

Topic 3: Neural Networks

3.1 Introduction

Neural networks are arguably the main reason why Machine Learning is making a big splash. Similar to logistic regression and random forests, their main purpose is prediction/classification. Classical applications of neural networks include:

1. Image classification. For example, determining the digit contained in an image, as in Figure 3.1.
2. Natural language processing. For example, interpreting voice commands, like Siri and Alexa do.
3. Stock market prediction.

Each of these problems can be rephrased mathematically as finding an unknown (and highly complex) function f^* , such that given certain data \mathbf{x} , $f^*(\mathbf{x})$ gives the desired classification/prediction. In our examples:

1. Image classification. In our example, \mathbf{x} is a vectorized image of a digit, and $f^*(\mathbf{x}) \in \{0, \dots, 9\}$ is the digit depicted in such figure. For the images $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_5$ in Figure 3.1, $f^*(\mathbf{x}_1) = 5, f^*(\mathbf{x}_2) = 0, \dots, f^*(\mathbf{x}_5) = 9$.
2. Natural language processing. In our example, \mathbf{x} is a vector containing the values in a voice signal, as in Figure 3.2, and $f^*(\mathbf{x})$ is the interpretation, in this case “turn off the t.v.”.
3. Stock market prediction. Here \mathbf{x} could be a sequence of stock market prices at times $t = 1, \dots, T$, and $f^*(\mathbf{x})$ could be the prediction of such stock market price at time $T + 1$.

All of these tasks can be tremendously challenging. For example, whether an image contains a 0, or a 1, or any other digit, depends not only on the values of isolated pixels, but on the way that pixels *interact* with one another in complex manners.

The main idea behind ANNs is to use a sequence of simpler functions that interact with one another in a networked way, so that combined, they approximate f^* with arbitrary precision.

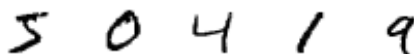


Figure 3.1: Digit images from the MNIST dataset. The goal is to automatically determine the digit contained in each image, in this case $\{5, 0, 4, 1, 9\}$.

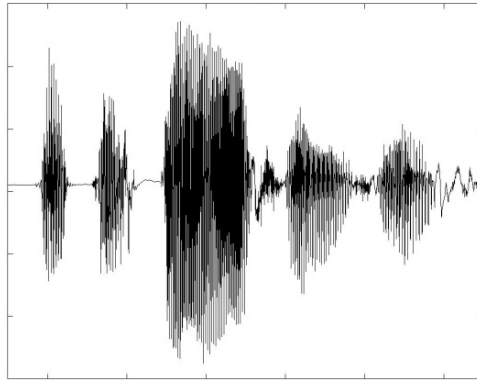
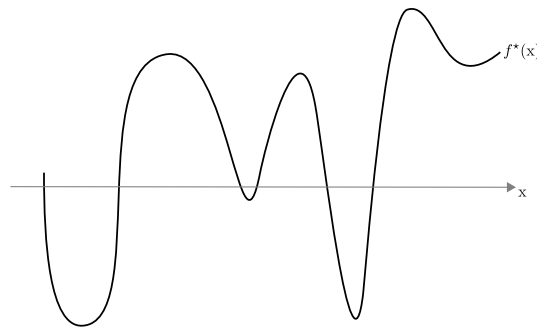


Figure 3.2: Voice signal corresponding to the command “turn off the t.v.”.

3.2 Intuition

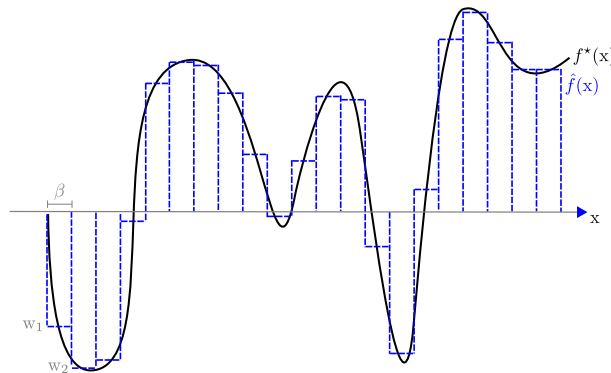
To build some intuition, suppose we want to estimate the following function f^* :

Figure 3.3: Function $f^*(\mathbf{x})$ that we want to estimate.

One option is to use a Riemann-type approximation $\hat{f} = \sum_{\ell=1}^L g_{\ell}$ where g_{ℓ} are step functions of the form:

$$g_{\ell}(x) = \begin{cases} w_{\ell} & \text{if } \beta(\ell-1) \leq x < \beta\ell \\ 0 & \text{otherwise.} \end{cases}$$

Here w_{ℓ} indicates the height of f^* in the ℓ^{th} interval $[\beta(\ell-1), \beta\ell)$. This type of approximation would look something like:

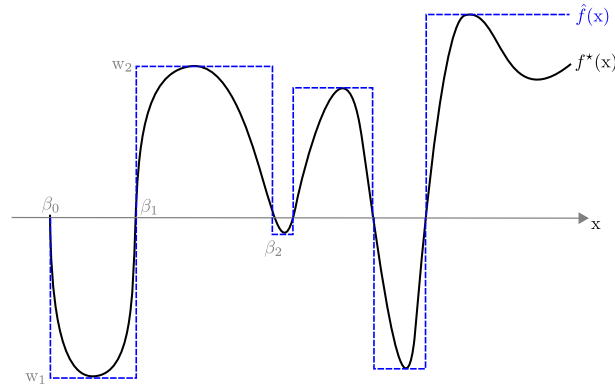


Notice that the parameters w_1, \dots, w_ℓ , together with the interval width β determine \hat{f} . Unfortunately, estimating this type of function would require to observe *at least* one sample x on each interval $[\beta(\ell-1), \beta\ell)$. For the figure above, this would require at least 21 samples. This may not sound like much. However, as we move to higher dimensions, this number scales exponentially, and quickly becomes prohibitively large: if f^* is a multivariate function of $\mathbf{x} \in \mathbb{R}^d$ (instead of a single variable function of $x \in \mathbb{R}$), then a similar approximation would require 21^d samples. With d as little as 10, this would mean more than 1 trillion samples. As we will see, in most of the big data applications that we are interested in, d is often in the order of hundreds, and often thousands or millions, whence this approach is completely infeasible.

Another option is to allow varying size intervals, such that

$$g_\ell(x) = \begin{cases} w_\ell & \text{if } \beta_{\ell-1} \leq x < \beta_\ell \\ 0 & \text{otherwise.} \end{cases}$$

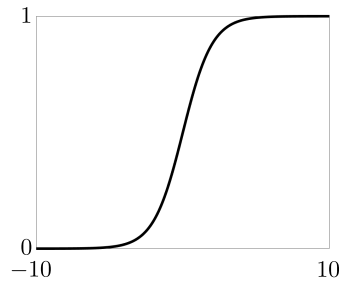
In this case, the parameters are the weights w_1, \dots, w_L and the interval boundaries β_0, \dots, β_L . This type of approximation would look more like:



and will potentially require far fewer samples. The downside is that now \hat{f} may be quite inaccurate (on top of being discontinuous). To address this, we can make each g_ℓ a simple nonlinear yet continuous function. For example, a sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}, \quad (3.1)$$

which looks as follows:

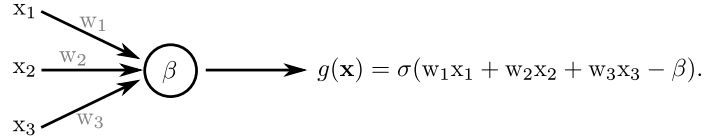


Notice that $\sigma(wx - \beta)$ *shifts* the function by β , and *squeezes* the function by w . So if we let $g(x) = \sigma(wx - \beta)$, and we choose the right parameters w and β , and add a bunch of these functions, then we can approximate

f^* much more accurately. Moreover, if f^* is a function of more than one variable, say of a vector $\mathbf{x} \in \mathbb{R}^d$, then we can generalize g in a very natural way:

$$g(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} - \beta), \quad (3.2)$$

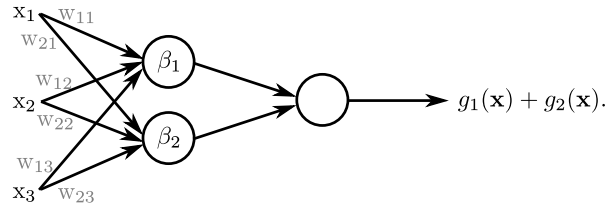
where now $\mathbf{w} \in \mathbb{R}^d$ is a parameter vector. Just as the step function is the building block of the Riemann-type approximation above, this type of function g is precisely the building block (often called *neuron*) of a neural network, usually depicted as follows:



The main insight of a neural network is that by adding a bunch of these building blocks, one can approximate f^* with arbitrary accuracy. For example, with two neurons we would obtain the following:

$$\hat{f}(\mathbf{x}) = g_1(\mathbf{x}) + g_2(\mathbf{x}),$$

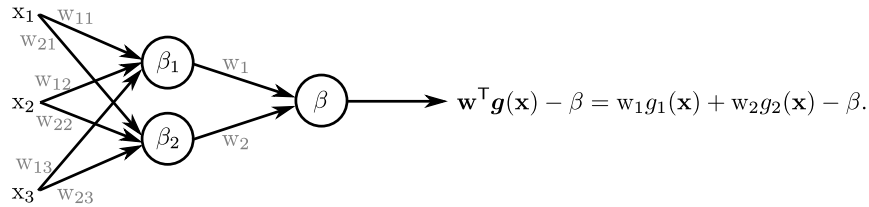
which we can depict as follows:



There is no reason why we cannot include weights to each g , nor a final shift, to obtain:

$$\hat{f}(\mathbf{x}) = w_1 g_1(\mathbf{x}) + w_2 g_2(\mathbf{x}) - \beta = \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \end{bmatrix} - \beta =: \mathbf{w}^\top \mathbf{g}(\mathbf{x}) - \beta,$$

where w_1 , w_2 , and β are new parameters. This function would be represented like:



There is no reason to stop there. We can add more neurons, and more *layers* to obtain more powerful networks, capable of approximating more complex functions. Let L be the number of layers, and let n_ℓ denote the number of neurons in the ℓ^{th} layer. For $\ell = 2, 3, \dots, L$, let $\mathbf{W}^\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ be the matrix formed by transposing and stacking the weight vectors of the n_ℓ neurons in the ℓ^{th} layer, and let $\boldsymbol{\beta}^\ell \in \mathbb{R}^{n_\ell}$ be the vector containing the shifting coefficients (often called *bias* coefficients) of the neurons in the ℓ^{th} layer, such that the output at the second layer is given by:

$$\mathbf{g}^2(\mathbf{x}) = \sigma(\mathbf{W}^2 \mathbf{x} - \boldsymbol{\beta}_2),$$

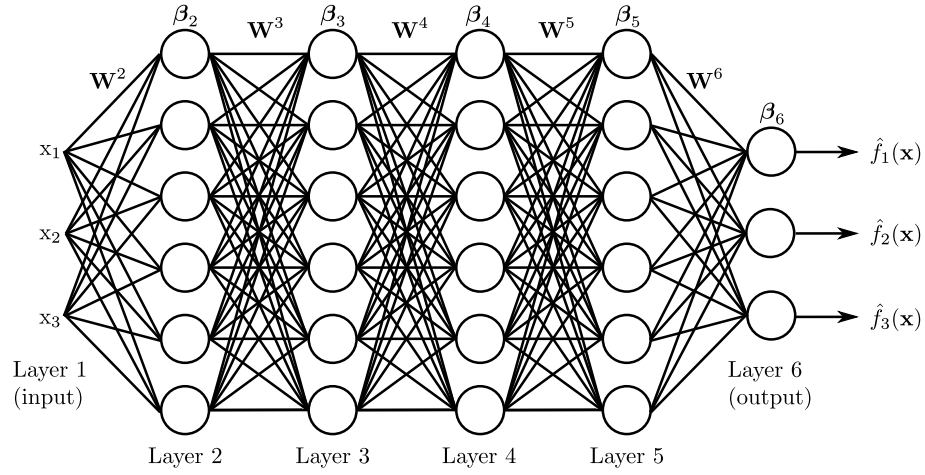


Figure 3.4: Components of a neural network.

and for $\ell = 3, \dots, L$, the output at the ℓ^{th} layer is given by:

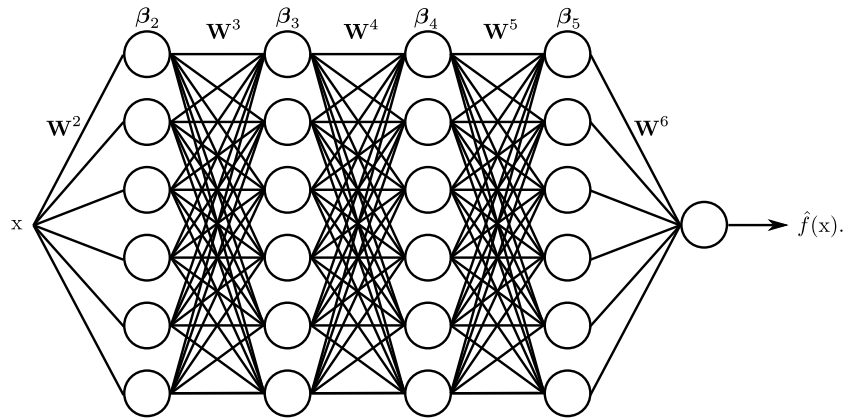
$$\mathbf{g}^\ell(\mathbf{x}) = \sigma(\mathbf{W}^\ell \mathbf{g}^{\ell-1}(\mathbf{x}) - \beta_\ell),$$

where $\mathbf{g}^L(\mathbf{x})$ is the final output of the network, also denoted as:

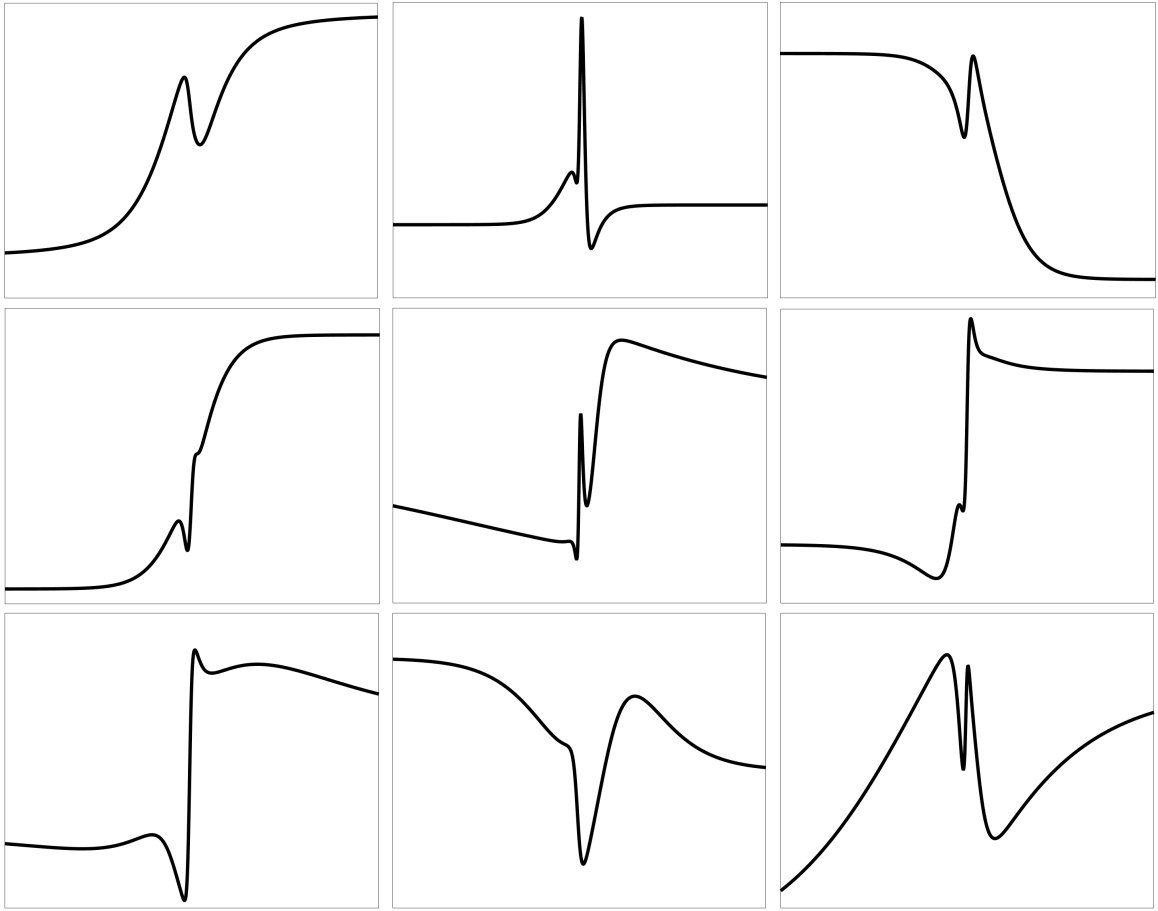
$$\hat{\mathbf{f}}(\mathbf{x}) := \sigma(\mathbf{W}^L \sigma(\mathbf{W}^{L-1} \dots \sigma(\mathbf{W}^3 \sigma(\mathbf{W}^2 \mathbf{x} - \beta^2) - \beta^3) \dots - \beta^{L-1}) - \beta^L). \quad (3.3)$$

Notice that $\hat{\mathbf{f}}(\mathbf{x}) \in \mathbb{R}^{n_L}$ may be a vector (if $n_L > 1$), as opposed to a scalar (if $n_L = 1$), and so neural networks allow to infer vector functions. See Figure 3.4 to build some intuition.

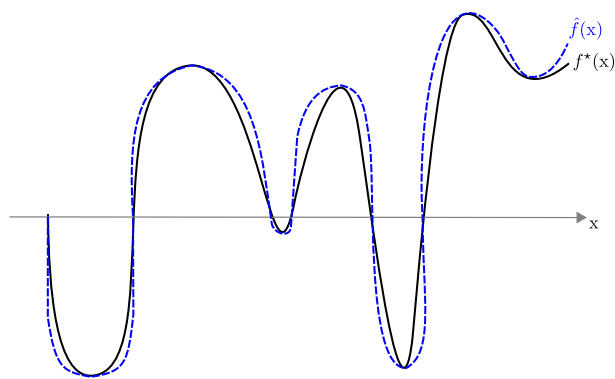
Going back to estimating the function in Figure 3.3, we can construct the following network:



Depending on the parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$ that we choose, may produce functions like the following:



An approximation of the function in Figure 3.3 with a neural network like this would look like:



However, this will require to identify the right parameters $\{\mathbf{W}^\ell, \boldsymbol{\beta}^\ell\}_{\ell=2}^L$ that produce this approximation.

3.3 Identifying the Right Parameters

As discussed before, with the right parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$, the function $\hat{\mathbf{f}}(\mathbf{x})$ in (3.3) can approximate any function $\mathbf{f}^*(\mathbf{x})$ with arbitrary accuracy. The challenge is to find the right parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$.

Hence, our goal can be summarized as follows: find the parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$ of the function $\hat{\mathbf{f}}$ in (3.3), such that $\hat{\mathbf{f}}(\mathbf{x}) \approx \mathbf{f}^*(\mathbf{x})$. Recall that we do not know \mathbf{f}^* . However, we do have *training* data. More precisely, we have N samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ of N , as well as their *response* $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$. For example, if we were studying diabetes, \mathbf{x}_i could contain demographic information about the i^{th} person in our study, and \mathbf{y}_i could be a binary variable indicating whether this person is diabetic or not.

In other words, we know that $\mathbf{f}^*(\mathbf{x}_i) = \mathbf{y}_i$ for every $i = 1, 2, \dots, N$. Intuitively, this means that we get to see N snapshots of \mathbf{f}^* at points \mathbf{x}_i . Our goal is to exploit this information to find the parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$ such that $\hat{\mathbf{f}} \approx \mathbf{f}^*$, such that the network can reproduce the response \mathbf{y} whenever a new vector \mathbf{x} is fed to the network. This can be done by minimizing the error (over all training data) between the network's prediction $\hat{\mathbf{f}}(\mathbf{x}_i)$ and its corresponding observation \mathbf{y}_i . Mathematically, we can achieve this by solving the following optimization

$$\min_{\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{f}}(\mathbf{x}_i)\|^2. \quad (3.4)$$

The most widely used technique to solve (3.4) is through stochastic gradient descent and back-propagation. The key insight is to observe that the gradient of each weight matrix can be computed *backwards* in terms of the gradients of the weight matrices of subsequent layers. To see this, first define $\mathbf{z}_i^1 = \mathbf{y}_i^1 = \mathbf{x}_i$, and then for $\ell = 2, 3, \dots, L$, recursively define $\mathbf{z}_i^\ell := \mathbf{W}^\ell \mathbf{y}_i^{\ell-1} - \beta^\ell$, where $\mathbf{y}_i^\ell := \mathbf{g}^\ell(\mathbf{x}_i) = \sigma(\mathbf{z}_i^\ell)$ denotes the output at the ℓ^{th} layer, so that $\hat{\mathbf{f}}(\mathbf{x}_i) = \mathbf{y}_i^L$. Defining the *cost* of the i^{th} sample as $\mathbf{c}_i := \mathbf{y}_i - \hat{\mathbf{f}}(\mathbf{x}_i)$, and

$$\begin{aligned} \delta_i^L &:= \mathbf{c}_i \odot \mathbf{f}'(\mathbf{z}_i^L), \\ \delta_i^\ell &:= [(\mathbf{W}^{\ell+1})^\top \delta_i^{\ell+1}] \odot \mathbf{f}'(\mathbf{z}_i^\ell), \end{aligned} \quad 2 \leq \ell \leq L-1,$$

where \odot represents the Hadamard product, with a simple chain rule we obtain the following gradients:

$$\nabla_i \mathbf{W}^\ell := \frac{\partial \|\mathbf{c}_i\|^2}{\partial \mathbf{W}^\ell} = -2 \delta_i^\ell (\mathbf{y}_i^{\ell-1})^\top, \quad \nabla_i \beta^\ell := \frac{\partial \|\mathbf{c}_i\|^2}{\partial \beta^\ell} = -2 \delta_i^\ell, \quad 2 \leq \ell \leq L. \quad (3.5)$$

Training the neural network is equivalent to identifying the parameters $\{\mathbf{W}^\ell, \beta^\ell\}_{\ell=2}^L$ such that $\hat{\mathbf{f}} \approx \mathbf{f}^*$. This can be done using stochastic gradient descent, which iteratively updates the parameters of interest according to (3.5).

More precisely, at each training time t we select a random subsample $\Omega_t \subset \{1, 2, \dots, N\}$ of the training data (hence the term *stochastic*), and we update our parameters as follows

$$\begin{aligned} \mathbf{W}_t^\ell &= \mathbf{W}_{t-1}^\ell - \eta \sum_{i \in \Omega_t} \nabla_i \mathbf{W}_{t-1}^\ell, \\ \beta_t^\ell &= \beta_{t-1}^\ell - \eta \sum_{i \in \Omega_t} \nabla_i \beta_{t-1}^\ell, \end{aligned}$$

where η is the gradient *step* parameter to be tuned. This iteration is repeated until convergence to obtain the *final* parameters $\{\hat{\mathbf{W}}^\ell, \hat{\beta}^\ell\}_{\ell=2}^L$.

3.4 Neural Networks Flavors

At this point you may be wondering whether $\mathbf{g}^\ell(\mathbf{x}) = \sigma(\mathbf{W}^\ell \mathbf{g}^{\ell-1}(\mathbf{x}) - \boldsymbol{\beta})$ is the only option for \mathbf{g} , where σ is the sigmoid function in (3.1). The answer is: no. For instance, if you use

$$\sigma(x) = x,$$

then you obtain a *linear* layer/network. If you use

$$\sigma(x) = \max(0, x),$$

then you obtain the so-called *rectified linear unit* (ReLU) layer/network. σ is usually called *activation function*. Similarly, if instead of $\mathbf{g}^\ell(\mathbf{x}) = \sigma(\mathbf{W}^\ell \mathbf{g}^{\ell-1}(\mathbf{x}) - \boldsymbol{\beta})$ you use

$$\mathbf{g}^\ell(\mathbf{x}) = \sigma(\mathbf{W}^\ell * \mathbf{g}^{\ell-1}(\mathbf{x}) - \boldsymbol{\beta}),$$

where $*$ denotes the convolution operator, then you obtain a *convolutional* layer/network.

The *topology* (shape) of the network is defined by the number of layers (L) and the number of neurons in each layer (n_1, n_2, \dots, n_L). Networks with large L are often called *deep*, as opposed to *shallow* networks that have small L .

Depending on the problem at hand, one choice of L , n_1, \dots, n_L , σ , and \mathbf{g} (or combinations) may be better than another. For example, for image classification, people usually like to mix cascades of one linear layer, followed by a ReLU layer, followed by a convolutional layer, with large n_ℓ in each case, and sometimes with fixed (predetermined) \mathbf{W} 's in the convolutional layers. The possibilities are endless, and deciding on the best choice of topology remains more of an art than science.

3.5 A Word of Warning

We have mentioned before that using neural networks we can approximate any function \mathbf{f}^* with arbitrary accuracy. Put another way, for every \mathbf{f}^* , and every $\epsilon > 0$ there exists a function $\hat{\mathbf{f}}_\theta$ with parameters $\theta := \{\mathbf{W}^\ell, \boldsymbol{\beta}^\ell\}_{\ell=2}^L$ such that,

$$\|\mathbf{f}^*(\mathbf{x}) - \hat{\mathbf{f}}_\theta(\mathbf{x})\|_2 < \epsilon \quad \text{for every } \mathbf{x} \text{ in the domain of } \mathbf{f}^*. \quad (3.6)$$

The challenge is to find the parameter θ for which (3.6) is true, which we aim to do using gradient descent to minimize the following *cost* function:

$$\mathbf{c}(\theta) := \sum_{i=1}^N \|\mathbf{f}^*(\mathbf{x}_i) - \hat{\mathbf{f}}_\theta(\mathbf{x}_i)\|_2.$$

The wrinkle is that the cost function that we are trying to minimize may be non-convex, which implies that we may *never* find the right parameter θ . In other words, for every \mathbf{f}^* there will always exist *a* neural network (parametrized by θ) that approximates \mathbf{f}^* arbitrarily well. However, we may be unable to find such network.