# An Overview of Neural Networks

Harley Clifton, Becky Catlett, Natasha Gesker, and Eliot Liucci

2023-10-15

## Contents

# Introduction

## Deep Learning (Becky)

## The Neural Network Model (Eliot)

A neural network is a type of deep learning algorithm that makes use of a web of nodes to predict or classify data. The complexity of a neural network can vary greatly based on what task is required. As the name suggests, they are similar in function to neurons in a brain. The model takes in information through the *input layer*, which then activates various nodes in the *hidden layers*, and then a result is produced.

### The Input Layer

The input layer is where data can be input into the model. If we have $p$ input variables, which we will denote $X = X_1, X_2, ..., X_p$, then our network will have $p$ input nodes. Each node in future layers will depend on the value that $X_i$ holds.

### Nonlinear Activation Functions

Before we get into the hidden layer, it is important to understand what is happening at each hidden layer node. Each hidden layer node is computed by taking a weighted linear combination of the input layer and then applying a *nonlinear activation function* so that the *activation*, which is the value the node will take based on input vector $X$, will be between 0 and 1.

We will discuss two of the most common activation functions. For simpler networks, the *sigmoid* function is effective. The sigmoid function is defined as

$$S(x) = \frac{1}{1 + e^{-x}}$$

As discussed previously, the purpose of the activation function is to bring the range of values for the input layer down to any value between 0 and 1.

Another activation function that is more common in networks that require more "training" is the rectified linear activation unit function, or ReLU for short. The ReLU function is defined as

$$R(x) = \left\{ \begin{array}{ll} 0 & \text{if x} < 0 \\ x & \text{otherwise} \end{array} \right.$$

The benefits of using ReLU over Sigmoid is that ReLU can be better used for *backpropagation*, which is the main technique used to train networks. We will get more into this later.

### Hidden Layers

Hidden layers are the bread and butter of neural network models. Take, for example, the network pictured below with 4 input nodes and 2 hidden nodes.
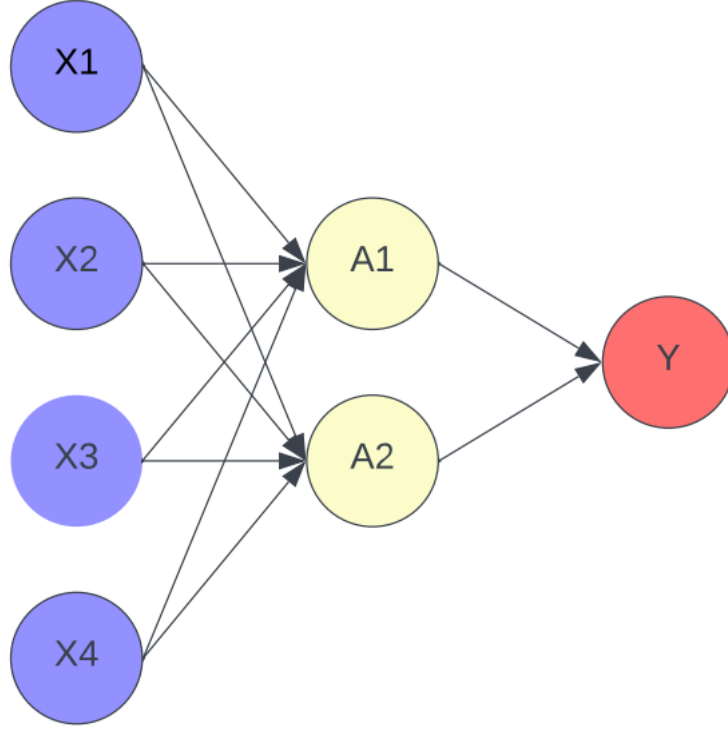
Figure 1: Example of a Simple Neural Network with 4 Input Nodes and 2 Hidden Nodes

Mathematically, we can write the *activation* of the 1st node of the hidden layer, $A_1$, as

$$A_1 = h_k(X) = g(w_{0,1} + \sum_{j=1}^{p} w_{j,1} \cdot X_j)$$

Where $g(.)$ is the nonlinear activation function of choice and $w_{j,1}$ is the weight associated with activation 1 and input node $j$. The value of $w_{0,1}$ is called the "bias" and can be added to offset the activation so that the minimum value matches what it is expected to be. For each activation of the hidden layer, we are taking a weighted sum of all nodes in the input layer. The activation function restricts the range of the values the activation can hold. For Sigmoid, it would be between 0 and 1 while ReLU would just be greater than or equal to 0. This can be generalized further for $k$ activations

$$A_k = g(w_{0,k} + \sum_{j=1}^{p} w_{j,k} \cdot X_j)$$

**Output Layer**

The output layer is what we would be predicting. For a quantitative response, we would have a single node that would hold the value we predict based on the input vector $X$. For a categorical response with $q$ levels, we would have $q$ output nodes. The output can be thought of as a linear regression model fit using the hidden layer nodes as inputs. This can be formally written as

$$f(x) = \beta_0 + \sum_{k=1}^{k} A_k \cdot \beta_k$$

## Applications of Principal Component Analysis with Neural Networks

Neural networks can grow in complexity very quickly. Given a dataset with 20 input variables, we could end up requiring many nodes in the Hidden Layer. The training process can be timely and computationally expensive.

Principal Component Analysis would allow for those 20 input variables to be trimmed down to 2 or 3 principal components. This would also theoretically cut down on the number of nodes in the Hidden Layer, thus reducing the computational cost of fitting the model while maintaining the accuracy of the model.

# Single-Layer Neural Networks (Harley)

## Example

For an example of a single-layer network, we used the `airquality` dataset from the `datasets` package in R. We trained a model to predict Ozone levels using the other predictors available. We then trained a model to predict Ozone levels using the first 3 principal components and compared the results.

| Ozone | Solar.R | Wind | Temp | Month | Day |
|---|---|---|---|---|---|
| 0.5675676 | 0.0192192 | 0.1981982 | 0.012012 | 0.000000 | 0.5675676 |
| 0.3513514 | 0.0210210 | 0.2132132 | 0.012012 | 0.003003 | 0.3513514 |
| 0.4444444 | 0.0348348 | 0.2192192 | 0.012012 | 0.006006 | 0.4444444 |
| 0.9369369 | 0.0315315 | 0.1831832 | 0.012012 | 0.009009 | 0.9369369 |
| 0.8948949 | 0.0228228 | 0.1921922 | 0.012012 | 0.018018 | 0.8948949 |

The network using the observed variables in the dataset can be seen below. The network was trained quickly and predicted Ozone levels with 5% error 90% of the time and predicted with 1% error 65% of the time. Impressive!
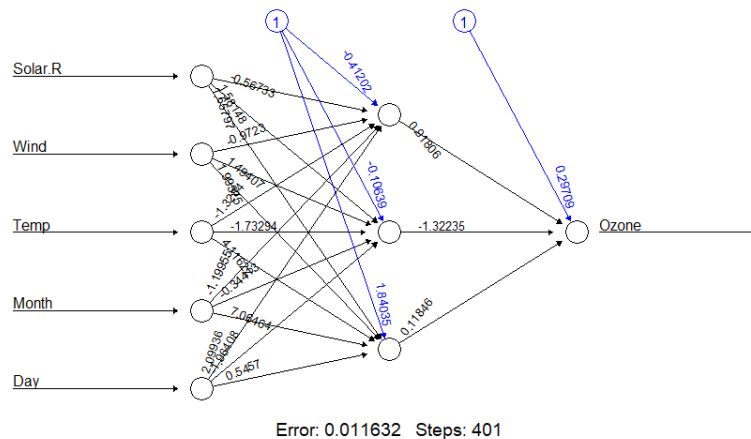


Figure 2: Neural network using observed variables to predict Ozone levels

Comparing those results to the network trained using the first 3 principal components, we got predictions within 5% of the observed values 73% of the time and predictions within 1% of the observed values 62% of the time. The performance of the network using principal components is similar with a 1% error rate but falls behind within a 5% error rate.
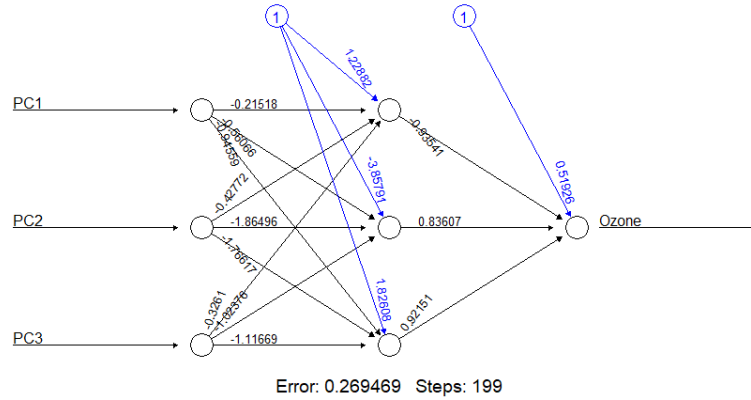
Figure 3: Neural network using first 3 principal components to predict Ozone levels

## Multi-Layer Neural Networks (Natasha)

### Example

For an example of a multi-layer network, we used the `Iris` dataset from the `datasets` package in R. We trained a model to predict Species using Sepal Length, Sepal Width, Petal Length, and Petal Width and will compare its performance to a model using the first 2 principal components to predict Species.

| Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|
| 0.6410256 | 0.4358974 | 0.1666667 | 0.0128205 | setosa |
| 0.6153846 | 0.3717949 | 0.1666667 | 0.0128205 | setosa |
| 0.5897436 | 0.3974359 | 0.1538462 | 0.0128205 | setosa |
| 0.5769231 | 0.3846154 | 0.1794872 | 0.0128205 | setosa |
| 0.6282051 | 0.4487179 | 0.1666667 | 0.0128205 | setosa |

The network that uses the 4 observed features of a flower performed well, with an accuracy of 98%. Below we can see the structure of the network.
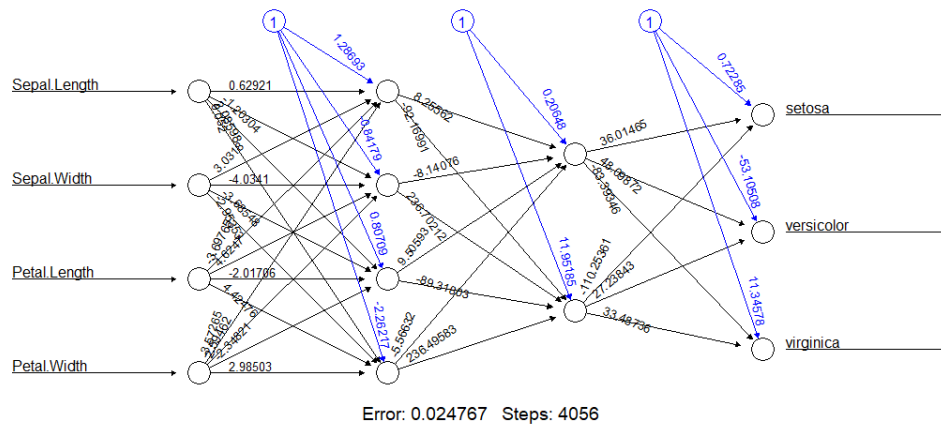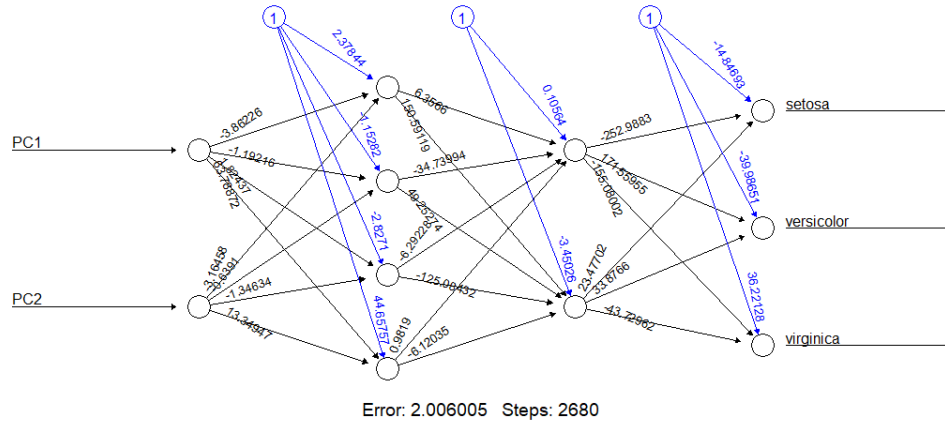


Figure 4: Neural network using observed variables to predict species

5

The network using the first 3 principal components also obtaining an accuracy of 98% with species classification, while only taking only $\frac{2}{3}$s of the time to train. We can see that network's structure below.



Error: 2.006005   Steps: 2680

The structure of these networks was the same, so it is unsurprising that the performance was also similar. A takeaway from this example is that PCA can be used to allow for simpler input layers while achieving similar results. When a data set includes thousands of observations to be used for training, reducing the time to train by $\frac{1}{3}$ can make a huge difference.

# References

https://www.datacamp.com/tutorial/neural-network-models-r

3Blue1Brown. 2017. "Neural Networks." YouTube. 2017. https://www.youtube.com/watch?v=aircAruvnK k&list=PLZHQObOWTQDNU6R1_67000Dx_ZCJB-3pi.

Abdi, Hervé, Dominique Valentin, and Betty Edelman. 1999. *Neural Networks*. 124. Sage.

Anderson, James A. 1995. *An Introduction to Neural Networks*. MIT press.

James, Gareth, Daniela Witten, Trevor Hastie, Robert Tibshirani, et al. 2013. *An Introduction to Statistical Learning*. Vol. 112. Springer.

Pramoditha, Rukshan. 2022. "Using PCA to Reduce Number of Parameters in a Neural Network by 30x Times."