

# Econ 425 Week 7

Manu Navjeevan

February 16, 2021

## 1 Neural Net Review

As we went over last time, a neural network consists of nodes and edges, arranged into layers.

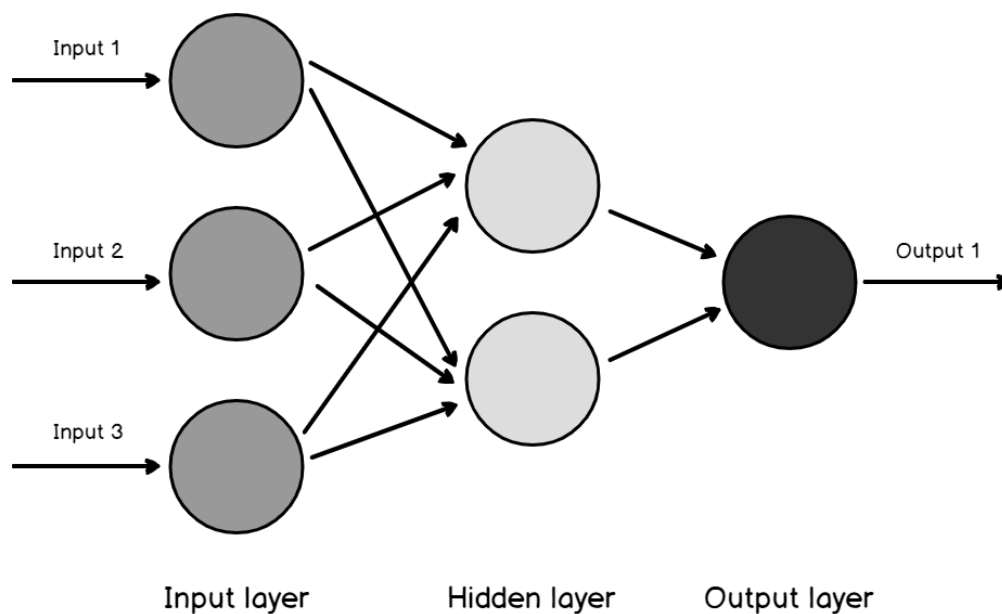


Figure 1: A simple neural network

Each neuron takes some inputs (the outputs of the neurons that are pointing to it),  $a_1, \dots, a_p$  and outputs a vector of the form:

$$\pi \left( \underbrace{\beta_0}_{\text{Bias Term}} + \beta_1 a_1, \dots, \beta_p a_p \right).$$

where  $\pi(\cdot)$  is some pre-specified activation function. The bias term  $\beta_0$  and each of the weights  $\beta_1, \dots, \beta_p$  are parameters that we'd like to “learn” using our data.

Specifically, we would like to “learn” the parameters that minimize some cost function of our output. An example of this would be to choose the parameters that minimize the mean squared error on our training data:

$$MSE(\cdot) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2.$$

In order to do this minimization, we essentially need to do a gradient descent on our model parameters. However, computing this gradient is not straightforward. This brings us to back propogation, which is a

way of computing the gradient that takes advantage of the feed-forward nature of the neural network.

## 2 Back Propagation

To begin with, consider a very simple neural network.

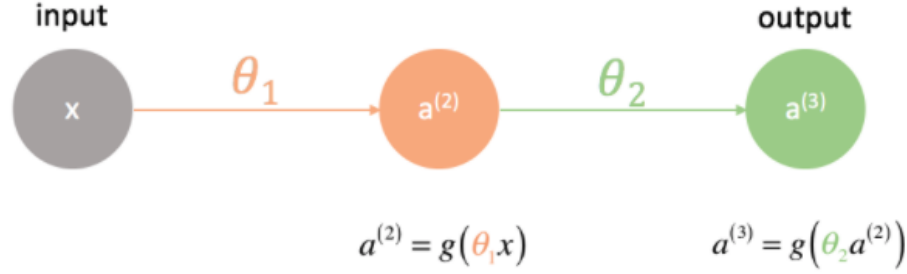


Figure 2: A very simple neural net

The output can be written as  $\hat{y}_i = g(\theta_2 \cdot g(\theta_1 \cdot x_i))$ . My cost function is the MSE:

$$\begin{aligned}
 \mathcal{C}(\theta_1, \theta_2) &= \sum_{i=1}^n (\hat{y}_i - y_i)^2 \\
 &= \sum_{i=1}^n \left( g(\theta_2 \cdot g(\theta_1 \cdot x_i)) - y_i \right)^2
 \end{aligned}$$

I would like to compute the gradient with respect to  $\theta_1$  and  $\theta_2$ . Since this is fairly simple in this setting we do it by hand and then examine the results to gain some intuition for back propagation:

$$\begin{aligned}
 \frac{\partial \mathcal{C}}{\partial \theta_2} &= 2 \sum_{i=1}^n \left( \left( g(\theta_2 \cdot g(\theta_1 \cdot x_i)) - y_i \right) \right) g'(\theta_2 \cdot g(\theta_1 \cdot x_i)) \cdot g(\theta_1 \cdot x_i) \\
 \frac{\partial \mathcal{C}}{\partial \theta_1} &= 2 \sum_{i=1}^n \left( \left( g(\theta_2 \cdot g(\theta_1 \cdot x_i)) - y_i \right) \right) g'(\theta_2 \cdot g(\theta_1 \cdot x_i)) \cdot g'(\theta_1 \cdot x_i) x_i
 \end{aligned}$$

Comparing the partial derivatives with respect to  $\theta_2$  and  $\theta_1$  note that the partial derivative with respect to  $\theta_1$  can be computed iteratively.

1. First compute the derivative of the activation function in the last layer:  $g'(\theta_2 \cdot g(\theta_1 \cdot x_1))$ . Use these to evaluate the derivative with respect to  $\theta_2$ .
2. Compute the derivative with respect to the activation function in the first layer:  $g'(\theta_1 \cdot x_2)$ . Multiply these with the derivative with the with respect to the activation function in the last layer from the first step and use the resulting product to compute the derivative with respect to  $\theta_1$ .

This provides the intuition for how back-propagation works to compute the total derivative with respect to all the weight/bias terms. To compute the total derivative we will start at the output layer, compute the derivatives of the activation function with respect to our parameters at this layer, and use these derivatives to calculate the derivative of the cost function with respect to the parameters in the last layer. To compute the derivative of the cost function with respect to parameters in the next layer back we will multiple these activation function derivatives from the final layer with the derivatives of the activation function from the next layer back layer. We then move to the second to last layer, calculate the derivatives of the activation

function at each neuron and multiply these with the derivatives in the prior two layers and so on until we have the derivative with respect to every parameter.

Some notes:

- Note that the first term inside the summand in each of these derivatives is  $\left(g(\theta_2 \cdot g(\theta_1 \cdot x_i)) - y_i\right)$ . This means that observations for which our neural network prediction is far from the actual value will have a larger influence on our gradient, as we will try to correct for this and “fix” the loss function.
- In the very toy example we have drawn up here we are dealing with a (trivially) “fully connected” neural network. In practice, the multiplications at each step will depend on which nodes are connected to which.

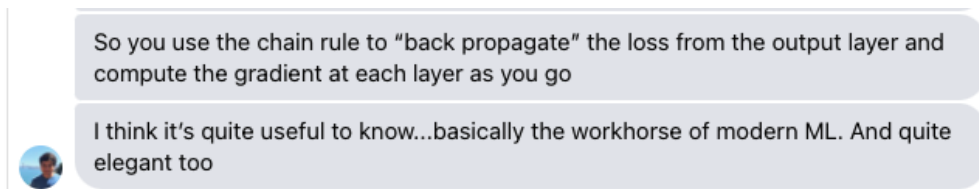


Figure 3: Some thoughts from my college friend Shouvik, who is working in ML.