

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

Choices of activation function

Our linear model is entirely **defined** by its input: the number of **weights** in a neuron is just the number of **inputs** m .

But our **activation** function is up to us to decide: what works best?

Trying out linear activation

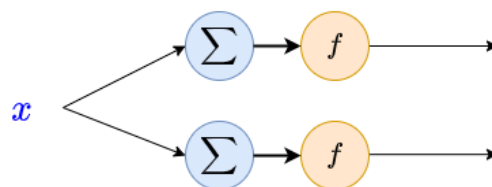
The simplest assumption would be to just use the **identity** function

$$f(z) = z \quad (1)$$

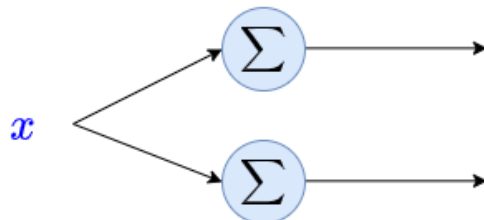
We might hope that we can combine a bunch of simple, **linear** models, and get a more sophisticated model. Why bother having a **nonlinear** activation at all?

Well, it turns out, combining **multiple** linear layers doesn't make our model any stronger. Let's try an example: we'll take a network with 2 layers, two neurons each.

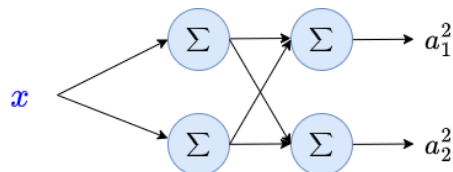
Let's look at layer 1:



Since the activation function has **no effect** on our result, we can **omit** it:

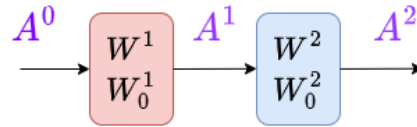


And now, we can show our **full** network:



Linear Layers: An example

We'll assume **two** inputs $A_0 = [x_1, x_2]^T$. For our sanity, we'll lump all of the weights in each layer:



We'll leave out W_0 terms to make it more readable, but the same will apply.

Layer 1:

$$A^1 = (W^1)^T A_0 \quad (2)$$

Layer 2:

$$A^2 = \overbrace{(W^2)^T (W^1)^T}^{\text{Weight matrices}} A_0 \quad (3)$$

The full function for this equation is two matrices, **multiplied** by our input vector.

Let's take an arbitrary example:

$$W^1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad W^2 = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \quad (4)$$

Our equation becomes:

$$A^2 = \overbrace{\begin{bmatrix} 5 & 7 \\ 6 & 8 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}}^{\text{Transposed matrices}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (5)$$

We created this function by applying two matrices separately. But, can't we **combine** them?

$$A^2 = \begin{bmatrix} 19 & 43 \\ 22 & 50 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (6)$$

Wait, but this looks like a **one-layer** network with those weights! The second layer is **pointless**, we could have represented it with a single layer...

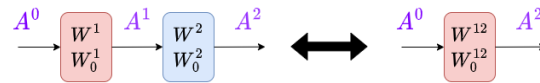
$$(W^{12})^T = \begin{bmatrix} 19 & 43 \\ 22 & 50 \end{bmatrix} \quad (7)$$

The problem with linear networks

In fact, this is true in general: we can always take our **two** linear layers and combine them into **one**.

$$(\mathbf{W}^2)^\top (\mathbf{W}^1)^\top = \mathbf{W}^{12} \quad (8)$$

Our network is **equivalent** to the supposedly "simpler" one-layer network.



What if we have more layers? Well, we can just combine them one-by-one. At the end, we're just left with one layer:

$$(\mathbf{W}^L)^\top (\mathbf{W}^{L-1})^\top \dots (\mathbf{W}^2)^\top (\mathbf{W}^1)^\top = \mathbf{W}^X \quad (9)$$

And so, we can't just use linear layers: we **need a nonlinear** activation function.

Concept 1

Having multiple consecutive **linear layers** (i.e. layers with linear **activation** functions) is **equivalent** to having one linear layer in its place.

This means that we do not expand our **hypothesis** class by using more linear layers: we have to use **nonlinear** activation functions.

If we use something **nonlinear**...

$$\mathbf{A}^2 = \mathbf{f} \left((\mathbf{W}^2)^\top \mathbf{A}^1 \right) \quad (10)$$

We get something that doesn't **simplify**:

$$\mathbf{A}^2 = \mathbf{f} \left((\mathbf{W}^2)^\top \overbrace{\mathbf{f} \left((\mathbf{W}^1)^\top \mathbf{x} \right)}^{\mathbf{A}^1} \right) \quad (11)$$

This is ugly, but we don't have to worry about the details.

Which is what we want!