

## Topic 12: Neural Networks

**DO NOT POLLUTE!** AVOID PRINTING, OR PRINT 2-SIDED MULTIPAGE.

## 12.1 Introduction

Neural networks are one of the main reasons why Machine Learning is making a big splash. Similar to logistic regression and support vector machines (SVMs), their main purpose is prediction/classification. Classical applications of neural networks include:

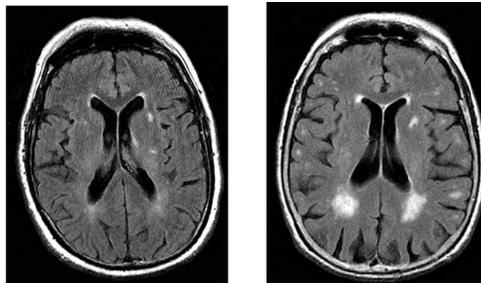
1. **Image classification.** For example, digit classification (determining the digit contained in an image), or medical diagnostics (deciding whether a patient is healthy or not based on an MRI).
2. **Natural language processing.** For example, interpreting voice commands, like Siri and Alexa do.
3. **Stock market prediction.**

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

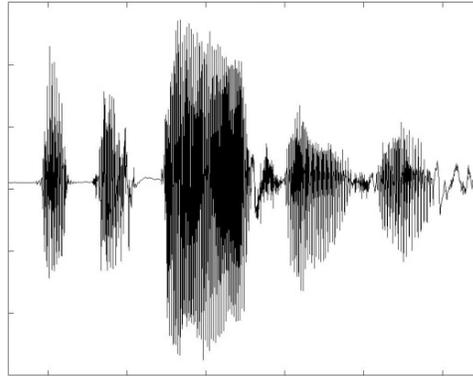
1. **Image classification.** In digit classification  $\mathbf{x}$  is the vector containing the pixels in the image of a digit, and  $\mathbf{y} \in \{0, \dots, 9\}$  is the digit depicted in such image:

5 0 4 1 9

In medical diagnostics,  $\mathbf{x}$  could contain the pixels in an MRI, and  $\mathbf{y} \in \{0, 1\}$  would be the diagnostic, 0 corresponding to healthy, and 1 to Alzheimers.

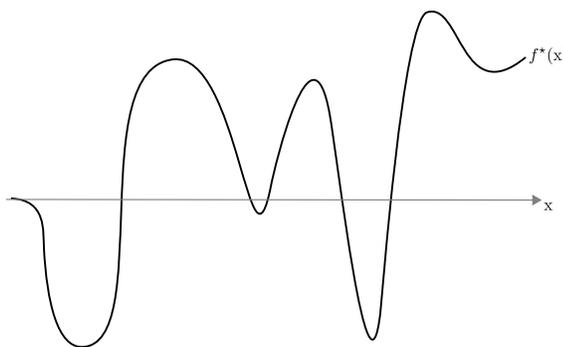


**Natural language processing.** Here  $\mathbf{x}$  is a vector containing the values in a voice signal, and  $\mathbf{y}$  is the interpretation, in this case “turn off the t.v.”:



**Stock market prediction.** Here  $\mathbf{x}$  could be a sequence of stock market prices at times  $t = 1, \dots, T$ , and  $\mathbf{y}$  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. Consequently,  $\mathbf{f}^*$  could be quite complex.



Neural networks arise as an alternative to approximate complex functions like these. The main idea behind neural networks is to use a sequence of simpler functions that interact with one another in a networked way, so that combined, they approximate  $\mathbf{f}^*$  with arbitrary precision.

## 12.2 The Building Blocks of Neural Networks

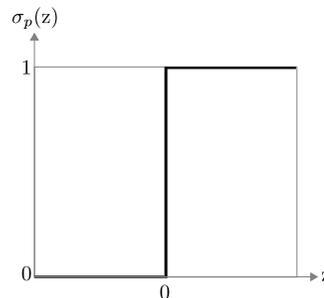
The *perceptron* was the first building block used in neural networks, inspired by the behavior of neurons in the brain. This function would have the form

$$g_p(\mathbf{x}) = \sigma_p(\boldsymbol{\theta}^T \mathbf{x} - b) \quad (12.1)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^D$  and  $b \in \mathbb{R}$  were the parameters of the perceptron, and the perceptron's *activation function*  $\sigma_p$  is given by

$$\sigma_p(z) = \begin{cases} 1 & \text{if } z > 0, \\ 0 & \text{otherwise.} \end{cases}$$

If you recall from SVMs, this is essentially the function that defines a hyperplane. The main intuition here is that a perceptron would *fire up*, like a neuron, if it receives the right input.



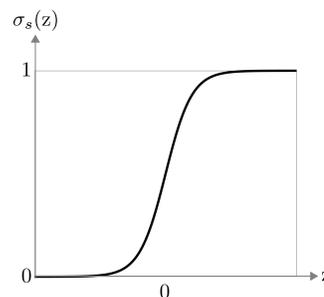
Unfortunately, such a discontinuous activation function can be quite unstable, as a tiny change in the input can produce a massive change in the output, which is a behavior we often don't want. For example, imagine this type of function were used to decide whether a self-driving car should stop or move. It sounds a bit unreasonable that if  $\boldsymbol{\theta}^T \mathbf{x} - b = -0.0001$ , then the car stops, and if  $\boldsymbol{\theta}^T \mathbf{x} - b = 0.0001$ , then the car moves. Often we would like a softer transition, which is one of the main reasons people decided to use the *sigmoid neuron*

$$g_s(\mathbf{x}) = \sigma_s(\boldsymbol{\theta}^T \mathbf{x} - b)$$

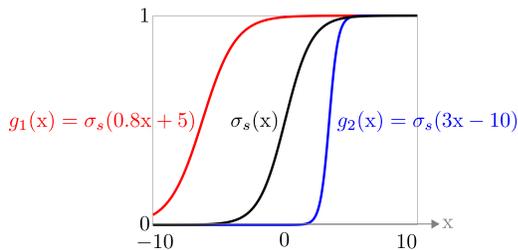
which has the same form as (12.1), only instead of  $\sigma_p$ , it uses the *logistic sigmoid* activation function

$$\sigma_s(z) := \frac{1}{1 + e^{-z}},$$

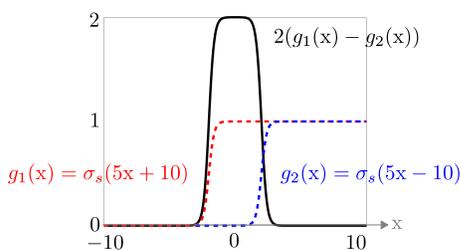
which looks like this:



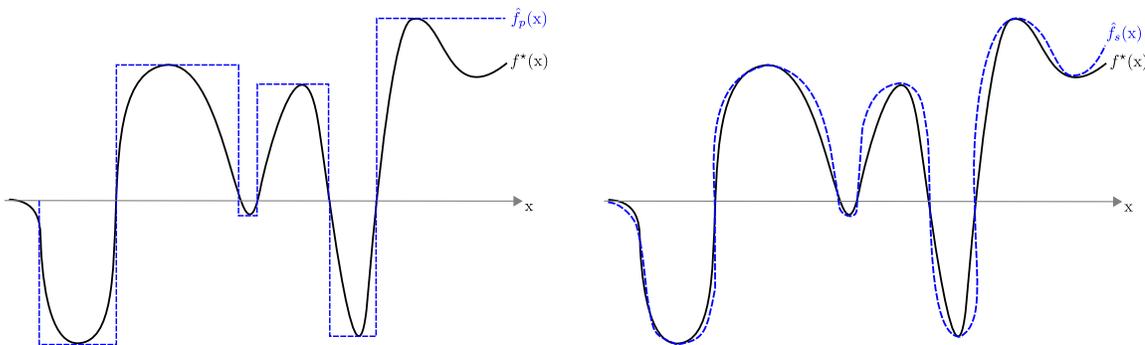
Notice that  $g(\mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x} - b)$  *shifts* and *squeezes* the function  $\sigma(x)$  by  $b$  and  $\boldsymbol{\theta}$ , for example:



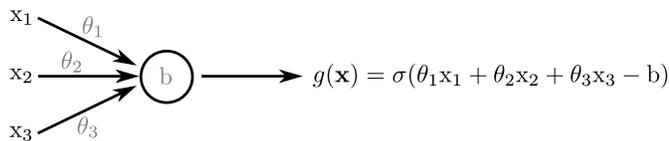
So by adding and scaling bunch of these functions  $g$ , each with appropriate parameters  $\theta$  and  $b$ , we can create *soft step* functions, for example:



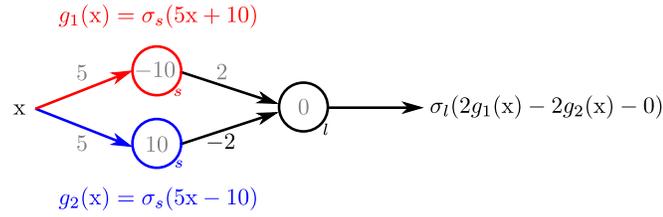
Hence, by adding and scaling enough of these functions  $g$  (either perceptrons  $g_p$  or sigmoid neurons  $g_s$ ), we can approximate any function  $f^*$  to arbitrary precision, kind of like a Riemann-type approximation:



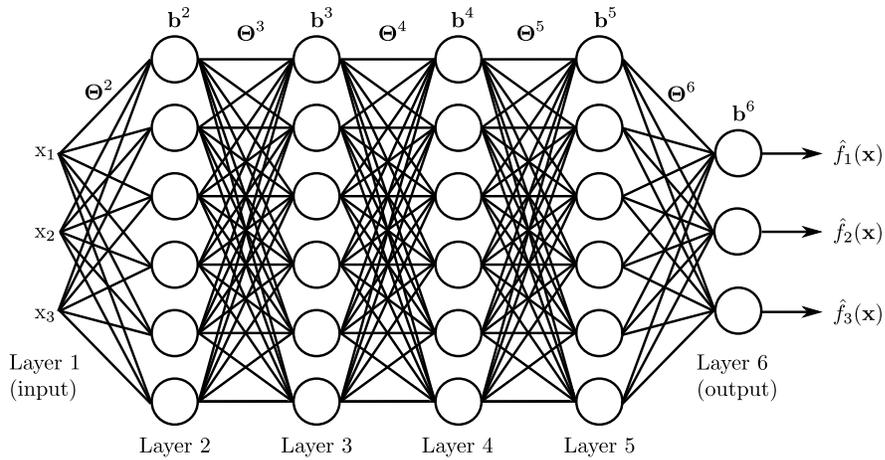
The composition of functions  $g$  is often represented through a network, where each function  $g$  is depicted by a node (a.k.a. neuron) like this:



Notice that the activation function  $\sigma$  does not have to be restricted to  $\sigma_p$  or  $\sigma_s$ . One can use other activation functions, such as the *linear* activation function  $\sigma_l(z) := z$ , or the so-called *rectified linear unit* (ReLU)  $\sigma_r(z) := \max(0, z)$ , or the *tanh*  $\sigma_h(z) := 2\sigma_s(2z) - 1$ , or combinations. For example, a representation of the soft step function above would look like:



where the neurons on the left have sigmoid activation functions, and the neuron on the right has a linear activation function. More generally, one can add more neurons to obtain more powerful networks, capable of approximating more complex functions (the downside of larger networks is that they will have more parameters that will have to be learnt, and one runs the risk of overfitting). In general, we consider networks like the following:



Here we use  $L$  to denote the number of layers (in the figure above  $L = 6$ ), and  $n_\ell$  to denote the number of neurons in the  $\ell^{\text{th}}$  layer. For  $\ell = 2, 3, \dots, L$ ,  $\Theta^\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$  is the matrix formed by transposing and stacking the weight vectors of the  $n_\ell$  neurons in the  $\ell^{\text{th}}$  layer;  $\mathbf{b}^\ell \in \mathbb{R}^{n_\ell}$  is the vector containing the shifting coefficients (often called *bias* coefficients) of the neurons in the  $\ell^{\text{th}}$  layer; and  $\sigma_\ell$  is the activation function of the  $\ell^{\text{th}}$  layer. This way the output at the second layer is given by:

$$\mathbf{g}^2(\mathbf{x}) = \sigma_2(\Theta^2 \mathbf{x} - \mathbf{b}^2). \tag{12.2}$$

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

$$\mathbf{g}^\ell(\mathbf{x}) = \sigma_\ell(\Theta^\ell \mathbf{g}^{\ell-1}(\mathbf{x}) - \mathbf{b}^\ell), \tag{12.3}$$

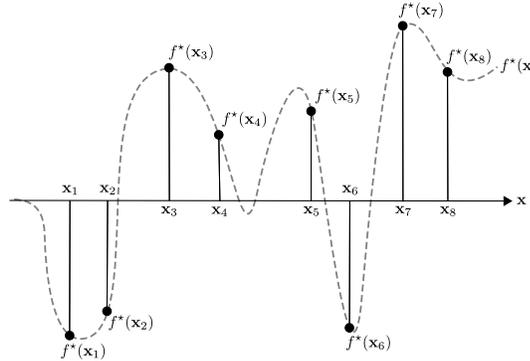
and the final output of the network,  $\mathbf{g}^L(\mathbf{x})$ , is also denoted as:

$$\hat{\mathbf{f}}(\mathbf{x}) := \sigma_L(\underbrace{\Theta^L \sigma_{L-1}(\underbrace{\Theta^{L-1} \dots \sigma_3(\underbrace{\Theta^3 \sigma_2(\Theta^2 \mathbf{x} - \mathbf{b}^2) - \mathbf{b}^3} - \mathbf{b}^{L-1}) - \mathbf{b}^L)}_{\mathbf{g}^L(\mathbf{x})}) \tag{12.4}$$

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.

### 12.3 Learning $\hat{f}$

By the arguments above, with the right parameters  $\Theta := \{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$ , the function  $\hat{f}$  in (12.4) can approximate any function  $f^*$  with arbitrary accuracy. The challenge is to find the right parameters  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$ . The process of finding such parameters is often known as *learning*. To this end, we use *training* data, that is, samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ , 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  $i^{\text{th}}$  person in our study, and  $\mathbf{y}_i$  could be a binary variable indicating whether this person is diabetic or not. Since  $\mathbf{y}_i = f^*(\mathbf{x}_i)$  for every  $i = 1, 2, \dots, N$ , this means that we don't get to see *all* of  $f^*$ , but we get to see  $N$  *snapshots* of  $f^*$  at the samples  $\mathbf{x}_i$ :



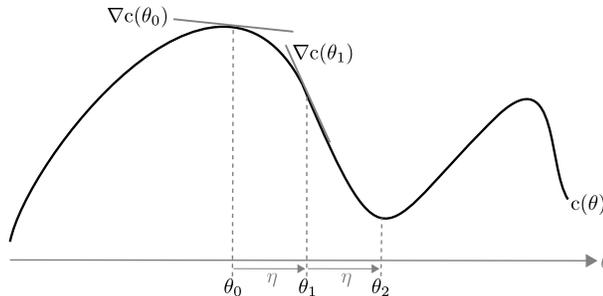
Our goal is to exploit this information to find the parameters  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$  such that  $\hat{f} \approx f^*$ , such that the function (network) can reproduce the response  $\mathbf{y}$  whenever a new vector  $\mathbf{x}$  is fed to the function (network). This can be done by minimizing the error over all training data (often called *cost function*) between the network's prediction  $\hat{f}(\mathbf{x}_i)$  and its corresponding observation  $\mathbf{y}_i$ . Mathematically, we can achieve this by solving the following optimization

$$\min_{\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L} \sum_{i=1}^N \|y_i - \hat{f}(\mathbf{x}_i)\|^2. \tag{12.5}$$

Notice that the dependency of  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$  is hidden in  $\hat{f}$  (see (12.4)).

### 12.4 Backpropagation

The most widely used technique to minimize the *cost*  $c$  in (12.5) is through *stochastic gradient descent* and *backpropagation*. Recall that gradient descent iteratively moves a step of size  $\eta$  in the negative direction of the gradient:



This, of course, requires computing the gradient of each parameter. The key insight to achieve this is that the gradient of the parameters at the  $\ell^{\text{th}}$  layer can be computed *backwards* in terms of the gradients of the parameters 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  $\mathbf{z}_i^\ell := \Theta^\ell \mathbf{y}_i^{\ell-1} - \mathbf{b}^\ell$ , where  $\mathbf{y}_i^\ell := \mathbf{g}^\ell(\mathbf{x}_i) = \sigma_\ell(\mathbf{z}_i^\ell)$  denotes the output at the  $\ell^{\text{th}}$  layer, so that  $\hat{\mathbf{f}}(\mathbf{x}_i) = \mathbf{y}_i^L$ . Define the *cost* of the  $i^{\text{th}}$  sample as  $\mathbf{c}_i := \mathbf{y}_i - \hat{\mathbf{f}}(\mathbf{x}_i)$ , and

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

where  $\odot$  represents the Hadamard product, and  $\sigma'_\ell$  represents the derivative of  $\sigma_\ell$ . Then with a simple chain rule we obtain the following gradients:

$$\nabla_i \Theta^\ell := \frac{\partial \|\mathbf{c}_i\|^2}{\partial \Theta^\ell} = -2 \delta_i^\ell (\mathbf{y}_i^{\ell-1})^\top, \quad \nabla_i \mathbf{b}^\ell := \frac{\partial \|\mathbf{c}_i\|^2}{\partial \mathbf{b}^\ell} = -2 \delta_i^\ell, \quad 2 \leq \ell \leq L. \quad (12.6)$$

With these results, one can use gradient descent, or, to reduce computational burden an increase convergence speed, *stochastic* gradient descent, which at each training time  $t$  selects a random subsample  $\Omega_t \subset \{1, 2, \dots, N\}$  of the training data (hence the term *stochastic*), and updates the parameters according to this subsample:

$$\begin{aligned}\Theta_t^\ell &= \Theta_{t-1}^\ell - \eta \sum_{i \in \Omega_t} \nabla_i \Theta_{t-1}^\ell, \\ \mathbf{b}_t^\ell &= \mathbf{b}_{t-1}^\ell - \eta \sum_{i \in \Omega_t} \nabla_i \mathbf{b}_{t-1}^\ell.\end{aligned}$$

These iterations are repeated until convergence to obtain the *final* parameters  $\{\hat{\Theta}^\ell, \hat{\mathbf{b}}^\ell\}_{\ell=2}^L$ .

## 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  $\hat{\mathbf{f}}$  with parameters  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$  such that,

$$\|\mathbf{f}^*(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x})\|_2 < \epsilon \quad \text{for every } \mathbf{x} \text{ in the domain of } \mathbf{f}^*. \quad (12.7)$$

The challenge is to find the parameters  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$  for which (12.7) is true. We aim to find such parameters by minimizing the cost function in (12.5). The wrinkle is that such cost may be non-convex, which implies that we may *never* find the right parameters  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$ . In other words, for every  $\mathbf{f}^*$  there will always exist a neural network (parametrized by  $\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L$ ) that approximates  $\mathbf{f}^*$  arbitrarily well. However, we may be unable to find such network.

## 12.5 Neural Networks Flavors

In practice, a specific network may work better than another for each problem. Unfortunately, options (number of layers, number of neurons in each layer, type of layers, connectivity, etc.) are endless, and deciding on the best choice remains more of an art than science. However, people have empirically identified certain strategies that tend to work well for certain tasks.

Starting from encoding, there are several creative and popular strategies. For example:

- **One-hot-encoding.** Nominal features are represented with binary vectors with a single 1. For example, for nucleotides, we could have:

$$A \mapsto \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad C \mapsto \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad G \mapsto \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad T \mapsto \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

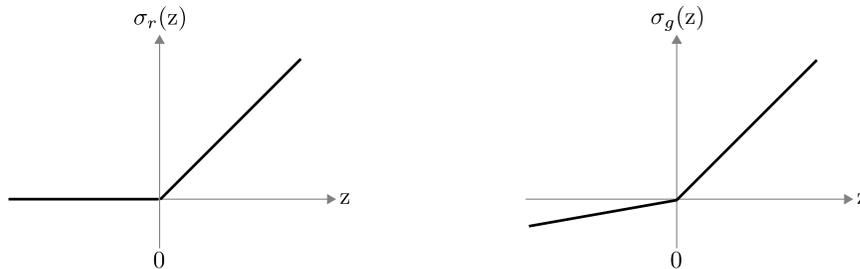
- **Thermometer encoding.** Ordinal features are represented with binary vectors that have an increasing number of ones. For example, for size we could have:

$$small \mapsto \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad medium \mapsto \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad large \mapsto \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

Often, the performance of a network depends on choosing the right feature encoding. Similarly, performance heavily depends on the choice of activation functions for each layer. Popular choices for *intermediate* layers (that is, all but the first and last, a.k.a. *hidden* layers) include combinations of the aforementioned linear, ReLU, sigmoid, tanh activation functions, and variants of these, for example the *generalized ReLU* with parameter  $\alpha$ :

$$\sigma_g(z) := \max(0, z) + \alpha \min(0, z).$$

In particular, when  $\alpha = 0$  we recover the standard ReLU, and when  $\alpha$  is small (e.g., 0.01), we obtain the so-called *Leaky ReLU*:



As to the output layer, arguably the most popular choice in classification problems is the *softmax* function:

$$\sigma_{sm}(\mathbf{z}) = \frac{e^{\mathbf{z}}}{\sum_j e^{z_j}}$$

The intuition here is that the softmax function produces a vector with nonnegative entries that add up to 1, and can thus be interpreted as the probability that the given input belongs to each class. Possibilities do not end with the choice of activation functions. One can also play around with the form of  $\mathbf{g}$ . For example, rather than the linear product in (12.2) and (12.3), one can use a convolution:

$$\begin{aligned} \mathbf{g}^2(\mathbf{x}) &= \sigma_2(\Theta^2 * \mathbf{x} - \mathbf{b}^2), \\ \mathbf{g}^\ell(\mathbf{x}) &= \sigma_\ell(\Theta^\ell * \mathbf{g}^{\ell-1}(\mathbf{x}) - \mathbf{b}^\ell), \end{aligned} \quad 3 \leq \ell \leq L,$$

resulting in the popular *convolution layers*. In addition, one can also play around with the cost function that one aims to optimize. For example, instead of minimizing the squared error in (12.5), one could try to maximize a data likelihood model:

$$\max_{\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L} \mathbb{P}(\mathbf{Y}, \mathbf{X} \mid \{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L),$$

or even a posterior probability:

$$\max_{\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L} \mathbb{P}\left(\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L \mid \mathbf{Y}, \mathbf{X}\right),$$

which (by Bayes rule) is equivalent to maximizing the log-likelihood with a prior regularization term:

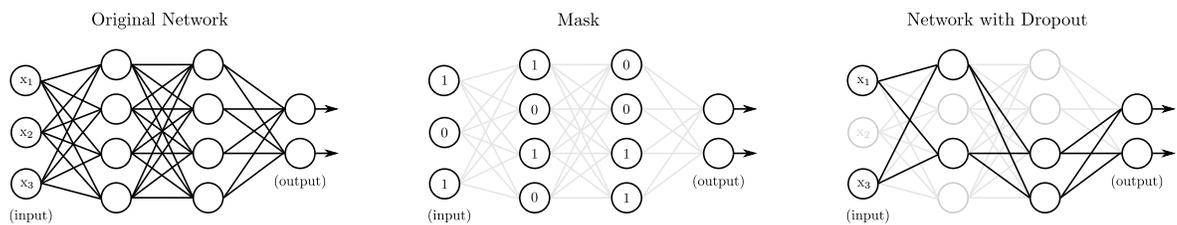
$$\max_{\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L} \log \mathbb{P}\left(\mathbf{Y}, \mathbf{X} \mid \{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L\right) + \log \mathbb{P}\left(\{\Theta^\ell, \mathbf{b}^\ell\}_{\ell=2}^L\right),$$

resulting in a bayesian framework. Similarly, one can play around with the cost function, adding regularization terms in the interest of favoring certain solutions, or achieving certain goals, like minimizing overfitting. For example, in addition to the squared error in (12.5), one could add a regularization term that favors sparse solutions:

$$\sum_{i=1}^N \|y_i - \hat{f}(\mathbf{x}_i)\|^2 + \lambda \sum_{\ell=2}^L \|\Theta^\ell\|_1.$$

where  $\lambda$  is a regularization parameter. Other alternatives to minimize overfitting include:

- **Early Stopping.** The main idea is to stop training as soon as the error on the validation data stops decreasing.
- **Dropout.** The main idea is to pretend at each update step, that a fraction of neurons don't exist (typically 20% of the input layer, and 50% of the intermediate layers). This is typically done by creating a random binary mask whose entries correspond to the input and intermediate neurons, then multiplying each neuron by its corresponding mask entry, and then proceed to do the update as usual.



**Data Augmentation.** The main idea is to generate additional samples to be used for training, by performing transformations on the original ones. For example, in image processing, typical approaches rotate, flip, translate, scale, crop, and add noise to the original samples:

