g(\theta)$  that is supposedly computing  $\frac{d}{d\theta}J(\theta)$ , we can now numerically verify its correctness by checking that

$$g(\theta) \approx \frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}.$$

The degree to which these two values should approximate each other will depend on the details of  $J$ . But assuming `EPSILON` =  $10^{-4}$ , you'll usually find that the left- and right-hand sides of the above will agree to at least 4 significant digits (and often many more).

Now, consider the case where  $\theta \in \mathbb{R}^n$  is a vector rather than a single real number (so that we have  $n$  parameters that we want to learn), and  $J : \mathbb{R}^n \mapsto \mathbb{R}$ . In our neural network example we used " $J(W, b)$ ," but one can imagine "unrolling" the parameters  $W, b$  into a long vector  $\theta$ . We now generalize our derivative checking procedure to the case where  $\theta$  may be a vector.

Suppose we have a function  $g_i(\theta)$  that purportedly computes  $\frac{\partial}{\partial \theta_i} J(\theta)$ ; we'd like to check if  $g_i$  is outputting correct derivative values. Let  $\theta^{(i+)} = \theta + \text{EPSILON} \times \vec{e}_i$ , where

$$\vec{e}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

is the  $i$ -th basis vector (a vector of the same dimension as  $\theta$ , with a "1" in the  $i$ -th position and "0"s everywhere else). So,  $\theta^{(i+)}$  is the same as  $\theta$ , except its  $i$ -th element has been incremented by `EPSILON`. Similarly, let  $\theta^{(i-)} = \theta - \text{EPSILON} \times \vec{e}_i$  be the corresponding vector with the  $i$ -th element decreased by `EPSILON`. We can now numerically verify  $g_i(\theta)$ 's correctness by checking, for each  $i$ , that:

$$g_i(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2 \times \text{EPSILON}}.$$

When implementing backpropagation to train a neural network, in a cor-

rect implementation we will have that

$$\begin{aligned}\nabla_{W^{(l)}} J(W, b) &= \left( \frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \\ \nabla_{b^{(l)}} J(W, b) &= \frac{1}{m} \Delta b^{(l)}.\end{aligned}$$

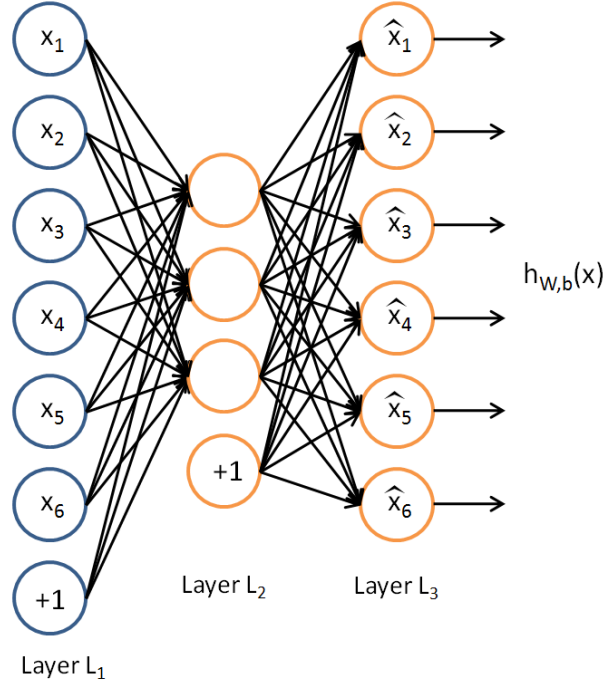
This result shows that the final block of psuedo-code in Section 2.2 is indeed implementing gradient descent. To make sure your implementation of gradient descent is correct, it is usually very helpful to use the method described above to numerically compute the derivatives of  $J(W, b)$ , and thereby verify that your computations of  $\left(\frac{1}{m} \Delta W^{(l)}\right) + \lambda W$  and  $\frac{1}{m} \Delta b^{(l)}$  are indeed giving the derivatives you want.

Finally, so far our discussion has centered on using gradient descent to minimize  $J(\theta)$ . If you have implemented a function that computes  $J(\theta)$  and  $\nabla_{\theta} J(\theta)$ , it turns out there are more sophisticated algorithms than gradient descent for trying to minimize  $J(\theta)$ . For example, one can envision an algorithm that uses gradient descent, but automatically tunes the learning rate  $\alpha$  so as to try to use a step-size that causes  $\theta$  to approach a local optimum as quickly as possible. There are other algorithms that are even more sophisticated than this; for example, there are algorithms that try to find an approximation to the Hessian matrix, so that it can take more rapid steps towards a local optimum (similar to Newton's method). A full discussion of these algorithms is beyond the scope of these notes, but one example is the **L-BFGS** algorithm. (Another example is **conjugate gradient**.) You will use one of these algorithms in the programming exercise. The main thing you need to provide to these advanced optimization algorithms is that for any  $\theta$ , you have to be able to compute  $J(\theta)$  and  $\nabla_{\theta} J(\theta)$ . These optimization algorithms will then do their own internal tuning of the learning rate/step-size  $\alpha$  (and compute its own approximation to the Hessian, etc.) to automatically search for a value of  $\theta$  that minimizes  $J(\theta)$ . Algorithms such as L-BFGS and conjugate gradient can often be much faster than gradient descent.

### 3 Autoencoders and sparsity

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 unlabeled training examples set  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots\}$ , 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.<sup>2</sup>

<sup>2</sup>In fact, this simple autoencoder often ends up learning a low-dimensional representation very similar to PCA's.

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.<sup>3</sup>

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  $x$  that led to that activation. Thus, we will write  $a_j^{(2)}(x)$  to denote the activation of this hidden unit when the network is given a specific input  $x$ . Further, let

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

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

---

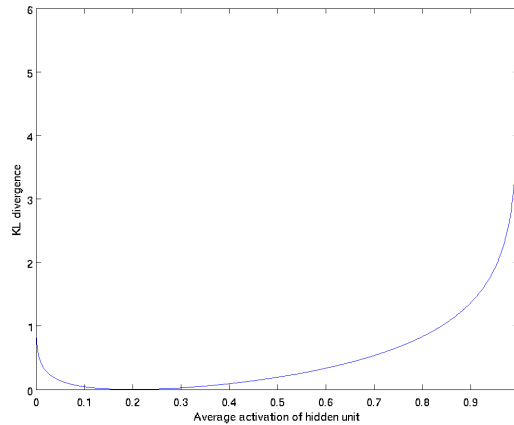
<sup>3</sup>This discussion assumes a sigmoid activation function. If you are using a tanh activation function, then we think of a neuron as being inactive when it outputs values close to -1.

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.

## 4 Visualization

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  $a_i^{(1)}$  as some non-linear feature of the input  $x$ . We ask: What input image  $x$  would cause  $a_i^{(1)}$  to be maximally activated? 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,  $j = 1, \dots, 100$ ) to

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

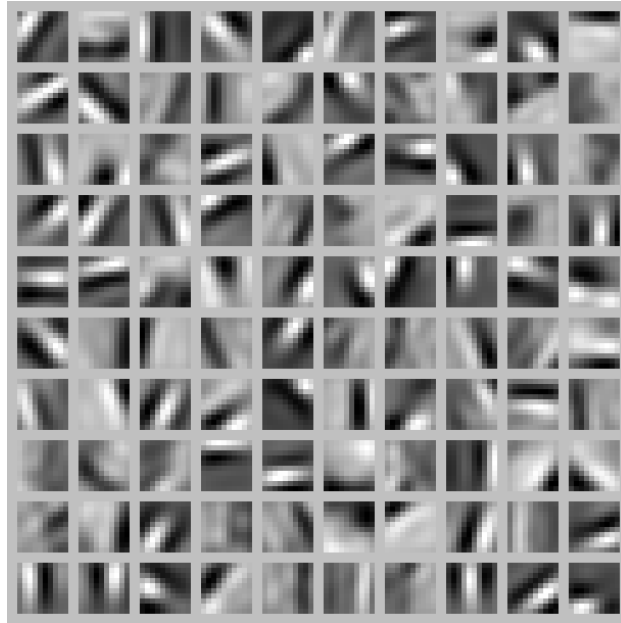
By displaying the image formed by these pixel intensity values, we can begin to understand what feature hidden unit  $i$  is looking for.

If we have an autoencoder with 100 hidden units (say), then 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<sup>4</sup>) we get the following result:

---

<sup>4</sup>The results below 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.



Each square in the figure above shows the (norm bounded) input image  $x$  that maximally activates 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.

## 5 Summary of notation

$x$	Input features for a training example, $x \in \mathbb{R}^n$ .
$y$	Output/target values. Here, $y$ can be vector valued. In the case of an autoencoder, $y = x$ .
$(x^{(i)}, y^{(i)})$	The $i$ -th training example
$h_{W,b}(x)$	Output of our hypothesis on input $x$ , using parameters $W, b$ . This should be a vector of the same dimension as the target value $y$ .
$W_{ij}^{(l)}$	The parameter associated with the connection between unit $j$ in layer $l$ , and unit $i$ in layer $l + 1$ .
$b_i^{(l)}$	The bias term associated with unit $i$ in layer $l + 1$ . Can also be thought of as the parameter associated with the connection between the bias unit in layer $l$ and unit $i$ in layer $l + 1$ .
$\theta$	Our parameter vector. It is useful to think of this as the result of taking the parameters $W, b$ and “unrolling” them into a long column vector.
$a_i^{(l)}$	Activation (output) of unit $i$ in layer $l$ of the network. In addition, since layer $L_1$ is the input layer, we also have $a_i^{(1)} = x_i$ .
$f(\cdot)$	The activation function. Throughout these notes, we used $f(z) = \tanh(z)$ .
$z_i^{(l)}$	Total weighted sum of inputs to unit $i$ in layer $l$ . Thus, $a_i^{(l)} = f(z_i^{(l)})$ .
$\alpha$	Learning rate parameter
$s_l$	Number of units in layer $l$ (not counting the bias unit).
$n_l$	Number layers in the network. Layer $L_1$ is usually the input layer, and layer $L_{n_l}$ the output layer.
$\lambda$	Weight decay parameter.
$\hat{x}$	For an autoencoder, its output; i.e., its reconstruction of the input $x$ . Same meaning as $h_{W,b}(x)$ .
$\rho$	Sparsity parameter, which specifies our desired level of sparsity
$\hat{\rho}_i$	The average activation of hidden unit $i$ (in the sparse autoencoder).
$\beta$	Weight of the sparsity penalty term (in the sparse autoencoder objective).