

Chapter 6

Deep Feedforward Networks

Deep feedforward networks, also called **feedforward multilayer perceptrons** (MLPs), are the quintessential feedforward neural network. The goal of a feedforward network is to approximate some function. For a classifier, $y = f^*(\mathbf{x})$ maps an input \mathbf{x} to a category y . The model defines a mapping $\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta})$ and learns the value of the parameters $\boldsymbol{\theta}$ in the best function approximation.

These models are called **feedforward** because information flows from the input function being evaluated from \mathbf{x} , through the intermediate layers to define f , and finally to the output \mathbf{y} . There are no **feedback** connections; outputs of the model are not fed back into itself. When feedforward networks are extended to include feedback connections, they are called **recurrent networks**, as presented in chapter 10.

Feedforward networks are of extreme importance to machine learning.

overall length of the chain gives the **depth** of the model. This arose from this terminology. The final layer of a feedforward network is called the **output layer**. During neural network training, we drive the network to produce the desired output. The training data provides us with noisy, approximate examples at different training points. Each example \mathbf{x} is accompanied by a target value y . The training examples specify directly what the output layer should produce for \mathbf{x} ; it must produce a value that is close to y . The behavior of the hidden layers is not directly specified by the training data. The learning algorithm learns how to use those layers to produce the desired output, but does not say what each individual layer should do. Instead, the learning algorithm decides how to use these layers to best implement an approximation of the target function. The training data does not show the desired output for each input, so the layers in between are called **hidden layers**.

Finally, these networks are called *neural* because they are inspired by neuroscience. Each hidden layer of the network is typically a vector of nodes. The dimensionality of these hidden layers determines the **width** of the network. Each element of the vector may be interpreted as playing a role in the computation. Rather than thinking of the layer as representing a single vector, we can also think of the layer as consisting of many **units**, each representing a vector-to-scalar function. Each unit receives input from the previous layer in the sense that it receives input from many other units and produces an activation value. The idea of using many layers of vector-to-scalar functions is drawn from neuroscience. The choice of the functions f used in these representations is also loosely guided by neuroscience, specifically the functions that biological neurons compute. Modern neural network design, however, is guided by many mathematical and engineering considerations. The goal of neural networks is not to perfectly model the brain.

nonlinear transformation. Equivalently, we can apply the kernel trick from section 5.7.2, to obtain a nonlinear learning algorithm based on the ϕ mapping. We can think of ϕ as providing a set of features for \mathbf{x} , as providing a new representation for \mathbf{x} .

The question is then how to choose the mapping ϕ .

1. One option is to use a very generic ϕ , such as the infinite-dimensional mapping that is implicitly used by kernel machines based on the Gaussian kernel. If the space of high enough dimension, we can always have enough features to fit the training set, but generalization to the test set often fails. Generic feature mappings are usually based only on smoothness and do not encode enough prior information about the problem.
2. Another option is to manually engineer ϕ . Until the advent of deep learning, this was the dominant approach. It requires decades of domain knowledge for each separate task, with practitioners specializing in different domains such as speech recognition or computer vision, and with different feature sets for different domains.
3. The strategy of deep learning is to learn ϕ . In this approach, we model the output $y = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \phi(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{w}$. We now have parameters $\boldsymbol{\theta}$ for ϕ from a broad class of functions, and parameters \mathbf{w} for the linear combination to the desired output. This is an example of a deep feedforward network with ϕ defining a hidden layer. This approach is the only one that does not give up on the convexity of the training problem, but it has its own harms. In this approach, we parametrize the representation ϕ and use the optimization algorithm to find the $\boldsymbol{\theta}$ that

deterministic mappings from \mathbf{x} to \mathbf{y} that lack feedback connections. In the chapters presented later, apply these principles to learning stochastic mappings with feedback, and probability distributions over a single variable.

We begin this chapter with a simple example of a feedforward network. As we address each of the design decisions needed to deploy a neural network, we review the design decisions necessary for a linear model: choosing the activation function, and the form of the output units. We review these basic concepts of learning, then proceed to confront some of the design decisions specific to feedforward networks. Feedforward networks have introduced the concept of a hidden layer, and this requires us to choose the **activation function** to be used to compute the hidden layer values. We must also address the architecture of the network, including how many layers the network should have, how many layers should be connected to each other, and how many units should be in each layer. Learning in deep neural networks requires computing gradients of complicated functions. We present the **back-propagation** algorithm and its modern generalizations, which can be used to efficiently compute gradients. Finally, we close with some historical perspective.

6.1 Example: Learning XOR

To make the idea of a feedforward network more concrete, we will use as an example of a fully functioning feedforward network on a very simple task: learning the XOR function.

The XOR function (“exclusive or”) is an operation on two binary values, x_1 and x_2 . When exactly one of these binary values is equal to 1, the XOR function returns 1; otherwise, it returns 0.

appropriate cost function for modeling binary data. More are described in section 6.2.2.2.

Evaluated on our whole training set, the MSE loss func

$$J(\boldsymbol{\theta}) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \boldsymbol{\theta}))^2 .$$

Now we must choose the form of our model, $f(\mathbf{x}; \boldsymbol{\theta})$. So a linear model, with $\boldsymbol{\theta}$ consisting of \mathbf{w} and b . Our model is

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{x}^\top \mathbf{w} + b.$$

We can minimize $J(\boldsymbol{\theta})$ in closed form with respect to \mathbf{w} and b using the normal equations.

After solving the normal equations, we obtain $\mathbf{w} = \mathbf{0}$ and $b = 0.5$. This model simply outputs 0.5 everywhere. Why does this happen? The problem is how a linear model is not able to represent the XOR function. The solution to this problem is to use a model that learns a different feature space. A non-linear model is able to represent the solution.

Specifically, we will introduce a simple feedforward network with two layers. The first layer contains two hidden units. See figure 6.2 for an illustration. This feedforward network has a vector of hidden units \mathbf{h} that is computed by the function $f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c})$. The values of these hidden units are then used as input for a second layer. The second layer is the output layer of the network. This layer is still just a linear regression model, but now it is applied to the hidden units. The network now contains two functions chained together, $y = f^{(2)}(\mathbf{h}; \mathbf{w}, b)$, with the complete model being $f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b)$.

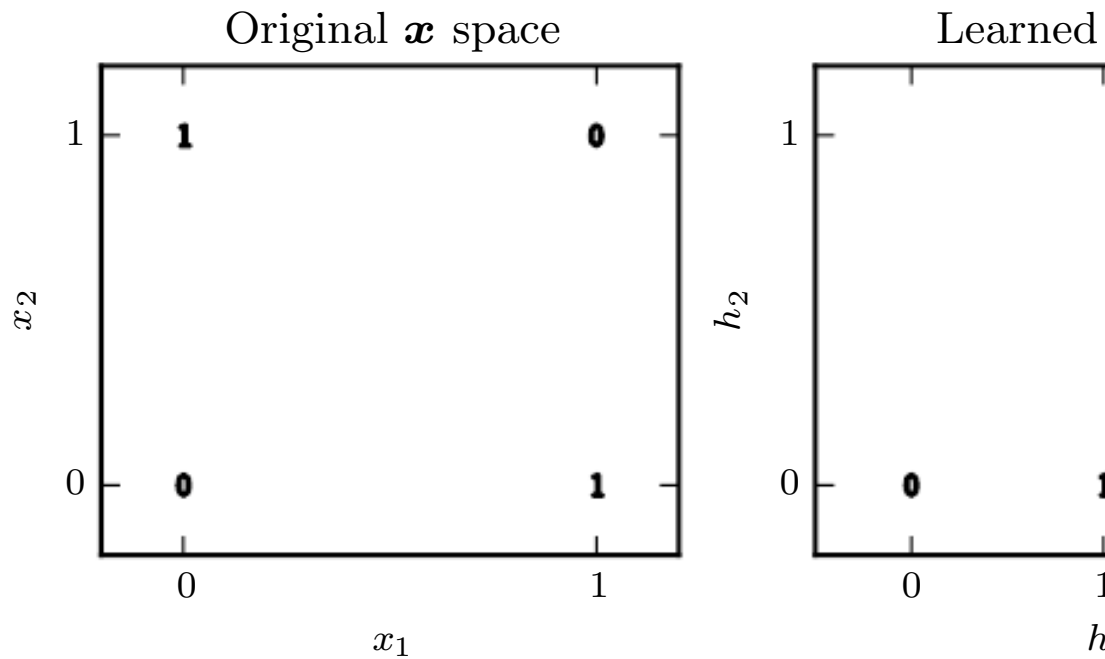


Figure 6.1: Solving the XOR problem by learning a representation. The values printed on the plot indicate the value that the learned function maps to. *(Left)* A linear model applied directly to the original input cannot solve this problem. When $x_1 = 0$, the model’s output must increase as x_2 increases, but when $x_1 = 1$, the model’s output must decrease as x_2 increases. A linear model cannot use the coefficient w_2 to x_2 . The linear model therefore cannot solve this problem. *(Right)* If the data is represented by the features extracted by a neural network, a linear model can solve the problem. In our example solution, the two points that must be separated have collapsed into a single point in feature space. In other words, the

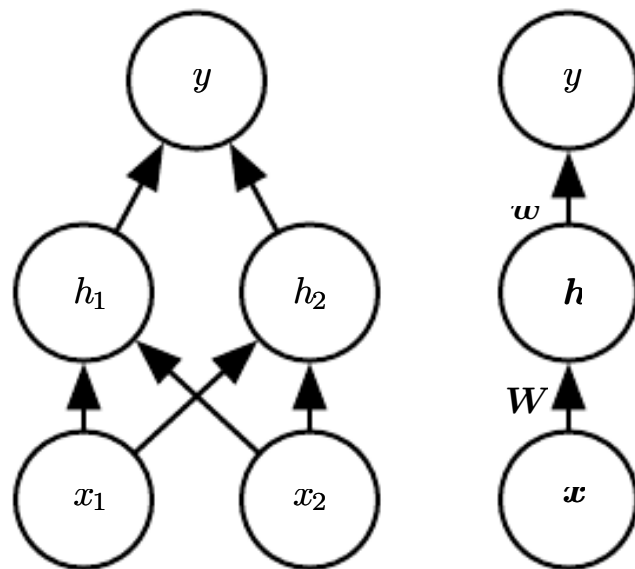
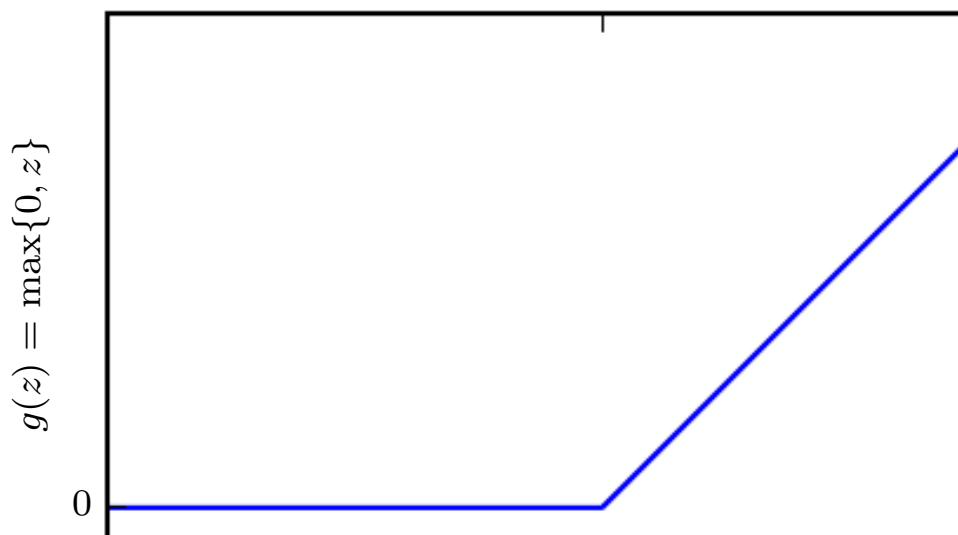


Figure 6.2: An example of a feedforward network, drawn in two different styles. (Left) This is the feedforward network we use to solve the XOR example in Section 6.1. It has an input layer containing two units. (Right) In this style, we draw every unit in a layer as a single node. This style is explicit and unambiguous, but for networks larger than this, it can consume too much space. (Right) In this style, we draw a node in a layer as a single node representing a layer's activations. This style is much more compact, but we annotate the edges in this graph with the name of the parameters that describe the relationship between two layers. Here, we indicate that a matrix W describes the mapping from x to h , and a vector w describes the mapping from h to y . We also indicate intercept parameters associated with each layer when labeling the edges.



affine transformation from an input vector to an output scalar. An affine transformation from a vector \mathbf{x} to a vector \mathbf{h} , so a vector of parameters is needed. The activation function g is typically a sigmoid function that is applied element-wise, with $h_i = g(\mathbf{x}^\top \mathbf{W}_{:,i} + c_i)$. In many cases, the default recommendation is to use the **rectified linear unit** (Glorot *et al.*, 2009; Nair and Hinton, 2010; Glorot *et al.*, 2011a), denoted ReLU , with the function $g(z) = \max\{0, z\}$, depicted in figure 6.3.

We can now specify our complete network as

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b$$

We can then specify a solution to the XOR problem. Let

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

and $b = 0$.

We can now walk through how the model processes a batch of data. Let \mathbf{X} be the design matrix containing all four points in the binary XOR example per row:

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}.$$

In this space, all the examples lie along a line with slope 1. On this line, the output needs to begin at 0, then rise to 1, then drop to 0. A linear model cannot implement such a function. To finish off, for each example, we apply the rectified linear transform

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}.$$

This transformation has changed the relationship between the inputs. They no longer lie on a single line. As shown in figure 6.1, they now form a shape that a linear model can solve the problem.

We finish with multiplying by the weight vector \mathbf{w} :

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}.$$

The neural network has obtained the correct answer for every example.

In this example, we simply specified the solution, then saw that it produced zero error. In a real situation, there might be billions of training examples, so one cannot simply guess the solution. Instead, a gradient-based optimization algorithm can find a solution that produces very little error. The solution we described to the XOR problem is a global minimum of the loss function, so gradient descent will find it. There are other equivalent solutions to the XOR problem.

The largest difference between the linear models we have seen and neural networks is that the nonlinearity of a neural network causes the loss functions to become nonconvex. This means that neural networks trained by using iterative, gradient-based optimizers that try to minimize the function to a very low value, rather than the linear equations solved by linear regression models or the convex optimization algorithms used to train logistic regression or SVMs, do not have a convergence guarantee starting from any initial parameters (in theory—but can encounter numerical problems). Stochastic gradient descent on nonconvex loss functions has no such convergence guarantee, and the values of the initial parameters matter. For feedforward neural networks, we initialize all weights to small random values. The biases may be initialized to small positive values. The iterative gradient-based optimizers used to train feedforward networks and almost all other deep learning models are described in detail in chapter 8, with parameter initialization in particular in section 8.4. For the moment, it suffices to understand that training is almost always based on using the gradient to descend the loss function in one way or another. The specific algorithms are improvements on the ideas of gradient descent, introduced in section 4.3, and the most often improvements of the stochastic gradient descent algorithm described in section 5.9.

We can of course train models such as linear regression and support vector machines with gradient descent too, and in fact this is common. The parameter set is extremely large. From this point of view, training a neural network is much different from training any other model. Computing the gradient is much more complicated for a neural network but can still be done. In Section 6.5 we describe how to obtain the gradient using backpropagation.

the same as those for other parametric models, such as linear models.

In most cases, our parametric model defines a distribution over \mathbf{y} , and we simply use the principle of maximum likelihood. The cost function is the cross-entropy between the training data and the model's predicted distribution.

Sometimes, we take a simpler approach, where rather than modeling a probability distribution over \mathbf{y} , we merely predict some scalar value on \mathbf{x} . Specialized loss functions enable us to train a predictor in this case.

The total cost function used to train a neural network is typically a combination of the primary cost functions described here with a regularization term. We have already seen some simple examples of regularization applied to linear models in section 5.2.2. The weight decay approach used for linear models is also applicable to deep neural networks and is among the most common regularization strategies. More advanced regularization strategies for deep networks are described in chapter 7.

6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

Most modern neural networks are trained using maximum likelihood estimation, meaning that the cost function is simply the negative log-likelihood, or equivalently, the cross-entropy between the training data and the model's predicted distribution. The cost function is given by

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x}).$$

The specific form of the cost function changes from model to model, but the principle remains the same.

mean squared error holds for a linear model, but in fact, regardless of the $f(\mathbf{x}; \boldsymbol{\theta})$ used to predict the mean of the \mathbf{C}

An advantage of this approach of deriving the cost function from the likelihood is that it removes the burden of designing cost functions. Specifying a model $p(\mathbf{y} \mid \mathbf{x})$ automatically determines a cost function.

One recurring theme throughout neural network design is that the cost function must be large and predictable enough to be used for the learning algorithm. Functions that saturate (become very small) are bad for this objective because they make the gradient become very small. This happens because the activation functions used to produce the outputs of hidden units or the output units saturate. The negative log-likelihood cost function avoids this problem for many models. Several output units in a neural network can saturate when its argument is very negative. The negative log-likelihood cost function undoes the \exp of some of these saturating units. We discuss the interaction between the cost function and the model in section 6.2.2.

One unusual property of the cross-entropy cost used in maximum likelihood estimation is that it usually does not have a minimum. This is true for the models commonly used in practice. For discrete output models are parametrized in such a way that they cannot reach a value of zero or one, but can come arbitrarily close to doing so. The softmax is an example of such a model. For real-valued output models we can control the density of the output distribution (for example, the variance parameter of a Gaussian output distribution) to avoid assigning extremely high density to the correct training set. The cross-entropy approaching negative infinity. Regularization techniques in chapter 7 provide several different ways of modifying the

rather than by having a specific parametric form. From this perspective, we can view the cost function as being a **functional** rather than a function. A functional is a mapping from functions to real numbers. This perspective on learning as choosing a function rather than merely choosing parameters. We can design our cost functional to have its minimum at the function we desire. For example, we can design the cost functional so that its minimum lie on the function that maps \mathbf{x} to the expected output \mathbf{y} . Solving an optimization problem with respect to a function rather than its parameters is a tool called **calculus of variations**, described in section 19.1. To understand calculus of variations to understand the context of this chapter. At the moment, it is only necessary to understand that calculus of variations is used to derive the following two results.

Our first result derived using calculus of variations is the minimization problem

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

yields

$$f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y}|\mathbf{x})}[\mathbf{y}],$$

so long as this function lies within the class we optimize over. If we could train on infinitely many samples from the true data generating process, minimizing the mean squared error cost function would give us the mean of \mathbf{y} for each value of \mathbf{x} .

Different cost functions give different statistics. A second result from calculus of variations is that

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|_1$$

yields a function that predicts the *median* value of \mathbf{y} for each value of \mathbf{x} . The function may be described by the family of functions we optimize over. This function is commonly called **mean absolute error**.

Unfortunately, mean squared error and mean absolute error both saturate when used with gradient-based optimization. Saturation produces very small gradients when combined with the sigmoid function. This is one reason that the cross-entropy cost function is more commonly used than squared error or mean absolute error, even when it is not the true underlying distribution $p(\mathbf{y} | \mathbf{x})$.

6.2.2 Output Units

The choice of cost function is tightly coupled with the choice of activation function.

6.2.2.1 Linear Units for Gaussian Output Distribution

One simple kind of output unit is based on an affine transformation and a Gaussian nonlinearity. These are often just called linear units.

Given features \mathbf{h} , a layer of linear output units produces

Linear output layers are often used to produce the parameters of a Gaussian distribution:

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \mathbf{I}).$$

Maximizing the log-likelihood is then equivalent to minimizing the squared error.

The maximum likelihood framework makes it straightforward to model the covariance of the Gaussian too, or to make the covariance a function of the input. However, the covariance must be constant and positive definite matrix for all inputs. It is difficult to satisfy such constraints with a single output layer, so typically other output units are used to parameterize the covariance. Approaches to modeling the covariance are described shortly.

Because linear units do not saturate, they pose little difficulty for gradient-based optimization algorithms and may be used with a wide variety of other algorithms.

6.2.2.2 Sigmoid Units for Bernoulli Output Distribution

Many tasks require predicting the value of a binary variable. Binary classification problems with two classes can be cast in this form.

The maximum likelihood approach is to define a Bernoulli distribution

outside the unit interval, the gradient of the output of the its parameters would be $\mathbf{0}$. A gradient of $\mathbf{0}$ is typically pr learning algorithm no longer has a guide for how to impr parameters.

Instead, it is better to use a different approach that en strong gradient whenever the model has the wrong answer. on using sigmoid output units combined with maximum lik

A sigmoid output unit is defined by

$$\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b),$$

where σ is the logistic sigmoid function described in section

We can think of the sigmoid output unit as having two uses a linear layer to compute $z = \mathbf{w}^\top \mathbf{h} + b$. Next, it uses function to convert z into a probability.

We omit the dependence on \mathbf{x} for the moment to di probability distribution over y using the value z . The sign by constructing an unnormalized probability distribution sum to 1. We can then divide by an appropriate const probability distribution. If we begin with the assumption tha probabilities are linear in y and z , we can exponentiate to ob probabilities. We then normalize to see that this yields a controlled by a sigmoidal transformation of z :

$$\begin{aligned} \log \tilde{P}(y) &= yz, \\ \tilde{P}(y) &= \exp(yz), \\ &= \frac{\exp(yz)}{\sum_y \exp(yz)} \end{aligned}$$

likelihood learning of a Bernoulli parametrized by a sigmoid

$$\begin{aligned} J(\boldsymbol{\theta}) &= -\log P(y \mid \boldsymbol{x}) \\ &= -\log \sigma((2y - 1)z) \\ &= \zeta((1 - 2y)z). \end{aligned}$$

This derivation makes use of some properties from section 2.2. In terms of the softplus function, we can see that $(1 - 2y)z$ is very negative. Saturation thus occurs only when y has the right answer—when $y = 1$ and z is very positive, or $y = 0$ and z is very negative. When z has the wrong sign, the argument to the softplus function, $(1 - 2y)z$, may be simplified to $|z|$. As $|z|$ becomes large while the softplus function asymptotes toward simply returning $|z|$, the derivative with respect to z asymptotes to $\text{sign}(z)$, so, in the case of an incorrect z , the softplus function does not shrink the gradient. This is useful because it means that gradient-based learning can still be useful for a mistaken z .

When we use other loss functions, such as mean squared error, the gradient can saturate anytime $\sigma(z)$ saturates. The sigmoid activation function saturates to 0 when z becomes very negative and saturates to 1 when z becomes very positive. The gradient can shrink too small to be useful for learning. This is a problem whether the model has the correct answer or the incorrect answer. The maximum likelihood is almost always the preferred approach for learning in output units.

Analytically, the logarithm of the sigmoid is always defined. Since the sigmoid returns values restricted to the open interval $(0, 1)$, the entire closed interval of valid probabilities $[0, 1]$. In softplus,

can be used inside the model itself, if we wish the model to produce n different options for some internal variable.

In the case of binary variables, we wished to produce a

$$\hat{y} = P(y = 1 \mid \mathbf{x}).$$

Because this number needed to lie between 0 and 1, and the logarithm of the number to be well behaved for gradient-descent on the log-likelihood, we chose to instead predict a number between 0 and 1. Exponentiating and normalizing gave us a Bernoulli distribution, which is the sigmoid function.

To generalize to the case of a discrete variable with n options, we wish to produce a vector $\hat{\mathbf{y}}$, with $\hat{y}_i = P(y = i \mid \mathbf{x})$. We require each element of $\hat{\mathbf{y}}$ to be between 0 and 1, but also that the entire vector represents a valid probability distribution. The same approach that worked for the Bernoulli distribution generalizes to the multinoulli distribution. A softmax layer predicts unnormalized log probabilities:

$$\mathbf{z} = \mathbf{W}^\top \mathbf{h} + \mathbf{b},$$

where $z_i = \log \tilde{P}(y = i \mid \mathbf{x})$. The softmax function can then be used to normalize \mathbf{z} to obtain the desired $\hat{\mathbf{y}}$. Formally, the softmax function is defined as

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}.$$

As with the logistic sigmoid, the use of the \exp function in the softmax makes training the softmax to output a target value y using maximum likelihood difficult. In this case, we wish to maximize $\log P(y = i \mid \mathbf{z}) = \log \text{softmax}(\mathbf{z})_i$.

that this term can be roughly approximated by $\max_j z_j$. based on the idea that $\exp(z_k)$ is insignificant for any z_k that $\max_j z_j$. The intuition we can gain from this approximation is that the log-likelihood cost function always strongly penalizes the incorrect prediction. If the correct answer already has the largest input, the $-z_i$ term and the $\log \sum_j \exp(z_j) \approx \max_j z_j = z_i$ term cancel out. This example will then contribute little to the overall training loss, being dominated by other examples that are not yet correctly classified.

So far we have discussed only a single example. Overall, maximizing the log-likelihood will drive the model to learn parameters that drive the predicted probabilities to match the fraction of counts of each outcome observed in the training data.

$$\text{softmax}(\mathbf{z}(\mathbf{x}; \boldsymbol{\theta}))_i \approx \frac{\sum_{j=1}^m \mathbf{1}_{y^{(j)}=i, \mathbf{x}^{(j)}=\mathbf{x}}}{\sum_{j=1}^m \mathbf{1}_{\mathbf{x}^{(j)}=\mathbf{x}}}$$

Because maximum likelihood is a consistent estimator, this is true as long as the model family is capable of representing the true distribution. In practice, limited model capacity and imperfect optimization mean the model is only able to approximate these fractions.

Many objective functions other than the log-likelihood can be used with the softmax function. Specifically, objective functions that do not undo the exp of the softmax fail to learn when the arguments are very negative, causing the gradient to vanish. In particular, the cross-entropy loss function for softmax units can fail to train the model even when the model makes highly confident incorrect predictions. To understand why these other loss functions can fail, we will look at the softmax function itself.

Like the sigmoid, the softmax activation can saturate.

The reformulated version enables us to evaluate softmax with errors, even when \mathbf{z} contains extremely large or extremely small values. Examining the numerically stable variant, we see that the softmax is driven by the amount that its arguments deviate from maximum.

An output $\text{softmax}(\mathbf{z})_i$ saturates to 1 when the corresponding $z_i = \max_j z_j$ and z_i is much greater than all the other z_j . $\text{softmax}(\mathbf{z})_i$ can also saturate to 0 when z_i is not maximal and much smaller. This is a generalization of the way that sigmoid can cause similar difficulties for learning if the loss function does not compensate for it.

The argument \mathbf{z} to the softmax function can be produced in many ways. The most common is simply to have an earlier layer of the network produce every element of \mathbf{z} , as described above using the linear layer. While straightforward, this approach actually overparametrizes the problem. The constraint that the n outputs must sum to 1 means that only $n-1$ are necessary; the probability of the n -th value may be obtained from the first $n-1$ probabilities from 1. We can thus impose a requirement that the last element of \mathbf{z} be fixed. For example, we can require that $z_n = 0$, which is what the sigmoid unit does. Defining $P(y = 1 \mid \mathbf{x}) = \sigma(z)$ is equivalent to $P(y = 1 \mid \mathbf{x}) = \text{softmax}(\mathbf{z})_1$ with a two-dimensional \mathbf{z} and $z_n = 0$. The n argument approaches the softmax argument but has a restricted set of probability distributions but have different learning dynamics. In practice, there is rarely much difference between using the overparametrized version and the restricted version, and it is simpler to implement the overparametrized version.

From a neuroscientific point of view, it is interesting to think of softmax as a way to create a form of competition between the units that produce the softmax outputs always sum to 1 so an increase in the value of one unit must be compensated by a decrease in the value of another.

It would perhaps be better to call the softmax function current name is an entrenched convention.

6.2.2.4 Other Output Types

The linear, sigmoid, and softmax output units described are common. Neural networks can generalize to almost any kind we wish. The principle of maximum likelihood provides a good cost function for nearly any kind of output layer.

In general, if we define a conditional distribution $p(\mathbf{y} \mid \mathbf{x})$, maximum likelihood suggests we use $-\log p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$ as our cost function.

In general, we can think of the neural network as representing a function $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\omega}$. The outputs of this function are not direct predictions of \mathbf{y} . $\boldsymbol{\omega}$ provides the parameters for a distribution over \mathbf{y} . The cost function can then be interpreted as $-\log p(\mathbf{y}; \boldsymbol{\omega}(\mathbf{x}))$.

For example, we may wish to learn the variance of a continuous variable y given \mathbf{x} . In the simple case, where the variance σ^2 is a constant, we can use a simple expression because the maximum likelihood estimator of σ^2 is the empirical mean of the squared difference between observation and predicted value. A computationally more expensive approach that does not require special-case code is to simply include the variance as one of the parameters of a distribution $p(\mathbf{y} \mid \mathbf{x})$ that is controlled by $\boldsymbol{\omega} = f(\mathbf{x}; \boldsymbol{\theta})$. The cost function $-\log p(\mathbf{y}; \boldsymbol{\omega}(\mathbf{x}))$ will then provide a cost function with the flexibility necessary to make our optimization procedure incrementally more general than the simple case where the standard deviation does not depend on \mathbf{x} . We can make a new parameter in the network that is copied directly to the variance parameter. The parameter might be σ itself or could be a parameter v representing the variance.

and logarithm operations is well behaved. By comparison, output in terms of variance, we would need to use division, which becomes arbitrarily steep near zero. While large gradients or arbitrarily large gradients usually result in instability. If we parameterize in terms of standard deviation, the log-likelihood would still increase as squaring. The gradient through the squaring operation is small, making it difficult to learn parameters that are squared. Regardless, if we use standard deviation, variance, or precision, we must ensure the covariance matrix of the Gaussian is positive definite. Because the eigenvalues of the precision matrix are the reciprocals of the eigenvalues of the covariance matrix, ensuring that the precision matrix is positive definite is equivalent to ensuring that the covariance matrix is positive definite. If the diagonal matrix, or a scalar times the diagonal matrix, then we need to enforce on the output of the model is positivity. If the output is the raw activation of the model used to determine the variance, we can use the softplus function to obtain a positive precision. The same strategy applies equally if using variance or standard deviation, precision or if using a scalar times identity rather than diagonal.

It is rare to learn a covariance or precision matrix with a diagonal. If the covariance is full and conditional, then a constraint must be chosen that guarantees positive definiteness of the predicted covariance. This can be achieved by writing $\Sigma(\mathbf{x}) = \mathbf{B}(\mathbf{x})\mathbf{B}^\top(\mathbf{x})$, where $\mathbf{B}(\mathbf{x})$ is a square matrix. One practical issue if the matrix is full rank is that the log-likelihood is expensive, with a $d \times d$ matrix requiring $O(d^3)$ for the determinant and inverse of $\Sigma(\mathbf{x})$ (or equivalently, and more efficiently, eigendecomposition or that of $\mathbf{B}(\mathbf{x})$).

We often want to perform multimodal regression, that is, to sample values from a conditional distribution $p(\mathbf{y} \mid \mathbf{x})$ that can

1. Mixture components $p(c = i \mid \mathbf{x})$: these form a mixture over the n different components associated with later. typically be obtained by a softmax over an n -dimensional vector that these outputs are positive and sum to 1.
2. Means $\boldsymbol{\mu}^{(i)}(\mathbf{x})$: these indicate the center or mean as Gaussian component and are unconstrained (typically all for these output units). If \mathbf{y} is a d -vector, then the an $n \times d$ matrix containing all n of these d -dimensional these means with maximum likelihood is slightly more learning the means of a distribution with only one component want to update the mean for the component that a observation. In practice, we do not know which component observation. The expression for the negative log-likelihood each example's contribution to the loss for each component that the component produced the example.
3. Covariances $\boldsymbol{\Sigma}^{(i)}(\mathbf{x})$: these specify the covariance matrix i . As when learning a single Gaussian component, we use the matrix to avoid needing to compute determinants. As with of the mixture, maximum likelihood is complicated by the partial responsibility for each point to each mixture component. Gradient descent will automatically follow the correct procedure specification of the negative log-likelihood under the

It has been reported that gradient-based optimization of mixtures (on the output of neural networks) can be unreliable. It gets divisions (by the variance) which can be numerically

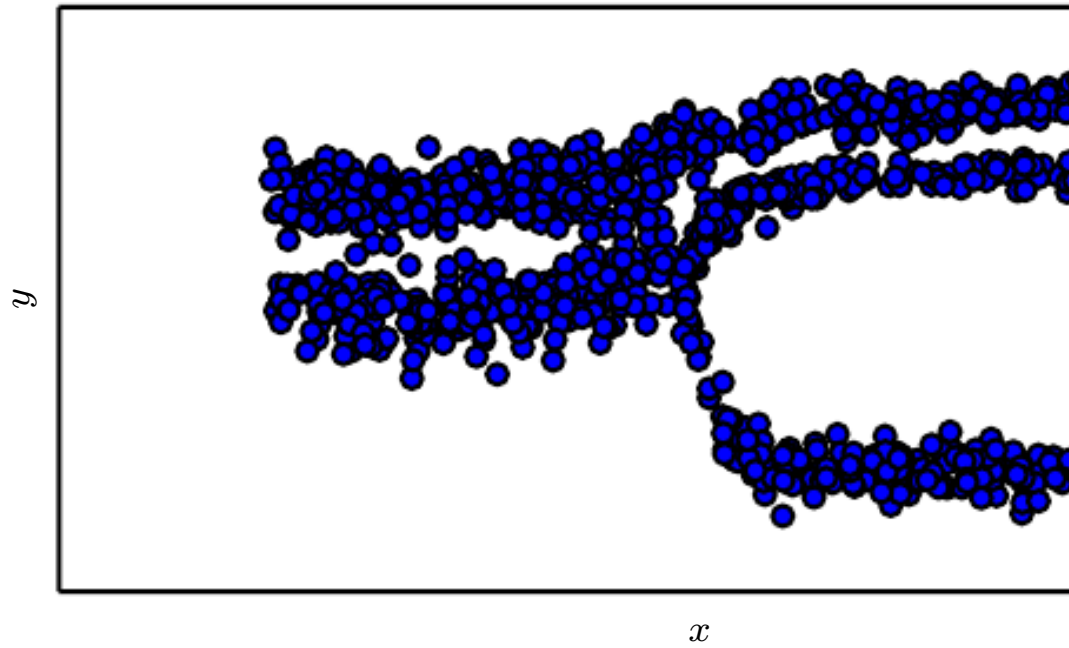


Figure 6.4: Samples drawn from a neural network with a mixture model. The input x is sampled from a uniform distribution, and the output y is sampled from the conditional distribution $p_{\text{model}}(y | x)$. The neural network is able to learn nonlinear mappings from x to y . The parameters of the output distribution are learned from the data. These parameters include the parameters governing which of three mixture components will generate the output, and the parameters for each mixture component. Each mixture component has its own predicted mean and variance. All these aspects of the output distribution are learned with respect to the input x , and to do so in nonlinear ways.

a high degree of quality in these real-valued domains. Another example of a density network is shown in figure 6.4.

In general, we may wish to continue to model larger vector-valued variables, and to impose richer and richer structures on these variables. For example, if we want our neural network to output a sequence of words that forms a sentence, we might continue to use the principle of maximum likelihood applied to our model $p(\mathbf{y} | \mathbf{x})$. In this case, the model we want

networks: how to choose the type of hidden unit to use in the model.

The design of hidden units is an extremely active area of research, yet there are no yet have many definitive guiding theoretical principles.

Rectified linear units are an excellent default choice of hidden unit types. Many other types of hidden units are available. It can be difficult to choose which kind (though rectified linear units are usually an excellent choice). We describe here some of the basic intuitions motivating each type. These intuitions can help decide when to try out which unit. Choosing which will work best is usually impossible. The design process involves trial and error, intuiting that a kind of hidden unit may work well, building a network with that kind of hidden unit and evaluating it on a validation set.

Some of the hidden units included in this list are not differentiable at all input points. For example, the rectified linear function is not differentiable at $z = 0$. This may seem like it invalidates a gradient-based learning algorithm. In practice, gradient descent works well enough for these models to be used for machine learning. The main part because neural network training algorithms do not usually find the minimum of the cost function, but instead merely reduce it. This is shown in figure 4.3. (These ideas are described further in chapter 4.) We do not expect training to actually reach a point where the gradient is zero. For the minima of the cost function to correspond to points where the gradient is zero. Hidden units that are not differentiable are usually nondifferentiable at a small number of points. In general, a function $g(z)$ has a subgradient at z by the slope of the function immediately to the left of z . The supergradient is defined by the slope of the function immediately to the right of z .

these usually do not apply to neural network training. That is, that in practice one can safely disregard the nondifferentiable activation functions described below.

Unless indicated otherwise, most hidden units can be computed by taking a vector of inputs \mathbf{x} , computing an affine transformation, and then applying an element-wise nonlinear function $g(\mathbf{z})$. The units are distinguished from each other only by the choice of the activation function $g(\mathbf{z})$.

6.3.1 Rectified Linear Units and Their Generalization

Rectified linear units use the activation function $g(z) = \max(0, z)$.

These units are easy to optimize because they are so simple. The only difference between a linear unit and a rectified linear unit is that a rectified linear unit outputs zero across half its domain. This makes the gradient of a rectified linear unit remain large whenever the unit is active. The gradients are not only large but also consistent. The second derivative of the rectified linear unit operation is 0 almost everywhere, and the derivative of the rectified linear unit is 1 everywhere that the unit is active. This means that the gradient is more useful for learning than it would be with activation functions that have second-order effects.

Rectified linear units are typically used on top of an affine transformation:

$$\mathbf{h} = g(\mathbf{W}^\top \mathbf{x} + \mathbf{b}).$$

When initializing the parameters of the affine transformation,

rectification fixes $\alpha_i = -1$ to obtain $g(z) = |z|$. It is used from images (Jarrett *et al.*, 2009), where it makes sense to be invariant under a polarity reversal of the input illumination. Rectified linear units are more broadly applicable. A **leaky rectification** (Goodfellow *et al.*, 2013) fixes α_i to a small value like 0.01, while a **parametric ReLU** treats α_i as a learnable parameter (He *et al.*, 2015).

Maxout units (Goodfellow *et al.*, 2013a) generalize further. Instead of applying an element-wise function $g(z)$, they apply it to groups of k values. Each maxout unit then outputs the maximum of these groups:

$$g(\mathbf{z})_i = \max_{j \in \mathbb{G}^{(i)}} z_j,$$

where $\mathbb{G}^{(i)}$ is the set of indices into the inputs for group i . This provides a way of learning a piecewise linear function that is convex in the directions in the input \mathbf{x} space.

A maxout unit can learn a piecewise linear, convex function. Maxout units can thus be seen as *learning the activation function* just the relationship between units. With large enough k , a maxout layer can approximate any convex function with arbitrary fidelity. A maxout layer with two pieces can learn to implement the same function as a traditional layer using the rectified linear activation function, value rectification function, or the leaky or parametric ReLU. A maxout layer can implement a totally different function altogether. The maxout layer can be parametrized differently from any of these other layer types. Its dynamics will be different even in the cases where maxout layer implements the same function of \mathbf{x} as one of the other layer types.

Each maxout unit is now parametrized by k weight vectors.

the past (Goodfellow *et al.*, 2014a).

Rectified linear units and all these generalizations of the principle that models are easier to optimize if their behavior is piecewise linear. This same general principle of using linear behavior to obtain good performance also applies in other contexts besides deep linear networks. In recurrent networks that learn from sequences and produce a sequence of states and outputs, for example, it is easier when some linear computations (with some directional bias and magnitude near 1) are involved. One of the best-performing recurrent architectures, the LSTM, propagates information through the network using a particular straightforward kind of linear activation. This is discussed in section 10.10.

6.3.2 Logistic Sigmoid and Hyperbolic Tangent

Prior to the introduction of rectified linear units, most neural networks used the logistic sigmoid activation function

$$g(z) = \sigma(z)$$

or the hyperbolic tangent activation function

$$g(z) = \tanh(z).$$

These activation functions are closely related because $\tanh(z) = 2\sigma(2z) - 1$.

We have already seen sigmoid units as output units for binary classification, representing the probability that a binary variable is 1. Unlike piecewise linear units, sigmoid units saturate across most of their domain—they saturate at 0 and 1.

$\mathbf{w}^\top \mathbf{U}^\top \mathbf{V}^\top \mathbf{x}$ as long as the activations of the network can be made to be positive. This makes training the tanh network easier.

Sigmoidal activation functions are more common in sequential forward networks. Recurrent networks, many probabilistic autoencoders have additional requirements that rule out linear activation functions and make sigmoidal units more attractive despite the drawbacks of saturation.

6.3.3 Other Hidden Units

Many other types of hidden units are possible but are used

In general, a wide variety of differentiable functions p
Many unpublished activation functions perform just as well a
provide a concrete example, we tested a feedforward network
on the MNIST dataset and obtained an error rate of less th
competitive with results obtained using more conventiona
During research and development of new techniques, it is
different activation functions and find that several variation
perform comparably. This means that usually new hidden u
only if they are clearly demonstrated to provide a significa
hidden unit types that perform roughly comparably to know
as to be uninteresting.

It would be impractical to list all the hidden unit types the literature. We highlight a few especially useful and dis

One possibility is to not have an activation $g(z)$ at all. this as using the identity function as the activation function

these low-rank relationships are often sufficient. Linear hidden units are an effective way of reducing the number of parameters in a network.

Softmax units are another kind of unit that is usually described in section 6.2.2.3) but may sometimes be used as a hidden unit. Softmax units naturally represent a probability distribution over a set of possible values, so they may be used as a kind of switch. Softmax units are usually only used in more advanced architectures that need to manipulate memory, as described in section 10.12.

A few other reasonably common hidden unit types include:

- **Radial basis function** (RBF), unit: $h_i = \exp(-\frac{\|x - c_i\|^2}{2\sigma^2})$. This function becomes more active as x approaches a template c_i . Since the function saturates to 0 for most x , it can be difficult to optimize.
- **Softplus**: $g(a) = \zeta(a) = \log(1 + e^a)$. This is a smooth approximation to the rectifier introduced by Dugas *et al.* (2001) for function approximation and Hinton (2010) for the conditional distributions of undirected graphical models. Glorot *et al.* (2011a) compared the softplus to the rectifier and got better results with the latter. The use of the softplus instead of the rectifier is very counterintuitive—one might expect it to have the performance of the rectifier due to being differentiable everywhere on the real line, but empirically it does not.
- **Hard tanh**. This is shaped similarly to the tanh and the rectifier. Unlike the latter, it is bounded, $g(a) = \max(-1, \min(1, a))$. It was used by Collobert (2004).

layer being a function of the layer that preceded it. In this case, the output of the first layer is given by

$$\mathbf{h}^{(1)} = g^{(1)} \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right);$$

the second layer is given by

$$\mathbf{h}^{(2)} = g^{(2)} \left(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right);$$

and so on.

In these chain-based architectures, the main architectural choices are choosing the depth of the network and the width of each layer. In fact, as we will see, a network with even one hidden layer is sufficient to approximate any continuous function. Deeper networks are often able to use far fewer units per layer, fewer parameters, as well as frequently generalizing to the test set better. However, they can be harder to optimize. The ideal network architecture for a given task is often found via experimentation guided by monitoring the validation set performance.

6.4.1 Universal Approximation Properties and the Bias-Variance Tradeoff

A linear model, mapping from features to outputs via matrix multiplication, by definition represent only linear functions. It has the advantage of being easy to train because many loss functions result in convex optimization problems when applied to linear models. Unfortunately, we often want to approximate nonlinear functions.

At first glance, we might presume that learning a nonlinear function requires designing a specialized model family for the kind of nonlinearity we want to approximate. Fortunately, feedforward networks with hidden layers provide a way to approximate any continuous function.

to another. While the original theorems were first stated for activation functions that saturate for both very negative and very positive arguments, universal approximation theorems have also been proved for a wider class of activation functions, which includes the now common ReLU unit (Leshno *et al.*, 1993).

The universal approximation theorem means that regardless of what function we are trying to learn, we know that a large MLP will be able to approximate that function. We are not guaranteed, however, that the trained MLP is able to *learn* that function. Even if the MLP is able to approximate the function, learning can fail for two different reasons. First, the optimization algorithm for training may not be able to find the value of the parameters that best approximate the desired function. Second, the training algorithm may converge to a function as a result of overfitting. Recall from section 5.2.1 that the universal approximation theorem shows that there is no universally superior machine learning algorithm. Feedforward networks provide a universal system for representing functions in the sense that, given a function, there exists a feedforward network that can approximate the function. There is no universal procedure for examining a function on specific examples and choosing a function that will generalize to a new training set.

According to the universal approximation theorem, there is no limit on how small enough to achieve any degree of accuracy we desire, but there is no way to say how large this network will be. Barron (1993) provides a bound on the size of a single-layer network needed to approximate a bounded function. Unfortunately, in the worst case, an exponential number of hidden units (one with one hidden unit corresponding to each input configuration that is distinguished) may be required. This is easiest to see in the case of a large number of possible binary functions on vectors $\mathbf{v} \in \{0, 1\}^n$.

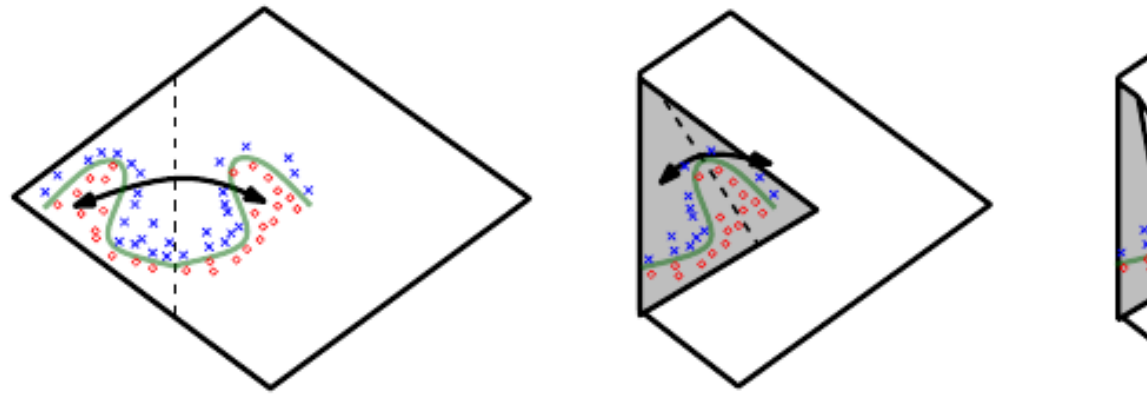


Figure 6.5: An intuitive, geometric explanation of the exponential width of rectifier networks formally by Montufar *et al.* (2014). (Left) An activation unit has the same output for every pair of mirror points in its input space. The axis of symmetry is given by the hyperplane defined by the weights of the unit. The green curve is the decision surface of a simpler pattern across that axis of symmetry. (Center) The space is folded by folding the space around the axis of symmetry. (Right) Another fold can be folded on top of the first (by another downstream unit) to create a more complex decision surface (which is now repeated four times, with two hidden layers). This illustrates the exponential width of rectifier networks (with permission from Montufar *et al.* (2014)).

These results were first proved for models that do not resemble the conventional neural networks used for machine learning but have since been extended to more general models. The first results were for circuits of logic gates (Minsky and Papert, 1969). The work extended these results to linear threshold units with bounded weights (Håstad and Goldmann, 1991; Hajnal *et al.*, 1993), and to networks with continuous-valued activations (Maass, 1992; Maass *et al.*, 1992). Since most neural networks use rectified linear units, Leshno *et al.* (1993) proved that shallow networks with a broad family of non-polynomial activation functions, including rectified linear units, have universal approximation properties. These results do not address the questions of depth or efficiency—whether a sufficiently wide rectifier network could represent any function.

The main theorem in Montufar *et al.* (2014) states that the number of regions carved out by a deep rectifier network with d input nodes and h hidden layers per hidden layer is

$$O \left(\binom{n}{d}^{d(l-1)} n^d \right),$$

that is, exponential in depth l . In the case of maxout network unit, the number of linear regions is

$$O\left(k^{(l-1)+d}\right).$$

Of course, there is no guarantee that the kinds of function applications of machine learning (and in particular for AI)

We may also want to choose a deep model for statistical learning. If we choose a specific machine learning algorithm, we are implicitly encoding a set of prior beliefs we have about what kind of function we want to learn. Choosing a deep model encodes a very general belief that the function we want to learn should involve composition of several simpler functions, which can be interpreted from a representation learning point of view as the learning problem consists of discovering a set of underlying features that can in turn be described in terms of other, simpler features. Alternately, we can interpret the use of a deep architecture as a belief that the function we want to learn is a computer program that involves multiple steps, where each step makes use of the previous steps. In this view, the intermediate outputs are not necessarily factors of variation, but are more analogous to counters or pointers that the network uses to coordinate its processing. Empirically, greater depth does seem to result in better performance.

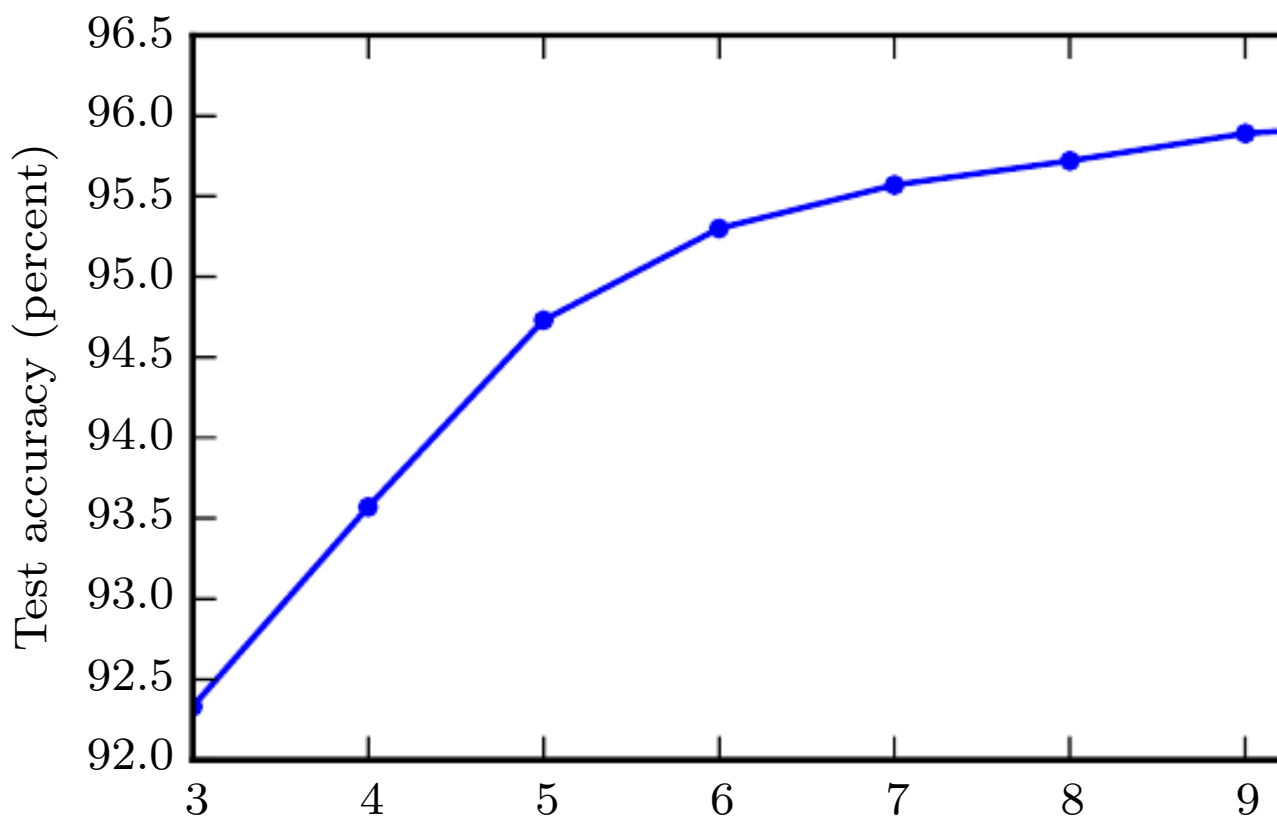


Figure 6.6: Effect of depth. Empirical results showing that deep models are better when used to transcribe multidigit numbers from photographs, from Goodfellow *et al.* (2014d). The test set accuracy consistently increases with depth. See figure 6.7 for a control experiment demonstrating that increasing model size do not yield the same effect.

Many neural network architectures have been developed. Specialized architectures for computer vision called convolutional neural networks are described in chapter 9. Feedforward networks may also be compared to recurrent neural networks for sequence processing, described in chapter 10. Each has their own architectural considerations.

In general, the layers need not be connected in a chain, which is the most common practice. Many architectures build a main chain

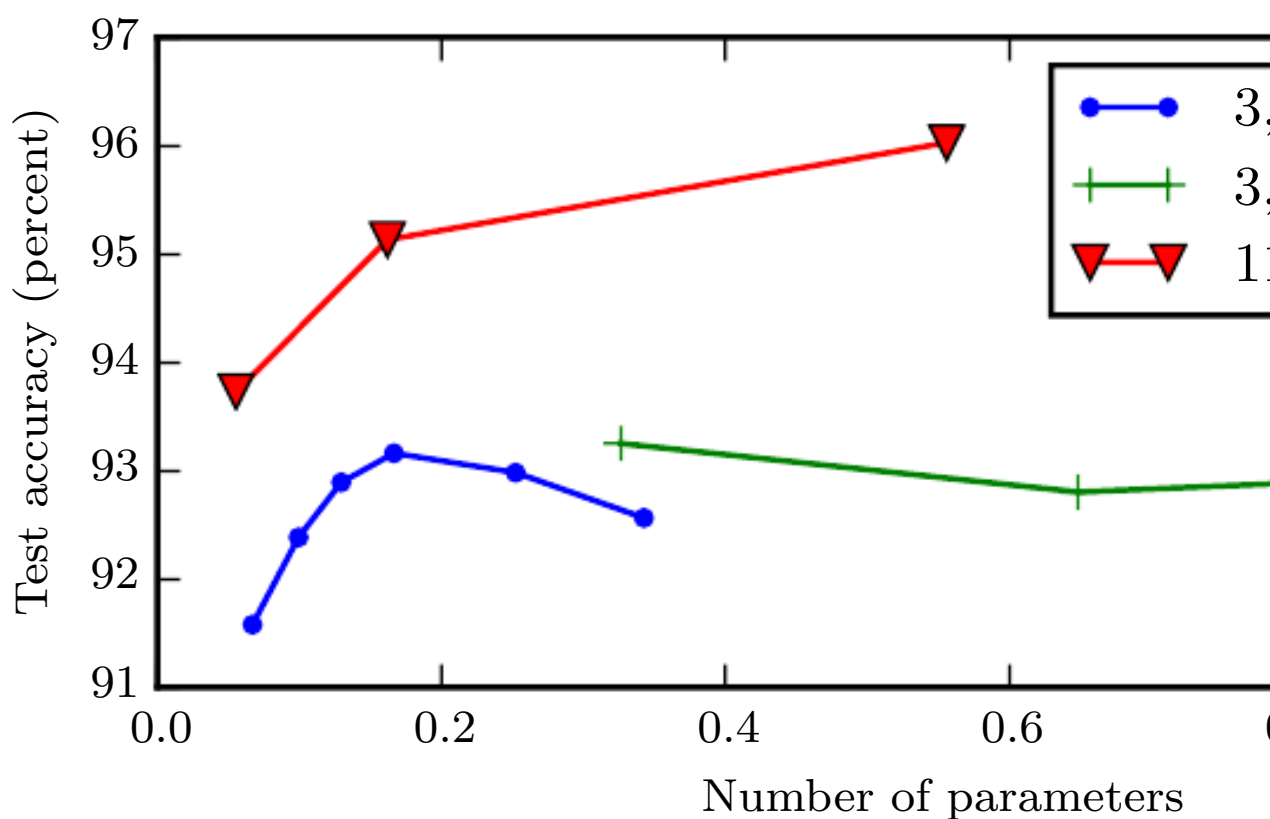


Figure 6.7: Effect of number of parameters. Deeper models benefit from more parameters. This is not merely because the model is larger. This experiment (2014d) shows that increasing the number of parameters in layers of width 1024 without increasing their depth is not nearly as effective at increasing accuracy as illustrated in this figure. The legend indicates the depth of each curve and whether the curve represents variation in the size of the input or the fully connected layers. We observe that shallow models benefit from around 20 million parameters while deep ones can benefit from

networks, described in chapter 9, use specialized patterns that are very effective for computer vision problems. In this chapter we give more specific advice concerning the architecture of a neural network. In subsequent chapters we develop the particular architectures that have been found to work well for different application domains.

6.5 Back-Propagation and Other Differentiation Algorithms

When we use a feedforward neural network to accept an input \mathbf{x} and produce an output $\hat{\mathbf{y}}$, information flows forward through the network. The initial information is the input \mathbf{x} that then propagates up to the hidden layers and finally produces $\hat{\mathbf{y}}$. This is called **forward propagation**. Forward propagation can continue onward until it produces the final output. The **back-propagation** algorithm (Rumelhart *et al.*, 1986), often abbreviated **backprop**, allows the information from the cost to then flow backward through the network in order to compute the gradient.

Computing an analytical expression for the gradient is often difficult. Numerically evaluating such an expression can be computationally expensive. The back-propagation algorithm does so using a simple and efficient method.

The term back-propagation is often misunderstood as a learning algorithm for multi layer neural networks. Actually, it refers only to the method for computing the gradient, which, such as stochastic gradient descent, is used to perform learning. Furthermore, back-propagation is often misunderstood as a method for multi layer neural networks, but in principle it can compute derivatives for any function.

f with multiple outputs. We restrict our description here to the most commonly used case, where f has a single output.

6.5.1 Computational Graphs

So far we have discussed neural networks with a relatively informal description. To describe the back-propagation algorithm more precisely, we need a more precise **computational graph** language.

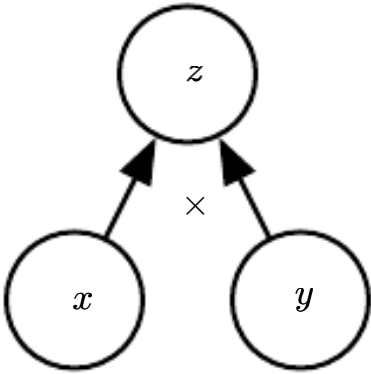
Many ways of formalizing computation as graphs are possible.

Here, we use each node in the graph to indicate a variable. A variable can be a scalar, vector, matrix, tensor, or even a variable of an arbitrary type.

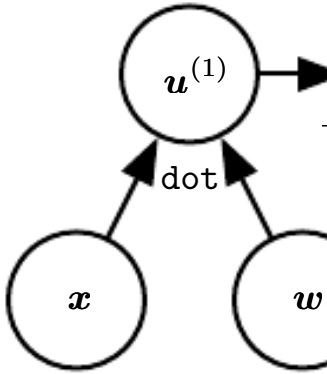
To formalize our graphs, we also need to introduce the concept of an operation. An operation is a simple function of one or more variables. Each operation is accompanied by a set of allowable operations. Functions not in this set may be described by composing operations together.

Without loss of generality, we define an operation to have a single output variable. This does not lose generality because the output can be a vector of multiple entries, such as a vector. Software implementations usually support operations with multiple outputs, but we restrict our description because it introduces many extra details that are not necessary for conceptual understanding.

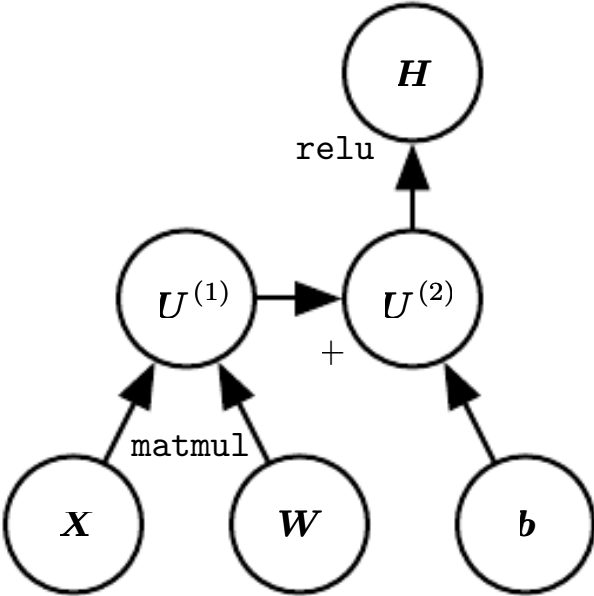
If a variable y is computed by applying an operation to a set of variables, we draw a directed edge from x to y . We sometimes annotate the edge with the name of the operation applied, and other times omit the name if the operation is clear from context.



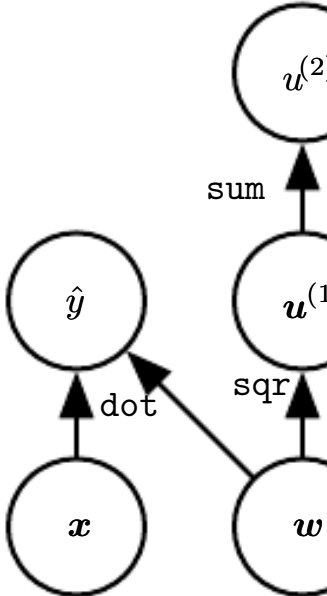
(a)



(b)



(c)



(d)

the chain rule states that

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}.$$

We can generalize this beyond the scalar case. Suppose g maps from \mathbb{R}^m to \mathbb{R}^n , and f maps from \mathbb{R}^n to \mathbb{R} . If $\mathbf{y} = g(\mathbf{x})$, then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}.$$

In vector notation, this may be equivalently written as

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z,$$

where $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is the $n \times m$ Jacobian matrix of g .

From this we see that the gradient of a variable \mathbf{x} can be computed by multiplying a Jacobian matrix $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ by a gradient $\nabla_{\mathbf{y}} z$. The back-propagation algorithm is simply a matter of performing such a Jacobian-gradient product for each operation in the network.

Usually we apply the back-propagation algorithm to tensors of arbitrary dimensionality, not merely to vectors. Conceptually, this is just back-propagation with vectors. The only difference is how the tensors are represented. A tensor is just a multi-dimensional array of numbers, which can be thought of as a vector flattened into a grid to form a tensor. We could imagine flattening a vector before we run back-propagation, computing a gradient, and then reshaping the gradient back into a tensor. In practice, back-propagation is still just multiplying Jacobians by gradients.

To denote the gradient of a value z with respect to a tensor \mathbf{X} , we write $\nabla_{\mathbf{X}} z$, just as if \mathbf{X} were a vector. The indices into \mathbf{X} now have multiple dimensions.

6.5.3 Recursively Applying the Chain Rule to

Using the chain rule, it is straightforward to write down an expression for the gradient of a scalar with respect to any node in the computational graph that produced that scalar. Actually evaluating that expression in a naive way introduces some extra considerations.

Specifically, many subexpressions may be repeated several times in the overall expression for the gradient. Any procedure that naively computes the gradient will need to choose whether to store these subexpressions or recompute them several times. An example of how these repeated subexpressions can arise is shown in figure 6.9. In some cases, computing the same subexpression multiple times can be wasteful. For complicated graphs, there can be exponentially many wasted computations, making a naive implementation of the chain rule inefficient. In other cases, computing the same subexpression twice can be useful to reduce memory consumption at the cost of higher runtime.

We begin with a version of the back-propagation algorithm that computes the actual gradient computation directly (algorithm 6.2 along with the associated forward computation), in the order it will actually be executed, as opposed to the recursive application of chain rule. One could either view this as a sequence of computations or view the description of the algorithm as a traversal of the computational graph for computing the back-propagation. The latter formulation does not make explicit the manipulation and transformation of the symbolic graph that performs the gradient computation. This is presented in section 6.5.6, with algorithm 6.5, where we allow for nodes that contain arbitrary tensors.

First consider a computational graph describing how to compute a scalar $u^{(n)}$ (say, the loss on a training example). This scalar

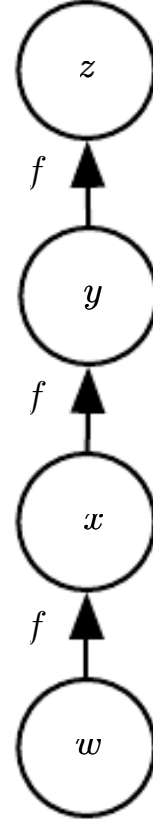


Figure 6.9: A computational graph that results in repeated subexp the gradient. Let $w \in \mathbb{R}$ be the input to the graph. We use the s as the operation that we apply at every step of a chain: $x = f$ To compute $\frac{\partial z}{\partial w}$, we apply equation 6.44 and obtain:

$$\begin{aligned}
 & \frac{\partial z}{\partial w} \\
 &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \\
 &= f'(y) f'(x) f'(w) \\
 &= f'(f(f(w))) f'(f(w)) f'(w).
 \end{aligned}$$

That algorithm specifies the forward propagation computation in a graph \mathcal{G} . To perform back-propagation, we can construct a graph that depends on \mathcal{G} and adds to it an extra set of edges. The subgraph \mathcal{B} with one node per node of \mathcal{G} . Computation in \mathcal{B} is the reverse of the order of computation in \mathcal{G} , and each node $u^{(i)}$ in \mathcal{B} has a derivative $\frac{\partial u^{(n)}}{\partial u^{(i)}}$ associated with the forward graph node $u^{(i)}$. The chain rule with respect to scalar output $u^{(n)}$:

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i: j \in Pa(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

as specified by algorithm 6.2. The subgraph \mathcal{B} contains exactly one edge from node $u^{(j)}$ to node $u^{(i)}$ of \mathcal{G} . The edge from $u^{(j)}$ to $u^{(i)}$ is labeled with the computation of $\frac{\partial u^{(i)}}{\partial u^{(j)}}$. In addition, a dot product is performed between the gradient already computed with respect to node $u^{(i)}$ and the vector containing the partial derivatives $\frac{\partial u^{(i)}}{\partial u^{(j)}}$ for all nodes $u^{(i)}$. To summarize, the amount of computation required for the back-propagation scales linearly with the number of nodes. The computation for each edge corresponds to computing a partial derivative (node with respect to one of its parents) as well as performing a dot product and one addition. Below, we generalize this analysis to tensors. A tensor is just a way to group multiple scalar values in the same

Algorithm 6.1 A procedure that performs the computation of $u^{(1)}$ to $u^{(n)}$ to an output $u^{(n)}$. This defines a computational graph where each node $u^{(i)}$ computes numerical value $u^{(i)}$ by applying a function $f^{(i)}$ to its parents $u^{(j)}$.

6.5.4 Back-Propagation Computation in Fully

To clarify the above definition of the back-propagation computation, we show the specific graph associated with a fully-connected multi-layer perceptron.

Algorithm 6.3 first shows the forward propagation, which computes the supervised loss $L(\hat{\mathbf{y}}, \mathbf{y})$ associated with a single (input, target) pair (\mathbf{x}, \mathbf{y}) , with $\hat{\mathbf{y}}$ the output of the neural network when \mathbf{x} is the input.

Algorithm 6.4 then shows the corresponding computation for the backward pass, applying the back-propagation algorithm to this graph.

Algorithms 6.3 and 6.4 are demonstrations chosen to be simple and easy to understand. However, they are specialized to one specific architecture.

Modern software implementations are based on the general back-propagation described in section 6.5.6 below, which can accommodate any computational graph by explicitly manipulating a data structure for the graph and the computation.

Algorithm 6.3 Forward propagation through a typical deep neural network to compute the cost function. The loss $L(\hat{\mathbf{y}}, \mathbf{y})$ is computed by passing the input \mathbf{x} through the network to obtain $\hat{\mathbf{y}}$ and on the target \mathbf{y} (see section 6.2.1.1 for examples). To compute the total cost J , the loss may be added to a regularization term. The network contains all the parameters (weights and biases). Algorithm 6.4 computes gradients of J with respect to parameters \mathbf{W} and \mathbf{b} . This demonstration uses only a single input example \mathbf{x} . Practically, one uses a minibatch. See section 6.5.7 for a more realistic demonstration.

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the network layers. $\mathbf{b}^{(i)}$ are the bias vectors.

Algorithm 6.4 Backward computation for the deep network. This algorithm, which uses, in addition to the input \mathbf{x} , a target \mathbf{y} , yields the gradients on the activations $\mathbf{a}^{(k)}$ for each layer, from the output layer and going backwards to the first hidden layer. This gradient, which can be interpreted as an indication of how each layer's output contributes to reduce error, one can obtain the gradient on the parameters. These gradients on weights and biases can be immediately used for a stochastic gradient update (performing the update right after the gradients are computed) or used with other gradient-based optimization algorithms.

After the forward computation, compute the gradient on the output:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ **do**

 Convert the gradient on the layer's output into a gradient on the nonlinearity activation (element-wise multiplication if the activation is nonlinear):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

 Compute gradients on weights and biases (including the regularization term where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

 Propagate the gradients w.r.t. the next lower-level hidden layer:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

6.5.5 Symbol-to-Symbol Derivatives

Algebraic expressions and computational graphs both operate on symbolic variables that do not have specific values. These algebraic

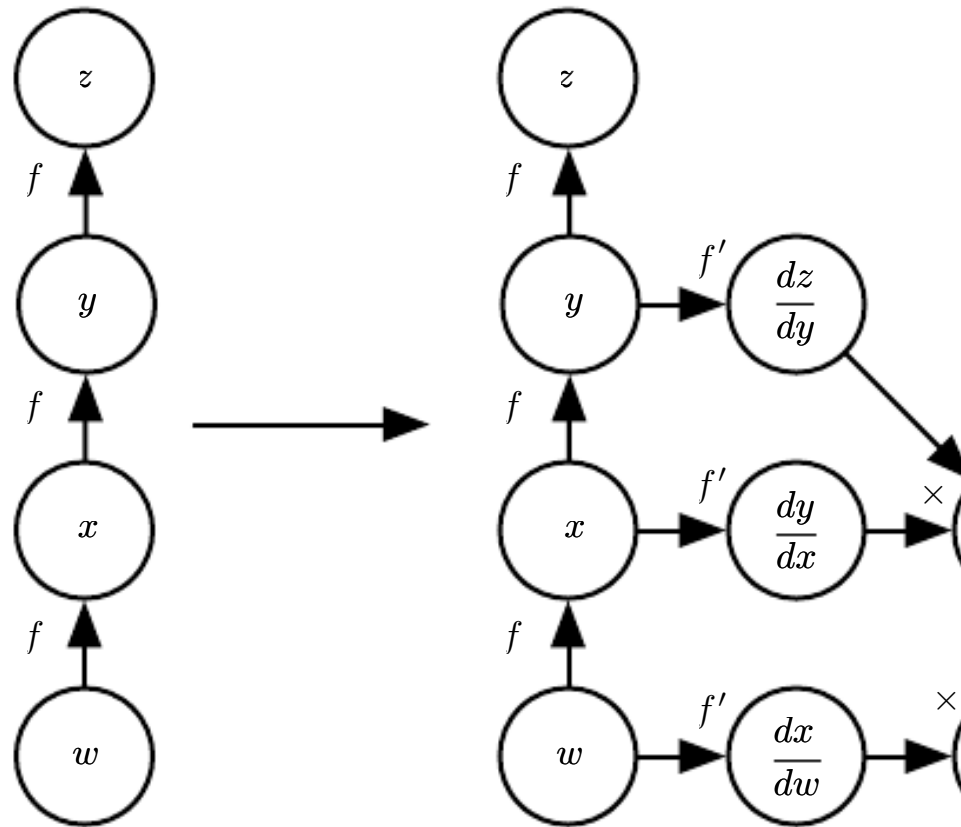


Figure 6.10: An example of the symbol-to-symbol approach to computing derivatives. In this approach, the back-propagation algorithm does not need to use specific numeric values. Instead, it adds nodes to a computational graph to compute these derivatives. A generic graph evaluation engine can then compute derivatives for any specific numeric values. (Left) In this example, the graph represents $z = f(f(f(w)))$. (Right) We run the back-propagation algorithm to construct the graph for the expression corresponding to $\frac{dz}{dw}$. We do not explain how the back-propagation algorithm works. The purpose is to show what the desired result is: a computational graph with a symbolic expression for the derivative.

This is the approach taken by Theano (Bergstra *et al.*, 2010; Bergstra and Breuleux, 2013) and TensorFlow (Abadi *et al.*, 2015). An example of how it works is shown in figure 6.10. The primary advantage of this approach is that the derivative is described in the same language as the original expression. This

performing exactly the same computations as are done in the symbol-to-symbol approach. The key difference is that the latter approach does not expose the graph.

6.5.6 General Back-Propagation

The back-propagation algorithm is very simple. To compute the gradient of a scalar z with respect to one of its ancestors \mathbf{x} in the graph, we assume that the gradient with respect to z is given by $\frac{dz}{dz} = 1$. We then compute the gradient with respect to each parent of z in the graph by multiplying the current gradient by the Jacobian of the operation that produced z . By multiplying by Jacobians, traveling backward through the graph, we reach \mathbf{x} . For any node that may be reached by going backward along two or more paths, we simply sum the gradients arriving from each path at that node.

More formally, each node in the graph \mathcal{G} corresponds to a variable. At maximum generality, we describe this variable as being a vector. Vectors in general can have any number of dimensions. They subsume scalars and matrices.

We assume that each variable \mathbf{V} is associated with the following information:

- **get_operation(\mathbf{V})**: This returns the operation that produced \mathbf{V} , represented by the edges coming into \mathbf{V} in the computation graph. There may be a Python or C++ class representing the operation, and the **get_operation** function. Suppose \mathbf{C} is created by matrix multiplication, $\mathbf{C} = \mathbf{A}\mathbf{B}$. The **get_operation** function returns a pointer to an instance of the corresponding

of a scalar z with respect to \mathbf{C} is given by \mathbf{G} . The matrix \mathbf{A} is responsible for defining two back-propagation rules, on \mathbf{A} and \mathbf{B} arguments. If we call the `bprop` method to request the gradient on \mathbf{A} given that the gradient on the output is \mathbf{G} , then the matrix multiplication operation must state that the gradient on \mathbf{A} is given by $\mathbf{G}\mathbf{B}^\top$. Likewise, if we call the `bprop` method to request the gradient with respect to \mathbf{B} , then the matrix operation is responsible for returning $\mathbf{G}\mathbf{A}^\top$. The `bprop` method and specifying that the desired gradient is on \mathbf{A} or \mathbf{B} in the back-propagation algorithm itself does not need to know any details of the operation, it only needs to call each operation's `bprop` rules with the right arguments. The `op.bprop(inputs, \mathbf{X} , \mathbf{G})` must return

$$\sum_i (\nabla_{\mathbf{X}} \text{op.f}(\text{inputs})_i) G_i,$$

which is just an implementation of the chain rule as expressed in equation 6.1. Here, `inputs` is a list of inputs that are supplied to the operation, `op.f` is the mathematical function that the operation implements, \mathbf{X} is the input on which we wish to compute, and \mathbf{G} is the gradient on the output of the operation.

The `op.bprop` method should always pretend that all inputs are independent from each other, even if they are not. For example, if the operation takes two copies of x to compute x^2 , the `op.bprop` method should return the derivative with respect to both inputs. The back-propagation algorithm should add both of these arguments together to obtain $2x$, which is the derivative on x .

Software implementations of back-propagation usually provide common operations and their `bprop` methods, so that users of deep learning frameworks are able to back-propagate through graphs built using common

Algorithm 6.5 The outermost skeleton of the back-propagation portion does simple setup and cleanup work. Most of the inner loop is in the `build_grad` subroutine of algorithm 6.6.

Require: \mathbb{T} , the target set of variables whose gradients must be computed

Require: \mathcal{G} , the computational graph

Require: z , the variable to be differentiated

Let \mathcal{G}' be \mathcal{G} pruned to contain only nodes that are ancestors of nodes in \mathbb{T} .

Initialize `grad_table`, a data structure associating tensors with gradients

`grad_table`[z] $\leftarrow 1$

for \mathbf{V} in \mathbb{T} **do**

`build_grad`(\mathbf{V} , \mathcal{G} , \mathcal{G}' , `grad_table`)

end for

Return `grad_table` restricted to \mathbb{T}

Algorithm 6.6 The inner loop subroutine `build_grad`(\mathbf{V} , \mathcal{G} , \mathcal{G}' , `grad_table`) is the back-propagation algorithm, called by the back-propagation algorithm in algorithm 6.5.

Require: \mathbf{V} , the variable whose gradient should be added to the gradient of \mathbf{V}

Require: \mathcal{G} , the graph to modify

Require: \mathcal{G}' , the restriction of \mathcal{G} to nodes that participate in the computation of \mathbf{V}

Require: `grad_table`, a data structure mapping nodes to their gradients

if \mathbf{V} is in `grad_table` **then**

 Return `grad_table`[\mathbf{V}]

end if

$i \leftarrow 1$

for \mathbf{C} in `get_producers`(\mathbf{V} , \mathcal{G}') **do**

of operations executed. Keep in mind here that we refer to a node as the fundamental unit of our computational graph, which might consist of several arithmetic operations (for example, we might have a graph node representing a multiplication as a single operation). Computing a gradient for a node will never execute more than $O(n^2)$ operations or store the results of $O(n^2)$ operations. Here we are counting operations in the computational graph, not the individual operations executed by the underlying hardware. Remember that the runtime of each operation may be highly variable. For example, multiplying two matrices that each contain millions of entries might be a single operation in the graph. We can see that computing a gradient for a node has at most $O(n^2)$ operations because the forward propagation stage visits all n nodes in the original graph (depending on which values are used, we may not need to execute the entire graph). The back-propagation stage adds one Jacobian-vector product, which should be expressed as a single edge in the original graph. Because the computational graph has at most $O(n^2)$ edges, the back-propagation graph it has at most $O(n^2)$ edges. For the kinds of graphs that we encounter in practice, the situation is even better. Most neural networks are roughly chain-structured, causing back-propagation to have a cost that is much better than the naive approach, which might need to execute all nodes. This potentially exponential cost can be seen by expressing the recursive chain rule (equation 6.53) nonrecursively:

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{\substack{\text{path}(u^{(\pi_1)}, u^{(\pi_2)}, \dots, u^{(\pi_t)}), \\ \text{from } \pi_1=j \text{ to } \pi_t=n}} \prod_{k=2}^t \frac{\partial u^{(\pi_k)}}{\partial u^{(\pi_{k-1})}}$$

Since the number of paths from node j to node n can grow exponentially with the length of these paths, the number of terms in the above sum

6.5.7 Example: Back-Propagation for MLP Tr

As an example, we walk through the back-propagation algorithm to train a multilayer perceptron.

Here we develop a very simple multilayer perceptron with one hidden layer. To train this model, we will use minibatch stochastic gradient descent. The back-propagation algorithm is used to compute the gradients of the cost function for a single minibatch. Specifically, we use a minibatch of examples from the training set formatted as a design matrix \mathbf{X} and a vector of associated target labels \mathbf{y} . The network computes a layer of hidden features $\mathbf{H} = \mathbf{XW}^{(1)}$. To simplify the presentation we do not use biases in this model. The graph language includes a `relu` operation that can compute the element-wise maximum of zero and its argument. The predictions of the unnormalized log probabilities are given by $\mathbf{HW}^{(2)}$. We assume that our graph language includes an operation that computes the cross-entropy between the target distribution and the predicted distribution defined by these unnormalized log probabilities. This cross-entropy defines the cost J_{MLE} . Minimizing this cross-entropy is equivalent to maximum likelihood estimation of the classifier. However, to make this model more robust we also include a regularization term. The total cost

$$J = J_{\text{MLE}} + \lambda \left(\sum_{i,j} \left(W_{i,j}^{(1)} \right)^2 + \sum_{i,j} \left(W_{i,j}^{(2)} \right)^2 \right)$$

consists of the cross-entropy and a weight decay term where λ is a regularization parameter. The computational graph is illustrated in figure 6.11.

The computational graph for the gradient of this example is shown in figure 6.12. It would be tedious to draw or to read. This demonstrates the power of the back-propagation algorithm, which is that it can compute the gradients of the cost function with respect to the weights of the network.

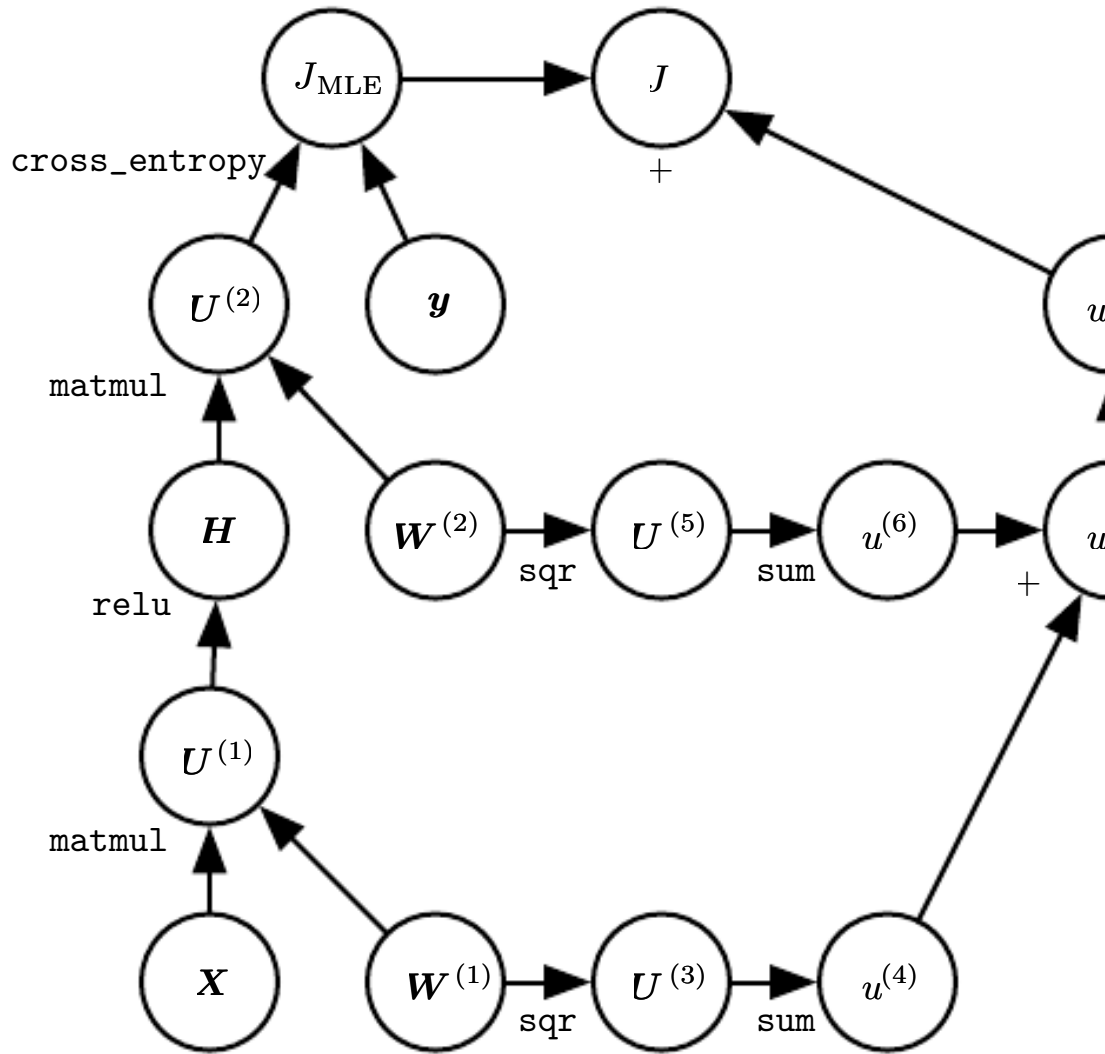


Figure 6.11: The computational graph used to compute the cost of a single-layer MLP using the cross-entropy loss and weight decay.

explore two different branches. On the shorter branch, the gradient on $W^{(2)}$, using the back-propagation rule for the matrix multiplication operation. The other branch continues the chain descending further along the network. First, the back-propagation computes $\nabla_H J = G W^{(2)\top}$ using the back-propagation rule for the matrix multiplication operation. Next, the **relu** operation uses the back-propagation rule to zero out components of the gradient corresponding to the

This value is stored from the time it is computed until it is returned to the same point. The memory cost is thus $O(n_h)$, where n_h is the number of examples in the minibatch and n_h is the number of hidden units.

6.5.8 Complications

Our description of the back-propagation algorithm here is simplified. Real-world implementations actually used in practice.

As noted above, we have restricted the definition of a tensor to a single value. Most software implementations support operations that can return more than one tensor. For example, to compute both the maximum value in a tensor and the index of the maximum value, it is best to compute both in a single pass through memory, so as to avoid the need to implement this procedure as a single operation with two outputs.

We have not described how to control the memory usage of back-propagation. Back-propagation often involves summation of many tensors. In the naive approach, each of these tensors would be computed separately, and all of them would be added in a second step. The naive approach has a high memory bottleneck that can be avoided by maintaining a buffer and adding each value to that buffer as it is computed.

Real-world implementations of back-propagation also need to handle different data types, such as 32-bit floating point, 64-bit floating point, and integer. The policy for handling each of these types takes special care.

Some operations have undefined gradients, and it is important to handle these cases and determine whether the gradient requested by the operation is defined.

Various other technicalities make real-world differentiation difficult.

mode accumulation. Other approaches evaluate the subrule in different orders. In general, determining the order results in the lowest computational cost is a difficult problem. The sequence of operations to compute the gradient is NP-complete in the sense that it may require simplifying algebraic expressions to an expensive form.

For example, suppose we have variables p_1, p_2, \dots, p_n representing probabilities and variables z_1, z_2, \dots, z_n representing unnormalized log probabilities. We define

$$q_i = \frac{\exp(z_i)}{\sum_i \exp(z_i)},$$

where we build the softmax function out of exponentiation, summation, and division operations, and construct a cross-entropy loss $J = -\sum_i p_i \log q_i$. A mathematician can observe that the derivative of J with respect to z_i has the simple form: $q_i - p_i$. The back-propagation algorithm is not aware of this, so it computes the gradient the hard way and will instead explicitly propagate the gradient through the logarithm and exponentiation operations in the original expression. Libraries such as Theano (Bergstra *et al.*, 2010; Bastien *et al.*, 2012) perform some kinds of algebraic substitution to improve over the pure back-propagation algorithm.

When the forward graph \mathcal{G} has a single output node and the partial derivative $\frac{\partial u^{(i)}}{\partial u^{(j)}}$ can be computed with a constant amount of computation, Algorithm 6.2 guarantees that the number of computations for the gradient is of the same order as the number of computations for the forward pass. This can be seen in algorithm 6.2, because each local partial derivative can be computed only once along with an associated multiplication. This is the recursive chain-rule formulation (equation 6.53). The

accumulation has been proposed for obtaining real-time computation in recurrent networks, for example ([Williams and Zipser, 1998](#)) avoids the need to store the values and gradients for the whole computation for memory. The relationship between forward and backward mode is analogous to the relationship between left- and right-multiplying a sequence of matrices, such as

$$\mathbf{A}\mathbf{B}\mathbf{C}\mathbf{D},$$

where the matrices can be thought of as Jacobian. For example, if \mathbf{D} is a column vector while \mathbf{A} has many rows, the graph will have a single output and starting the multiplications from the end and going backwards involves only matrix-vector products. This order corresponds to the backward mode. Starting to multiply from the left would involve a series of matrix-matrix products which makes the whole computation much more expensive. If \mathbf{D} has many columns, however, it is cheaper to run the multiplications in the order corresponding to the forward mode.

In many communities outside machine learning, it is more common to use automatic differentiation software that acts directly on traditional programming code, such as Python or C code, and automatically generates the gradients for functions written in these languages. In the deep learning community, computational graphs are usually represented by explicit code generated by specialized libraries. The specialized approach has the advantage of allowing the library developer to define the `bprop` methods for every operation, allowing the user of the library to only those operations that have been implemented. The specialized approach also has the benefit of allowing custom rules to be developed for each operation, enabling the developer to optimize for speed or stability in nonobvious ways that an automatic procedure

These libraries use the same kind of data structure to describe derivatives as they use to describe the original function being differentiated, which means that the symbolic differentiation machinery can be reused.

In the context of deep learning, it is rare to compute a second-order derivative of a scalar function. Instead, we are usually interested in products of the Hessian matrix. If we have a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then the Hessian is a $n \times n$ matrix. In typical deep learning applications, n will be the number of parameters in the model, which could easily number in the billions. The entire Hessian is thus infeasible to even represent.

Instead of explicitly computing the Hessian, the typical approach is to use **Krylov methods**. Krylov methods are a set of iterative algorithms for performing various operations, such as approximately inverting a matrix, finding approximations to its eigenvectors or eigenvalues, without needing to store more than other than matrix-vector products.

To use Krylov methods on the Hessian, we only need to compute the product between the Hessian matrix \mathbf{H} and an arbitrary vector \mathbf{v} . A straightforward technique (Christianson, 1992) for doing so is

$$\mathbf{H}\mathbf{v} = \nabla_{\mathbf{x}} \left[(\nabla_{\mathbf{x}} f(\mathbf{x}))^\top \mathbf{v} \right].$$

Both gradient computations in this expression may be computed using the appropriate software library. Note that the outer gradient is the gradient of a function of the inner gradient expression.

If \mathbf{v} is itself a vector produced by a computational graph, we can specify that the automatic differentiation software should not differentiate the graph that produced \mathbf{v} .

While computing the Hessian is usually not advisable, it

have long been used to solve optimization problems in closed form. Gradient descent was not introduced as a technique for iteratively approximating solutions to optimization problems until the nineteenth century (Cauchy, 1846).

Beginning in the 1940s, these function approximation techniques were used to motivate machine learning models such as the perceptron. Early machine learning models were based on linear models. Critics including Marvin Minsky pointed out several of the flaws of the linear model family, such as its inability to represent the XOR function, which led to a backlash against the entire new field of machine learning.

Learning nonlinear functions required the development of more powerful models like the perceptron and a means of computing the gradient through successive layers. Applications of the chain rule based on dynamic programming were used in the 1960s and 1970s, mostly for control applications (Kemeny and Denham, 1961; Dreyfus, 1962; Bryson and Ho, 1969; Dreyfus, 1972) and sensitivity analysis (Linnainmaa, 1976). Werbos (1981) pioneered the application of these techniques to training artificial neural networks. The idea of back-propagation in practice after being independently rediscovered in different contexts (Parker, 1985; Rumelhart *et al.*, 1986a). The book **Parallel Distributed Processing** presented the results of some of the first successful applications of back-propagation in a chapter (Rumelhart *et al.*, 1986b) that led to the popularization of back-propagation and initiated a vigorous research program in multilayer neural networks. The ideas put forward in this book, particularly by Rumelhart and Hinton, go much beyond the traditional view of neural networks. They include crucial ideas about the possible computational functions of the brain. Several central aspects of cognition and learning, which are central to the “connectionism” because of the importance this school of thought places on connections between neurons as the locus of learning and memory, are included in these ideas include the notion of distributed representation.

number of algorithmic changes have also improved the performance of feedforward networks noticeably.

One of these algorithmic changes was the replacement of the mean squared error loss with the cross-entropy family of loss functions. Mean squared error was popular in the 1980s and 1990s but was gradually replaced by cross-entropy. The principle of maximum likelihood as ideas spread between the statistics and the machine learning community. The use of cross-entropy loss significantly improved the performance of models with sigmoid and softmax activation functions that had previously suffered from saturation and slow learning near the boundaries of the squared error loss.

The other major algorithmic change that has greatly improved the performance of feedforward networks was the replacement of sigmoid hidden units with linear hidden units, such as rectified linear units. Rectification, the ReLU activation function, was introduced in early neural network models and discussed in the context of the cognitron and neocognitron (Fukushima, 1975, 1980). Early models did not use rectified linear units but instead applied rectification to the output. Despite the early popularity of rectification, it was largely ignored in the 1980s, perhaps because sigmoids perform better when the output is very small. As of the early 2000s, rectified linear units were not the standard, but a somewhat superstitious belief that activation functions with sharp points must be avoided. This began to change in about 2006 when LeCun et al. (2006) observed that “using a rectifying nonlinearity is the single most effective factor in improving the performance of a recognition system,” and it became a key factor of neural network architecture design.

For small datasets, Jarrett *et al.* (2009) observed that the use of rectified linearities is even more important than learning the weights. In fact, random weights are sufficient to propagate useful information through the network.

are completely inactive. (2) For some inputs, a biological response is proportional to its input. (3) Most of the time, biological systems are in a regime where they are inactive (i.e., they should have **sparsely** active units).

When the modern resurgence of deep learning began, feedforward networks continued to have a bad reputation. From about 1990 to 2006, it was widely believed that feedforward networks would not perform well without assistance by other models, such as probabilistic models. Today, it is clear that with the right resources and engineering practices, feedforward networks can perform very well. Today, gradient-based learning in feedforward networks is used as a tool to develop probabilistic models, such as the variational autoencoders and generative adversarial networks, described in chapter 11. Feedforward networks were viewed as an unreliable technology that must be supported by other models. Today, gradient-based learning in feedforward networks has been shown to be a powerful technology that can be applied to many other machine learning tasks. In 2006, the community used unsupervised learning to support supervised learning, and now, ironically, it is more common to use supervised learning to support unsupervised learning.

Feedforward networks continue to have unfulfilled potential. We expect they will be applied to many more tasks, and that advances in algorithms and model design will improve their performance. This chapter has primarily described the neural network family. In the subsequent chapters, we turn to how to use these models—how to design them, how to train them.