

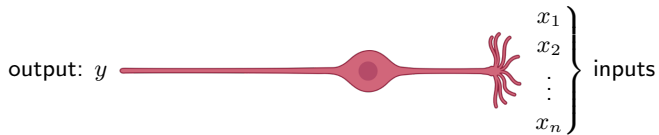
Neural Networks: Perceptrons

BIOE 498/598 PJ

Spring 2021

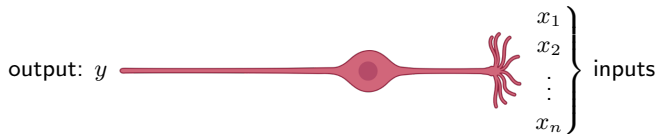
The artificial neuron

A neuron connects a series on inputs (dendrites) to an output (axon).



The artificial neuron

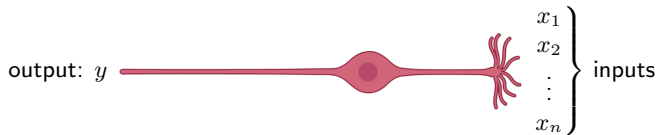
A neuron connects a series on inputs (dendrites) to an output (axon).



Our model needs to include two processes:

1. The cell body (soma) combines all of the n inputs.
2. If the combined input exceeds a threshold, the output fires.

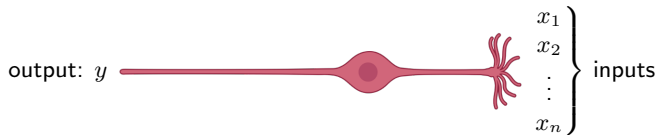
Modeling the artificial neuron



Let's model the combined input z as a linear combination of the inputs.

$$z = w_1x_1 + w_2x_2 + \cdots + w_nx_n = \mathbf{w} \cdot \mathbf{x}$$

Modeling the artificial neuron



Let's model the combined input z as a linear combination of the inputs.

$$z = w_1x_1 + w_2x_2 + \cdots + w_nx_n = \mathbf{w} \cdot \mathbf{x}$$

For now, let's assume the neuron "fires" based on the sign of z :

$$y = \begin{cases} +1, & \mathbf{w} \cdot \mathbf{x} > 0 \\ -1, & \mathbf{w} \cdot \mathbf{x} < 0 \end{cases}$$

Wait, what happened to the intercept?

Our perceptron fires using the rule

$$y = \begin{cases} +1, & \mathbf{w} \cdot \mathbf{x} > 0 \\ -1, & \mathbf{w} \cdot \mathbf{x} < 0 \end{cases}$$

Does this mean the perceptron hyperplane always passes through the origin?

Wait, what happened to the intercept?

Our perceptron fires using the rule

$$y = \begin{cases} +1, & \mathbf{w} \cdot \mathbf{x} > 0 \\ -1, & \mathbf{w} \cdot \mathbf{x} < 0 \end{cases}$$

Does this mean the perceptron hyperplane always passes through the origin?

No. We use a common ML trick to move the *bias* (intercept) into the weight vector and expand \mathbf{x} with a dummy dimension containing 1.

$$\mathbf{w} \cdot \mathbf{x} = b \Leftrightarrow \begin{pmatrix} w_1 \\ w_2 \\ -b \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} = 0$$

Summary (so far)

- ▶ A perceptron is a simplistic model of a single neuron.
- ▶ A perceptron can learn to perform simple classification tasks using an update rule.

Summary (so far)

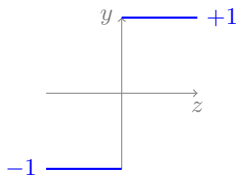
- ▶ A perceptron is a simplistic model of a single neuron.
- ▶ A perceptron can learn to perform simple classification tasks using an update rule.
- ▶ **Imagine what a network of millions of perceptrons can learn!**

Any nonlinearity will do

Any nonlinear function can be an activation function.

Sign/step activation

$$\text{sgn}(z) = \begin{cases} +1, & z > 0 \\ -1, & z < 0 \end{cases}$$

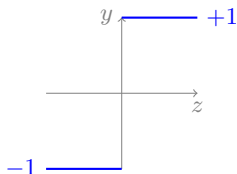


Any nonlinearity will do

Any nonlinear function can be an activation function.

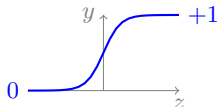
Sign/step activation

$$\text{sgn}(z) = \begin{cases} +1, & z > 0 \\ -1, & z < 0 \end{cases}$$



Sigmoid activation

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

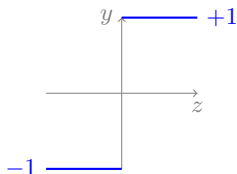


Any nonlinearity will do

Any nonlinear function can be an activation function.

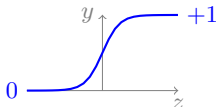
Sign/step activation

$$\text{sgn}(z) = \begin{cases} +1, & z > 0 \\ -1, & z < 0 \end{cases}$$



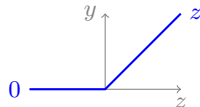
Sigmoid activation

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



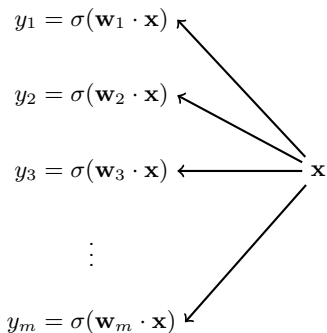
Rectified linear unit activation

$$\text{ReLU}(z) = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$$



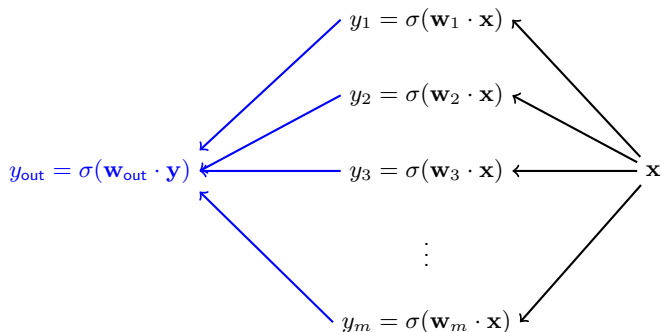
Multi-neuron (wide) perceptrons

Neural networks use multiple neurons to learn different features from the **same inputs**.



Multi-neuron (wide) perceptrons

Neural networks use multiple neurons to learn different features from the **same inputs**.



The outputs of each neuron are collected into a single neuron to predict the final class.

A matrix formalism for perceptrons

Consider a stack of m neurons that are all connected to the same input \mathbf{x} .

$$z_1 = \mathbf{w}_1 \cdot \mathbf{x}$$

$$z_2 = \mathbf{w}_2 \cdot \mathbf{x}$$

$$\vdots$$

$$z_m = \mathbf{w}_m \cdot \mathbf{x}$$

A matrix formalism for perceptrons

Consider a stack of m neurons that are all connected to the same input \mathbf{x} .

$$z_1 = \mathbf{w}_1 \cdot \mathbf{x}$$

$$z_2 = \mathbf{w}_2 \cdot \mathbf{x}$$

$$\vdots$$

$$z_m = \mathbf{w}_m \cdot \mathbf{x}$$

The stack can be written as the product of the input \mathbf{x} and a weight matrix

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

where each row in \mathbf{W} contains the weights for a single neuron

$$\mathbf{W} = \begin{pmatrix} \leftarrow \mathbf{w}_1 \rightarrow \\ \leftarrow \mathbf{w}_2 \rightarrow \\ \vdots \\ \leftarrow \mathbf{w}_m \rightarrow \end{pmatrix}.$$

Elementwise activation functions

Let's define an *elementwise* sigmoid activation function

$$\boldsymbol{\sigma}(\mathbf{z}) = \begin{pmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \vdots \\ \sigma(z_n) \end{pmatrix}.$$

Elementwise activation functions

Let's define an *elementwise* sigmoid activation function

$$\boldsymbol{\sigma}(\mathbf{z}) = \begin{pmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \vdots \\ \sigma(z_n) \end{pmatrix}.$$

A stack of m neurons can be written as

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{z})$$

Elementwise activation functions

Let's define an *elementwise* sigmoid activation function

$$\boldsymbol{\sigma}(\mathbf{z}) = \begin{pmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \vdots \\ \sigma(z_n) \end{pmatrix}.$$

A stack of m neurons can be written as

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{z})$$

or, more succinctly as

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{W}\mathbf{x})$$

Elementwise activation functions

Let's define an *elementwise* sigmoid activation function

$$\boldsymbol{\sigma}(\mathbf{z}) = \begin{pmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \vdots \\ \sigma(z_n) \end{pmatrix}.$$

A stack of m neurons can be written as

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{z})$$

or, more succinctly as

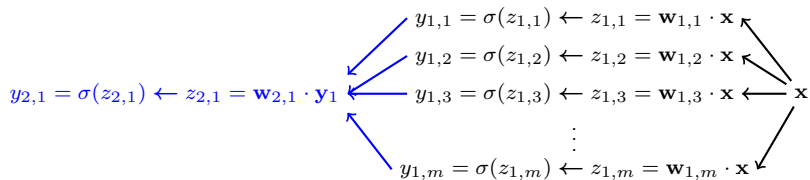
$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{W}\mathbf{x})$$

where

$$\dim(\mathbf{y}) = m \times 1, \quad \dim(\mathbf{z}) = m \times 1$$

$$\dim(\mathbf{W}) = m \times n, \quad \dim(\mathbf{x}) = n \times 1$$

Completing our matrix formalism



Completing our matrix formalism

Diagram illustrating the forward pass of a neural network layer:

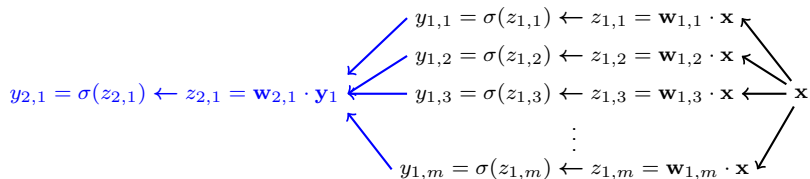
$$\begin{aligned} y_{1,1} &= \sigma(z_{1,1}) \leftarrow z_{1,1} = \mathbf{w}_{1,1} \cdot \mathbf{x} \\ y_{1,2} &= \sigma(z_{1,2}) \leftarrow z_{1,2} = \mathbf{w}_{1,2} \cdot \mathbf{x} \\ y_{1,3} &= \sigma(z_{1,3}) \leftarrow z_{1,3} = \mathbf{w}_{1,3} \cdot \mathbf{x} \\ &\vdots \\ y_{1,m} &= \sigma(z_{1,m}) \leftarrow z_{1,m} = \mathbf{w}_{1,m} \cdot \mathbf{x} \end{aligned}$$

The output vector \mathbf{y}_1 is then used to compute the next layer's input:

$$y_{2,1} = \sigma(z_{2,1}) \leftarrow z_{2,1} = \mathbf{w}_{2,1} \cdot \mathbf{y}_1$$

$$\mathbf{y}_2 = \boldsymbol{\sigma}(\mathbf{z}_2) \quad \leftarrow \quad \mathbf{z}_2 = \mathbf{W}_2 \mathbf{y}_1 \quad \leftarrow \quad \mathbf{y}_1 = \boldsymbol{\sigma}(\mathbf{z}_1) \quad \leftarrow \quad \mathbf{z}_1 = \mathbf{W}_1 \mathbf{x}$$

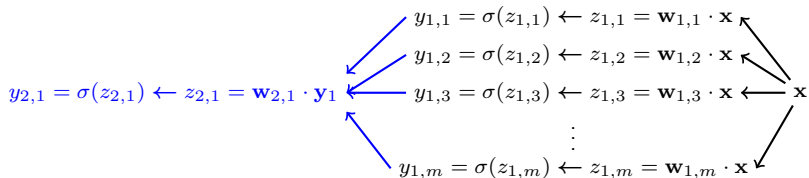
Completing our matrix formalism



$$\mathbf{y}_2 = \sigma(\mathbf{z}_2) \quad \leftarrow \quad \mathbf{z}_2 = \mathbf{W}_2 \mathbf{y}_1 \quad \leftarrow \quad \mathbf{y}_1 = \sigma(\mathbf{z}_1) \quad \leftarrow \quad \mathbf{z}_1 = \mathbf{W}_1 \mathbf{x}$$

$$\mathbf{y}_2 = \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1 \mathbf{x})))$$

Completing our matrix formalism



$$\mathbf{y}_2 = \sigma(\mathbf{z}_2) \quad \leftarrow \quad \mathbf{z}_2 = \mathbf{W}_2 \mathbf{y}_1 \quad \leftarrow \quad \mathbf{y}_1 = \sigma(\mathbf{z}_1) \quad \leftarrow \quad \mathbf{z}_1 = \mathbf{W}_1 \mathbf{x}$$

$$\mathbf{y}_2 = \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1 \mathbf{x})))$$

In general, a network with d layers is

$$\mathbf{y}_d = \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1 \mathbf{x}))))))).$$

Deep learning with neural networks

$$\mathbf{y}_d = \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x})))))))$$

- ▶ The number of neurons in each layer i is the *width* of the layer.
 - ▶ If the $(i - 1)$ th layer has n outputs and the i th layer has m outputs, the weight matrix \mathbf{W}_i has dimensions $m \times n$.
 - ▶ The dimensions of the inputs \mathbf{x} and outputs \mathbf{y}_d are fixed by the problem.
 - ▶ Layer 1 is called the *input layer*, and layer d is the *output layer*.
 - ▶ We can use as many nodes as we want in the *hidden* layers.

Deep learning with neural networks

$$\mathbf{y}_d = \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x})))))))$$

- ▶ The number of neurons in each layer i is the *width* of the layer.
 - ▶ If the $(i - 1)$ th layer has n outputs and the i th layer has m outputs, the weight matrix \mathbf{W}_i has dimensions $m \times n$.
 - ▶ The dimensions of the inputs \mathbf{x} and outputs \mathbf{y}_d are fixed by the problem.
 - ▶ Layer 1 is called the *input layer*, and layer d is the *output layer*.
 - ▶ We can use as many nodes as we want in the *hidden* layers.
- ▶ The number of layers d is the *depth* of the neural network.

Deep learning with neural networks

$$\mathbf{y}_d = \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x})))))))$$

- ▶ The number of neurons in each layer i is the *width* of the layer.
 - ▶ If the $(i - 1)$ th layer has n outputs and the i th layer has m outputs, the weight matrix \mathbf{W}_i has dimensions $m \times n$.
 - ▶ The dimensions of the inputs \mathbf{x} and outputs \mathbf{y}_d are fixed by the problem.
 - ▶ Layer 1 is called the *input layer*, and layer d is the *output layer*.
 - ▶ We can use as many nodes as we want in the *hidden* layers.
- ▶ The number of layers d is the *depth* of the neural network.
- ▶ *Deep learning* means $d > 2$.

The importance of nonlinearity

The nonlinear functions (like σ) sandwiched between the layers are critical to deep learning.

Let's imagine what would happen if we removed them:

$$\begin{aligned} \mathbf{y}_d &= \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x}))))))) \\ &= \mathbf{W}_d(\mathbf{W}_{d-1}(\cdots \mathbf{W}_2(\mathbf{W}_1\mathbf{x}))) \\ &= \mathbf{W}_d\mathbf{W}_{d-1}\cdots\mathbf{W}_2\mathbf{W}_1\mathbf{x} \\ &= \widetilde{\mathbf{W}}\mathbf{x} \end{aligned}$$

The importance of nonlinearity

The nonlinear functions (like σ) sandwiched between the layers are critical to deep learning.

Let's imagine what would happen if we removed them:

$$\begin{aligned} \mathbf{y}_d &= \sigma(\mathbf{W}_d(\sigma(\mathbf{W}_{d-1}(\cdots \sigma(\mathbf{W}_2(\sigma(\mathbf{W}_1\mathbf{x}))))))) \\ &= \mathbf{W}_d(\mathbf{W}_{d-1}(\cdots \mathbf{W}_2(\mathbf{W}_1\mathbf{x}))) \\ &= \mathbf{W}_d\mathbf{W}_{d-1}\cdots\mathbf{W}_2\mathbf{W}_1\mathbf{x} \\ &= \widetilde{\mathbf{W}}\mathbf{x} \end{aligned}$$

Without the activation functions, the entire neural network reduces to a single linear system!

Why do we want deep networks?

- ▶ The *Universal Approximation Theorem* states that given enough neurons, a 2-layer (input/output) perceptron can learn any reasonable function.
- ▶ Neural networks are therefore universal function approximators.

Why do we want deep networks?

- ▶ The *Universal Approximation Theorem* states that given enough neurons, a 2-layer (input/output) perceptron can learn any reasonable function.
- ▶ Neural networks are therefore universal function approximators.
- ▶ Unfortunately, the theorem does not tell us how many neurons we need to approximate a given function.
- ▶ For complicated functions, evidence suggests the number is enormous!

Why do we want deep networks?

- ▶ The *Universal Approximation Theorem* states that given enough neurons, a 2-layer (input/output) perceptron can learn any reasonable function.
- ▶ Neural networks are therefore universal function approximators.
- ▶ Unfortunately, the theorem does not tell us how many neurons we need to approximate a given function.
- ▶ For complicated functions, evidence suggests the number is enormous!
- ▶ Our brains are very deep, so it's reasonable to believe that deep networks learn more efficiently than wide ones.
- ▶ In practice this is almost certainly true.
- ▶ Deep learning reduces the total number of neurons needed to learn a function since each of the d layers needs fewer than $1/d$ -times the number of neurons.

Why do deeper networks learn better?

We can understand deep networks using examples from *feature engineering*.

Imagine you wanted to learn a Michaelis-Menten function.

Why do deeper networks learn better?

We can understand deep networks using examples from *feature engineering*.

Imagine you wanted to learn a Michaelis-Menten function.

Single Layer

The diagram illustrates the Michaelis-Menten equation as a single-layer model. On the left is the equation $\frac{V_{\max}[S]}{K_m + [S]}$. On the right are the parameters V_{\max} , $[S]$, and K_m . Three arrows point from these parameters to the equation: one from V_{\max} to the numerator's coefficient, one from $[S]$ to the numerator's variable, and one from K_m to the denominator's constant term.

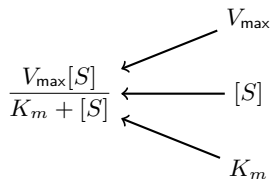
$$\frac{V_{\max}[S]}{K_m + [S]}$$

Why do deeper networks learn better?

We can understand deep networks using examples from *feature engineering*.

Imagine you wanted to learn a Michaelis-Menten function.

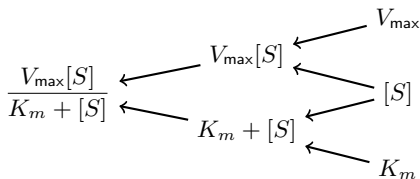
Single Layer



A diagram showing a single layer neural network. On the left is the output expression $\frac{V_{\max}[S]}{K_m + [S]}$. On the right are three inputs: V_{\max} , $[S]$, and K_m . Three arrows point from each input to the output expression.

$$\frac{V_{\max}[S]}{K_m + [S]}$$

Two Layers



A diagram showing a two layer neural network. On the left is the output expression $\frac{V_{\max}[S]}{K_m + [S]}$. In the middle are two hidden layer nodes: $V_{\max}[S]$ and $K_m + [S]$. On the right are three inputs: V_{\max} , $[S]$, and K_m . Arrows show the flow from inputs to hidden nodes and from hidden nodes to the output. Specifically, V_{\max} points to $V_{\max}[S]$, $[S]$ points to both $V_{\max}[S]$ and $K_m + [S]$, and K_m points to $K_m + [S]$. Then, $V_{\max}[S]$ points to the output, and $K_m + [S]$ also points to the output.

$$\frac{V_{\max}[S]}{K_m + [S]}$$

Why do deeper networks learn better?

We can understand deep networks using examples from *feature engineering*.

Imagine you wanted to learn a Michaelis-Menten function.

Single Layer

A diagram showing a single layer of a neural network. On the left is the output expression $\frac{V_{\max}[S]}{K_m + [S]}$. On the right are three input features: V_{\max} , $[S]$, and K_m . Three arrows point from each of these input features to the output expression, indicating that the output is a linear combination of these features.

$$\frac{V_{\max}[S]}{K_m + [S]}$$

Two Layers

A diagram showing a two-layer neural network. On the left is the output expression $\frac{V_{\max}[S]}{K_m + [S]}$. In the middle is a hidden layer with two nodes: $V_{\max}[S]$ and $K_m + [S]$. On the right are the original input features: V_{\max} , $[S]$, and K_m . Arrows show the flow of information: V_{\max} and $[S]$ feed into the $V_{\max}[S]$ node; $[S]$ and K_m feed into the $K_m + [S]$ node. Then, the two hidden nodes feed into the final output expression.

$$\frac{V_{\max}[S]}{K_m + [S]}$$

Each layer in the network only needs to improve the features for the next layer.

Summary

- ▶ Deep neural networks are built from layers in artificial neurons.
- ▶ Each neuron has the power of a linear classifier.
- ▶ Layers **must** be separated by nonlinear activation functions.
- ▶ Neural networks can learn nearly any function, but deep networks learn more efficiently.
- ▶ Each layer creates features for the subsequent layers to improve learning.

Summary

- ▶ Deep neural networks are built from layers in artificial neurons.
- ▶ Each neuron has the power of a linear classifier.
- ▶ Layers **must** be separated by nonlinear activation functions.
- ▶ Neural networks can learn nearly any function, but deep networks learn more efficiently.
- ▶ Each layer creates features for the subsequent layers to improve learning.
- ▶ **Next time:** Training a neural network to learn Q -factors.