

# Mathematics of Image and Data Analysis

## Math 5467

### Neural Networks

Instructor: Jeff Calder

Email: [jcalder@umn.edu](mailto:jcalder@umn.edu)

<http://www-users.math.umn.edu/~jwcalder/5467>

## Last time

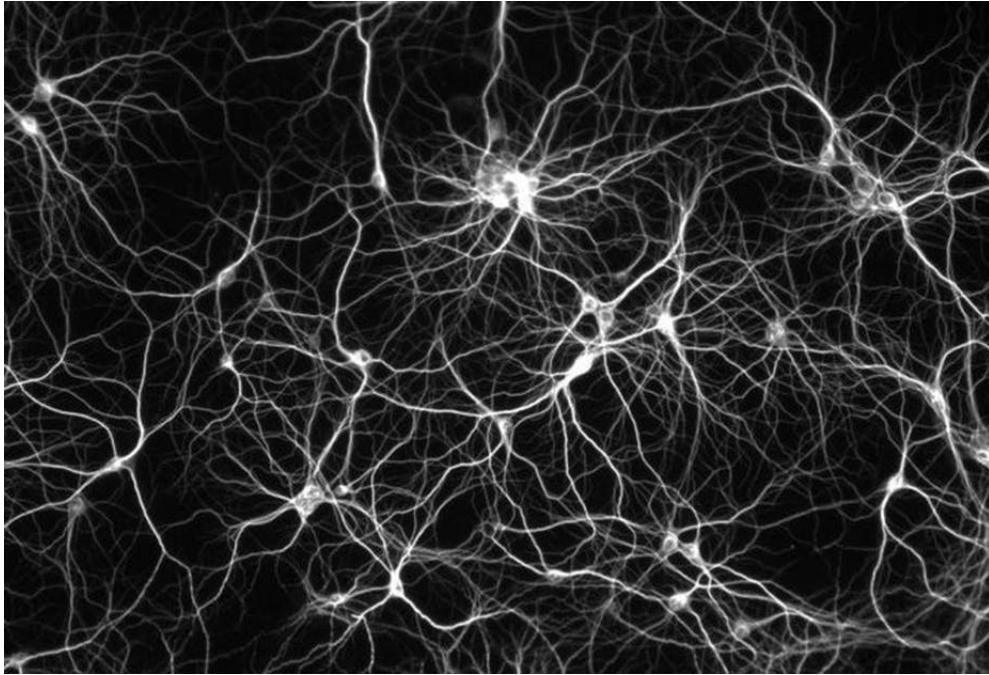
- Graph-based embeddings (spectral and t-SNE)

## Today

- Neural networks

# Artificial neural networks

Artificial neural networks (ANNs) are machine learning models (very) loosely built to model biological neural networks in the human brain.



# Artificial neural networks

(Artificial) neural networks are parameterized functions made up of simple building blocks: linear functions and simple nonlinearities. The basic building block is a neuron

$$f(x) = \sigma(\omega^T x + b),$$

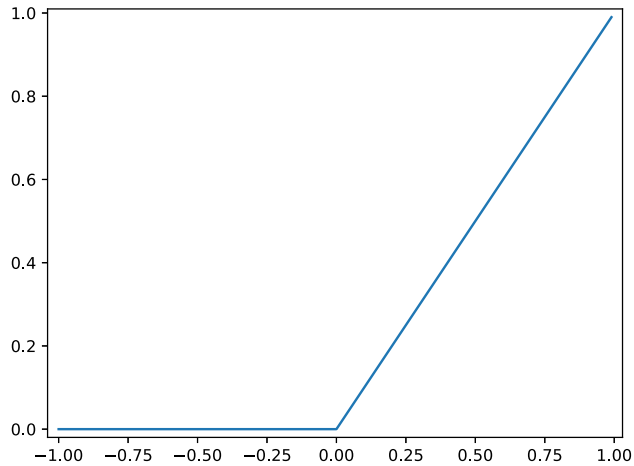
which is a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Common choices for the activation function  $\sigma$  is the rectified linear unit (ReLU)

$$(1) \quad \sigma(t) = \max\{t, 0\}.$$

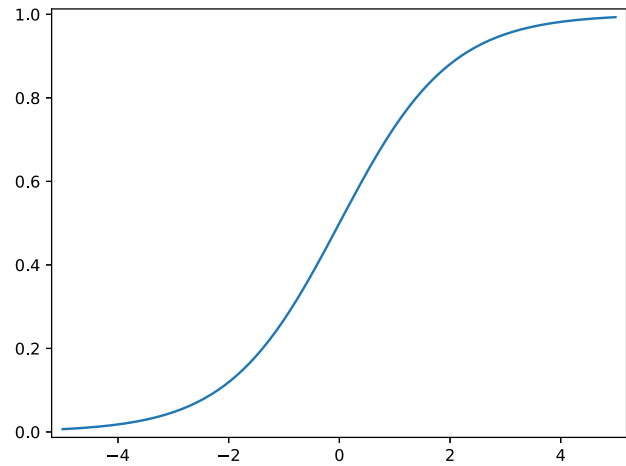
and the *sigmoid* activation function

$$(2) \quad \sigma(t) = \frac{1}{1 + e^{-t}}.$$

# Common activation functions



(a) ReLU



(b) Sigmoid

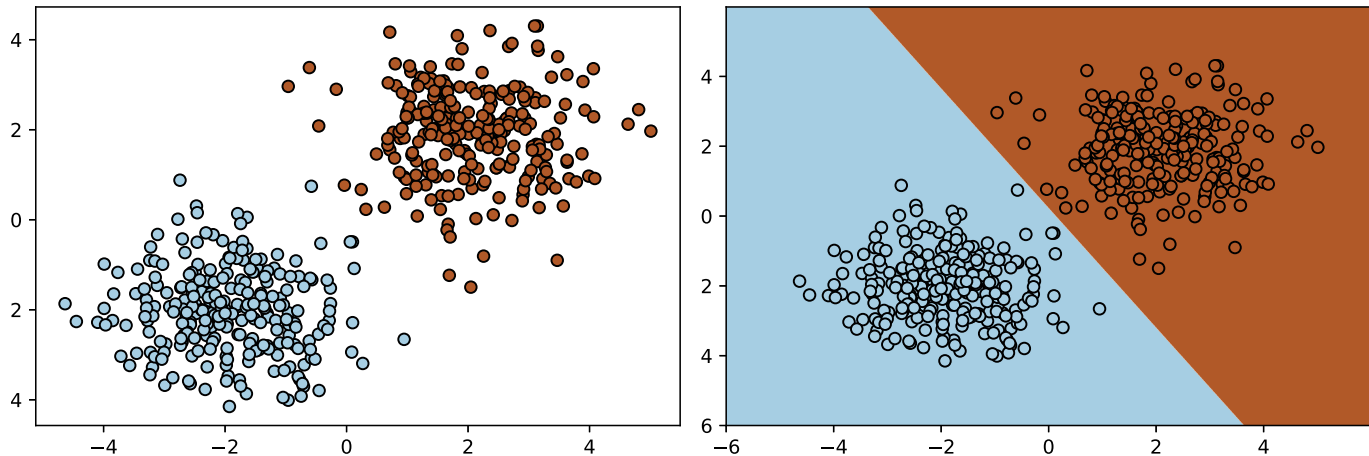
Figure 1: Plots of the ReLU and Sigmoid activation functions. Both activation functions have the behavior that they give zero, or close to zero, responses when the input is below a certain threshold, and give positive responses above.

# Connection to SVM

A single neuron

$$f(x) = \sigma(w^T x + b)$$

can implement a *support vector machine*. The set where the input  $w^T x + b$  is constant is a line in 2D or a hyperplane in  $\mathbb{R}^n$ .



# Non-linearly separable

Most real data is not linearly separable, and more than 1 neuron is required.

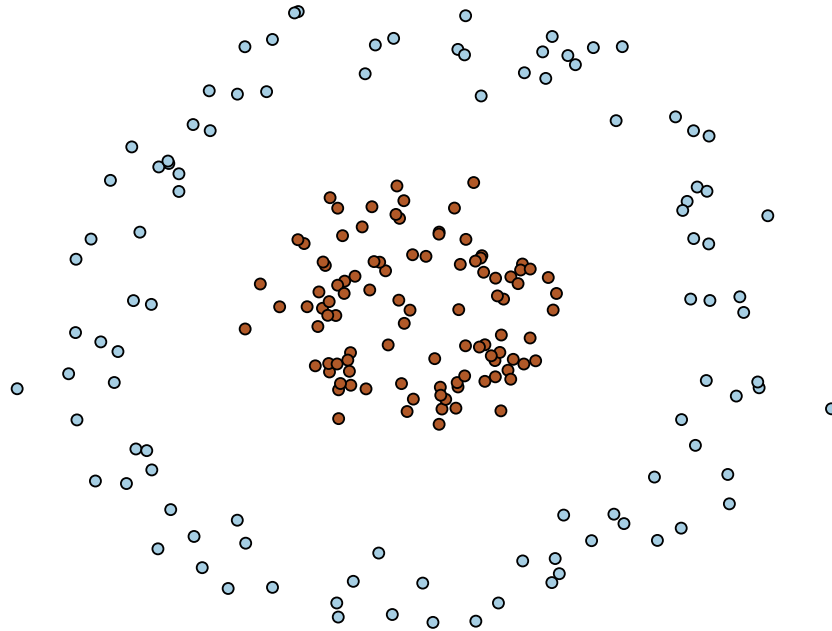


Figure 2: Synthetic data on rings consisting of two classes.

# Neural network

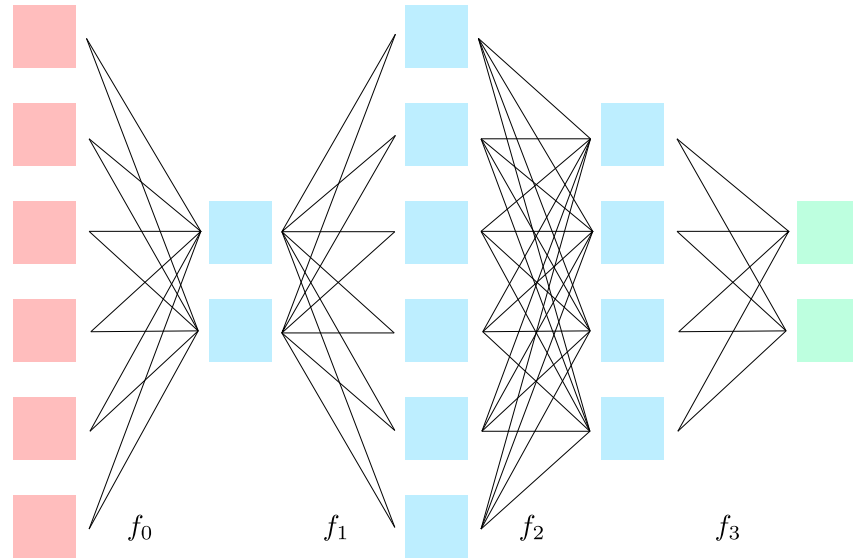


Figure 3: An example of a fully connected neural network with three hidden layers. The blue nodes are the hidden layers, the red is the input, and the green is the output. The hidden layers have width  $n_1 = 2$ ,  $n_2 = 6$ , and  $n_3 = 4$  and the number of input variables is  $n_0 = 6$ .



# Neural network

In more compact notation, we can write a fully connected neural network with  $L$  layers recursively as

$$(3) \quad f_k = \sigma_k(W_k f_{k-1} + b_k), \quad k = 1, \dots, L,$$

where

- $f_0 \in \mathbb{R}^{n_0}$  is the input to the network,
- $f_k \in \mathbb{R}^{n_k}$  for  $k = 1, \dots, L - 1$  are the values of the network at the hidden layers,
- $f_L$  is the output of the neural network,
- $n_k$  is the number of hidden nodes in the  $k^{\text{th}}$  layer,
- The weights  $W_k \in \mathbb{R}^{n_k \times n_{k-1}}$  and biases  $b_k \in \mathbb{R}^{n_k}$  are the learnable parameters in the neural network.

# Shallow Neural network

- *1-hidden layer:*

$$f_2(x) = \sigma_2(W_2\sigma_1(W_1x + b_1) + b_2).$$

- *2-hidden layers:*

$$f_3(x) = \sigma_3(W_3\sigma_2(W_2\sigma_1(W_1x + b_1) + b_2) + b_3).$$

- *3-hidden layers:*

$$f_4(x) = \sigma_4(W_4\sigma_3(W_3\sigma_2(W_2\sigma_1(W_1x + b_1) + b_2) + b_3) + b_4).$$

You will often see the last activation function  $\sigma_L$  omitted, and all activations the same, so a network with 3 hidden layers often looks like

$$f_4(x) = W_4\sigma(W_3\sigma(W_2\sigma(W_1x + b_1) + b_2) + b_3) + b_4.$$

# Loss

The output of the neural network  $f_L \in \mathbb{R}^{n_L}$  is typically fed into a loss function

$$\mathcal{L} : \mathbb{R}^{n_L} \rightarrow \mathbb{R},$$

which measures the performance of the network for the given learning task.

Typically the loss has the form

$$(4) \quad \mathcal{L}(W_1, b_1, \dots, W_L, b_L) = \sum_{i=1}^m \ell(f_L(x_i), y_i),$$

where  $(x_i, y_i)$  for  $i = 1, \dots, m$  are the training data. Here, we write  $f_L(x)$  to denote the value of the output of the network  $f_L$  given the input is  $f_0 = x$ .

Neural networks are trained by minimizing the loss function  $\mathcal{L}$  with gradient descent.

# Gradient descent

Let  $\frac{\partial \mathcal{L}}{\partial W_k}$  and  $\frac{\partial \mathcal{L}}{\partial b_k}$  denote the gradients of  $\mathcal{L}$  with respect to  $W_k$  and  $b_k$ , respectively.

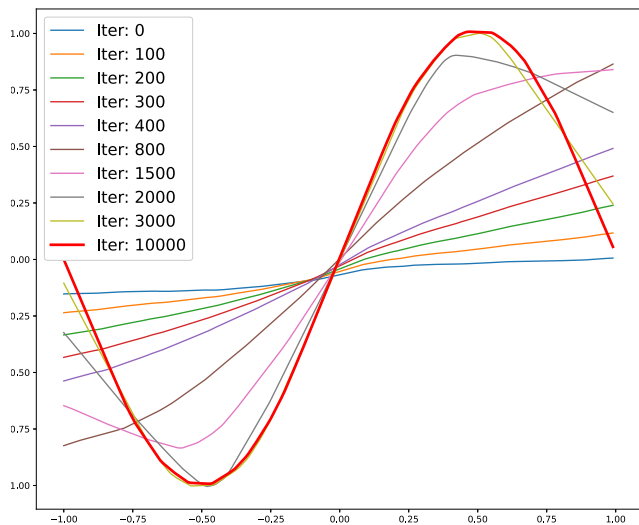
- The gradient  $\frac{\partial \mathcal{L}}{\partial W_k}$  is the  $n_k \times n_{k-1}$  matrix whose  $(i, j)$  entry is the partial derivative of  $\mathcal{L}$  in  $W_k(i, j)$ .
- The gradient  $\frac{\partial \mathcal{L}}{\partial b_k} \in \mathbb{R}^{n_k}$  is the vector whose  $i^{\text{th}}$  entry is the partial derivative of  $\mathcal{L}$  in  $b_k(i)$ .

Gradient descent for minimizing  $\mathcal{L}$  corresponds to updating the weights  $W_k$  and biases  $b_k$  according to

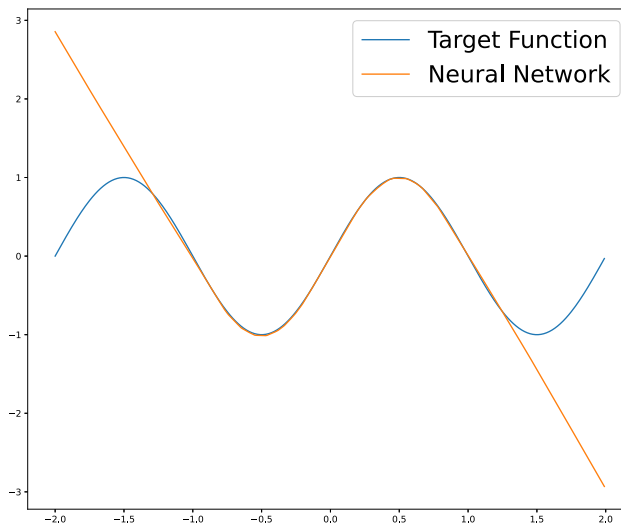
$$(5) \quad W_k^{j+1} = W_k^j - \alpha \frac{\partial \mathcal{L}}{\partial W_k} \quad \text{and} \quad b_k^{j+1} = b_k^j - \alpha \frac{\partial \mathcal{L}}{\partial b_k},$$

where  $\alpha > 0$  is the time step, also called the *learning rate*.

# Toy example



(a) Training iterations



(b) Neural network vs target

Figure 4: A toy example of fitting the function  $\sin(\pi x)$  with a 2-layer neural network with 100 hidden nodes. The loss is  $\mathcal{L} = \sum_{i=1}^m |f_L(x_i) - \sin(\pi x_i)|$  for evenly spaced points  $-1 = x_1 \leq x_2 \leq \dots \leq x_m = 1$ .

# Stochastic Gradient Descent (SGD)

For modern machine learning problems with very large training sets, it is sometimes impractical to compute the full gradients  $\frac{\partial \mathcal{L}}{\partial W_k}$  and  $\frac{\partial \mathcal{L}}{\partial b_k}$ , since the loss involves *all* of the training data.

*Stochastic gradient descent (SGD)* fixes this by computing the gradient of the loss over a random subset of the training data

$$\tilde{\mathcal{L}}(W_1, b_1, \dots, W_L, b_L) = \sum_{i \in I} \ell(f_L(x_i), y_i),$$

where  $I \subset \{1, 2, \dots, n\}$  is a random subset, called a *mini-batch*. The mini-batch changes at each iteration of SGD.

One pass over all the mini-batches in the dataset is called an *epoch*, and training usually proceeds for some number of epochs, say 100.

# Momentum descent

Various other tricks are used in the optimization, such as *momentum*

$$W_k^{j+1} = W_k^j - \alpha \frac{\partial \mathcal{L}}{\partial W_k} + \beta(W_k^j - W_k^{j-1}),$$

and

$$b_k^{j+1} = b_k^j - \alpha \frac{\partial \mathcal{L}}{\partial b_k} + \beta(b_k^j - b_k^{j-1}),$$

where  $\beta \in [0, 1]$  is the momentum parameter. Momentum can help to speed up convergence of gradient descent.

We have already analyzed gradient descent, and will study SGD and momentum descent later in the course.

# Back Propagation

For notational simplicity, we will write

$$(6) \quad z_k = W_k f_{k-1} + b_k,$$

so that  $f_k = \sigma_k(z_k)$ . Let  $\frac{\partial \mathcal{L}}{\partial z_k} \in \mathbb{R}^{n_k}$  denote the gradient of  $\mathcal{L}$  with respect to  $z_k$ . We also let  $D_k$  be the diagonal  $n_k \times n_k$  matrix with diagonal entries given by the vector  $\sigma'_k(z_k)$ . That is

$$D_k = \text{diag}(\sigma'_k(z_k)).$$

**Theorem 1** (Back propagation). *For  $k = 2, \dots, L$  we have*

$$(7) \quad \frac{\partial \mathcal{L}}{\partial z_{k-1}} = D_{k-1} W_k^T \frac{\partial \mathcal{L}}{\partial z_k},$$

$$(8) \quad \frac{\partial \mathcal{L}}{\partial W_k} = \frac{\partial \mathcal{L}}{\partial z_k} f_{k-1}^T, \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial b_k} = \frac{\partial \mathcal{L}}{\partial z_k}.$$

Proof:

$$z_k = W_k f_{k-1} + b_k$$



(8)

$$z_k(i) = \sum_{l=1}^{n_{k-1}} \underline{w_{k(i,l)}} f_{k-1}(l) + b_k(i)$$

$$\frac{\partial \mathcal{L}}{w_{k(i,j)}} = \frac{\partial \mathcal{L}}{\partial z_k(i)} \frac{\partial z_k(i)}{\partial w_{k(i,j)}} = \frac{\partial \mathcal{L}}{\partial z_k(i)} f_{k-1}(j)$$

Chain rule

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w_k} = \frac{\partial \mathcal{L}}{\partial z_k} f_{k-1}^T$$

Rank-one matrix.

$$\frac{\partial \mathcal{L}}{\partial z_{k-1}} = D_{k-1} W_k^T \frac{\partial \mathcal{L}}{\partial z_k},$$

$n_k \times 1$

$1 \times n_{k-1}$

(7)

$$z_k = W_k f_{k-1} + b_k$$

$$= W_k \sigma_{k-1}(z_{k-1}) + b_k$$

$$z_k(i) = \sum_{l=1}^{n_{k-1}} w_k(i,l) \sigma_{k-1}(z_{k-1}(l)) + b_k(i)$$

$$\frac{\partial z_k(i)}{\partial z_{k-1}(j)} = w_k(i,j) \sigma'_{k-1}(z_{k-1}(j))$$

$\downarrow$   $l=j$

Chain rule

$$\frac{\partial L}{\partial z_{k-1}(j)} = \sum_{i=1}^{n_k} \frac{\partial L}{\partial z_k(i)} \frac{\partial z_k(i)}{\partial z_{k-1}(j)}$$

$$= \sum_{i=1}^{n_k} \frac{\partial L}{\partial z_k(i)} w_k(i,j) \sigma'_{k-1}(z_{k-1}(j))$$

$$= \left[ w_k^T \frac{\partial L}{\partial z_k} \right] (j) \sigma'_{k-1}(z_{k-1}(j))$$

$$= \left( D_{k-1} w_k^T \frac{\partial L}{\partial z_k} \right) (j)$$

Since  $D_{k-1} = \text{diag}(\sigma'_{k-1}(1), \dots, \sigma'_{k-1}(n_{k-1}))$ .

$$\Rightarrow \frac{\partial L}{\partial z_{k-1}} = \left[ D_{k-1} w_k^T \frac{\partial L}{\partial z_k} \right]$$



































# Intro to Pytorch ([.ipynb](#))