

# Neural Networks For Binary Classification

mcc

3/13/2020

## Neural Networks For Binary Classification

“Machine learning is essentially a form of applied statistics with increased emphasis on the use of computers to statistically estimate complicated functions and a decreased emphasis on proving confidence intervals around these functions”

– Ian Goodfellow, et al<sup>1</sup>

### Introduction

If we discuss Neural Networks (NN), we should first consider the system we hope to emulate. Let us start with a simple count of neuronal cells in various organisms along the earth’s phylogenetic tree. We might get a better idea of the type of “computing power” these living creatures possess. See table 6.1.

**Table 6.1: Organisms Vs Number of Neurons In Each (Wikipedia)**

Organism	Common Name	Approximate Number of Neurons
<i>C. elegans</i>	roundworm	302
<i>Chrysaora fuscescens</i>	jellyfish	5,600
<i>Apis linnaeus</i>	honey bee	960,000
<i>Mus musculus</i>	mouse	71,000,000
<i>Felis silvestris</i>	cat	760,000,000
<i>Canis lupus familiaris</i>	dog	2,300,000,000
<i>Homo sapien sapien</i>	humans	100,000,000,000

This table portrays a high-level overview of the computing power of neuronal clusters and brains produced throughout evolution. However, there is one missing number worth noting. The table above does not describe the connectivity between neurons. The connectivity of neurons varies greatly from lower to higher organisms. For example, some simple animals, such as the roundworm, have only “four to eight separate branches,”<sup>2</sup> per nerve cell. While human neurons may have greater than 10,000 inter-connected synaptic junctions per neuron, thus resulting in a total of approximately 600 trillion synapses per human brain.<sup>3</sup>

Although neurons have differing morphologies, neurons in the human brain are extremely diverse. Indeed, size and shape may not be the definitive way of classifying neurons but instead by what neurotransmitters the cells secrete. “Neurotransmitters can be classified as either excitatory or inhibitory.”<sup>4</sup> Currently

<sup>1</sup>Ian Goodfellow, Yoshua Bengio, Aaron Courville, ‘Deep Learning’, MIT Press, 2016, <http://www.deeplearningbook.org>

<sup>2</sup><https://www.wormatlas.org/hermaphrodite/nervous/Neuroframeset.html>

<sup>3</sup>Shepherd, G. M. (2004), The synaptic organization of the brain (5th ed.), Oxford University Press, New York.

<sup>4</sup><https://www.kenhub.com/en/library/anatomy/neurotransmitters>

the NeuroPep database “holds 5949 non-redundant neuropeptide entries originating from 493 organisms belonging to 65 neuropeptide families.”<sup>5</sup>

---

<sup>5</sup><http://isyslab.info/NeuroPep/home.jsp>

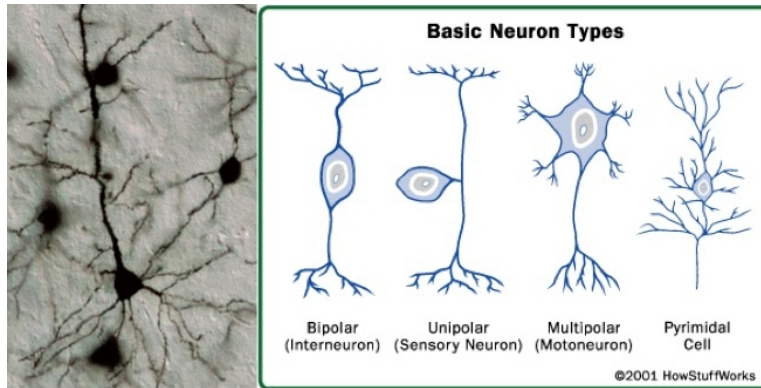


Figure 1: Basic Neuron Types and S.E.M. Image

6

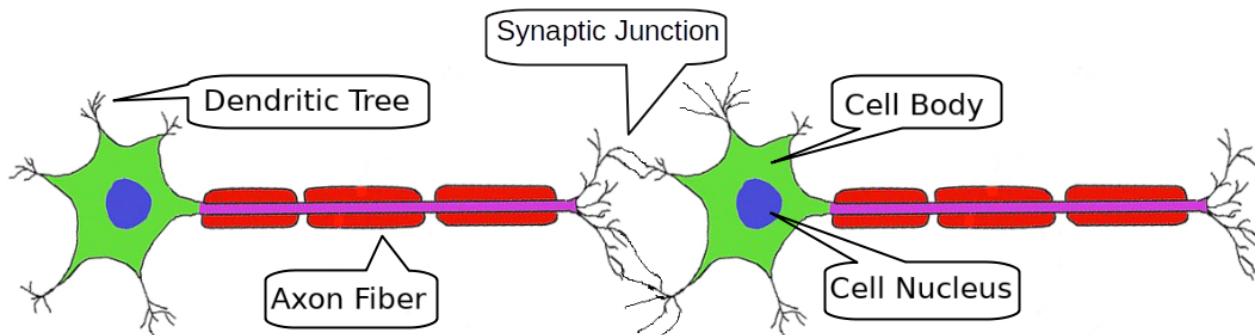


Figure 2: Two Neuron System (Image From The Public Domain)

Given an order of operation via:

*Dendrite(s) ⇒ Cell body ⇒ Fibrous Axon ⇒ Synaptic Junction or Synaptic Gap ⇒ Dendrite(s) ... Ad infinitum.*

However, nature is more subtle and intricate than to have neurons in a series, only blinking on and off, firing or not. NN are often programmed to classify dangerous road objects, as is the case of Tesla cars. The goal of a Tesla auto-piloted car is to use all available sensors to correctly classify all the conceivable circumstances on the road. On the road, a Tesla automobile uses dozens of sensors which the computer needs to evaluate and weigh the values of all these sensors to formulate a ‘decision.’ The altitude of the auto, derived from the GPS, may weigh less heavily than the speed of the vehicle or Lidar estimates on how close objects are. However, our goal of safe driving can be thwarted when an artificial intelligence system decides a truck is a sign and does not apply the brakes.<sup>7</sup>

<sup>6</sup><https://www.howstuffworks.com/>

<sup>7</sup><https://arstechnica.com/cars/2019/05/feds-autopilot-was-active-during-deadly-march-tesla-crash/>

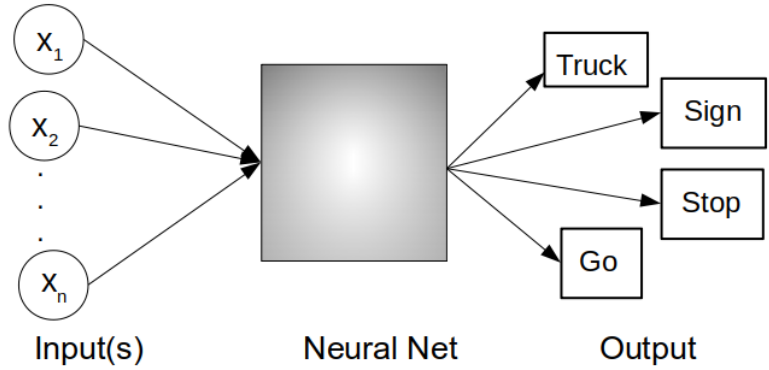


Figure 3: Goal of a Tesla Neural Networks is to generate the correct responses for its environment.

### The One Neuron System

If we investigate a one neuron system, *our* neuron could be diagrammed in four sections.<sup>8</sup>

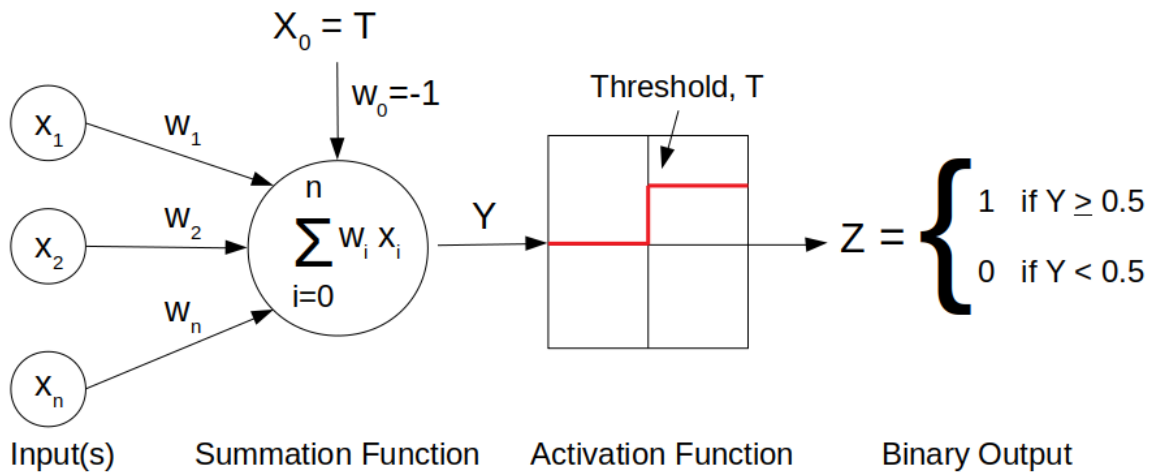


Figure 4: One Neuron Schema

If we investigate one neuron for a moment, we find two separate mathematical functions are being carried out by a single nerve cell.

### Summation Function

The first segment is a summation function. It receives the real number values from,  $x_1$  to  $x_n$ , all the branches of the dendritic trees, and multiplies them by a set of weights. These  $X$  inputs are multiplied by a set of corresponding unique weights from  $w_1$  to  $w_n$ . An analogy I prefer is of small or large rivers joining giving a total current. The current moves through the branches giving a total signal or current of sodium ions. Interestingly the summation in each neuron, while dealing with the vectors of inputs and weights, is carrying out the dot product of these vectors.

<sup>8</sup>Tom Mitchell, Machine Learning, McGraw-Hill, 1997, ISBN: 0070428077

Initially, the NN researchers used the Heaviside-Threshold Function, as shown in figure 5.4, *One Neuron System*. The benefits of step functions were their simplicity and high signal to noise ratio. While the detriments were, it is a discontinuous function, therefore not differentiable and a mathematical problem.

Let us take into account the product,  $x_0 \cdot w_0$ . If we assign  $x_0 = T$  and  $w_0 = -1$  this simply becomes a bias. This bias allows us the ability to shift our Activation Function and its inflection point in the positive or negative x-direction.

$$\hat{Y} = X^T \cdot W - Bias \equiv \sum_{i=0}^n x_i w_i - T \quad (1)$$

### Activation Functions

The second function is called an Activation Function. Once the Summation Function yields a value, its result is sent to the *Activation Function* or *Threshold Function*.

$$Z^{(1)} = f \left( \sum_{i=0}^n x_i w_i - T \right) = \{0, 1\} \quad (2)$$

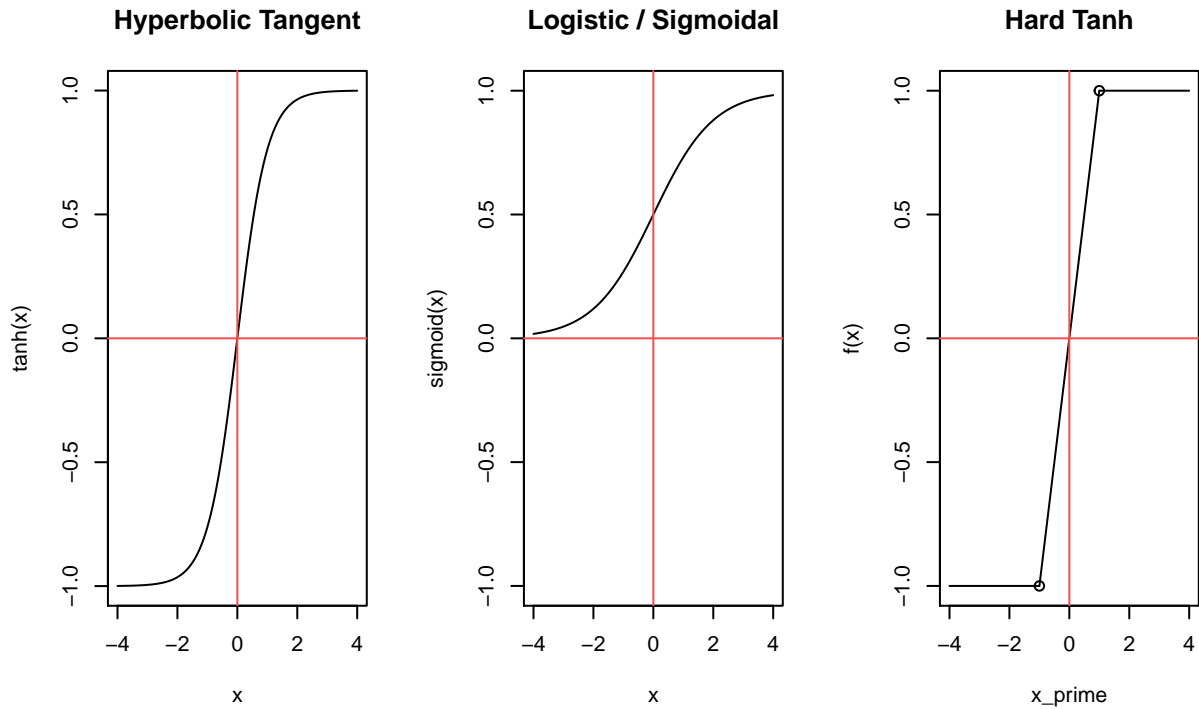
The function displayed in figure #6.4, One Neuron Schema, is a step function. However this step function has a problem mathematically, namely it is a discontinuous and therefore not differentiable. This fact is important.

Therefore several functions may be used in place of the step function. One is the hyperbolic tangent (*tanh*) function, the *sigmoidal* function, a *Hard Tanh*, a *reLU*, and *Softmax* Functions. These have certain advantages, namely they simplify the hyperbolic tangent function. Not only does the Hard Tanh and reLU simplify calculations it is useful for increasing the gain near the asymptotic limits of the sigmoidal and tanh functions. The derivatives of the sigmoidal and tanh functions are very small, near 0 and 1, while the reLU and Hard Tanh slopes are one or zero.

$$Z^{(2)} = \tanh(x) = \frac{1 - e^{-\alpha}}{1 + e^{-\alpha}} \quad : \quad \text{where } \alpha = \sum_{i=1}^n x_i w_i - T \quad (3)$$

$$Z^{(3)} = \text{sigmoid}(x) = \frac{1}{1 + e^{-\alpha}} \quad (4)$$

$$Z^{(4)} = \text{Hard Tanh}(x) = \begin{cases} 1 & x > 1 \\ x & -1 \leq x \leq 1 \\ -1 & x < -1 \end{cases} \quad (5)$$



Several alternative functions are useful for various reasons. The most common of which are Softmax and reLU functions.

Rectified Linear Activation Unit, (ReLU):

$$Z^{(5)} = ReLU = \begin{cases} x \geq 0 & y = x \\ x < 0 & y = 0 \end{cases} \quad (6)$$

### Binary Output Or Probability

In the case of real neurons, the output is off or on, zero or one. However, in the case of our electronic model, it is advantageous to calculate a probability for greater interpretability.

The Softmax function may appear like the Sigmoid function from above but it differs in major ways.<sup>9</sup>

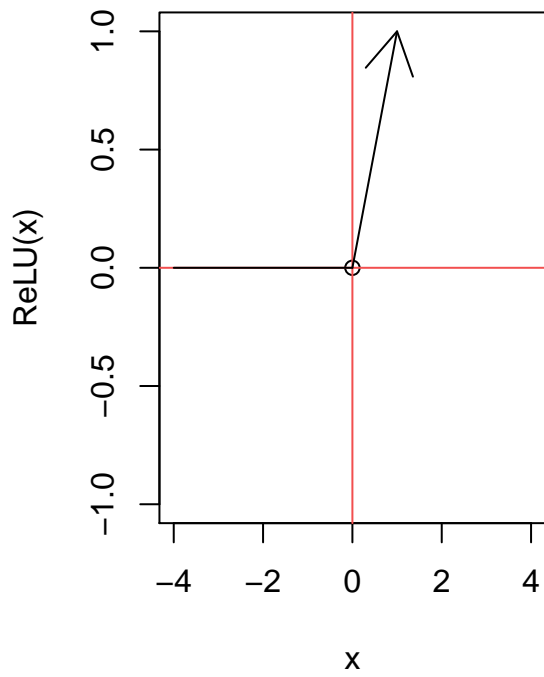
- The softmax activation function returns the probability distribution over mutually exclusive output classes.
- The calculated probabilities will be in the range of 0 to 1.
- The sum of all the probabilities is equals to 1.

Typically the Softmax Function is used in binary or multiple classification logistic regression models and in building the final output layer of NN.

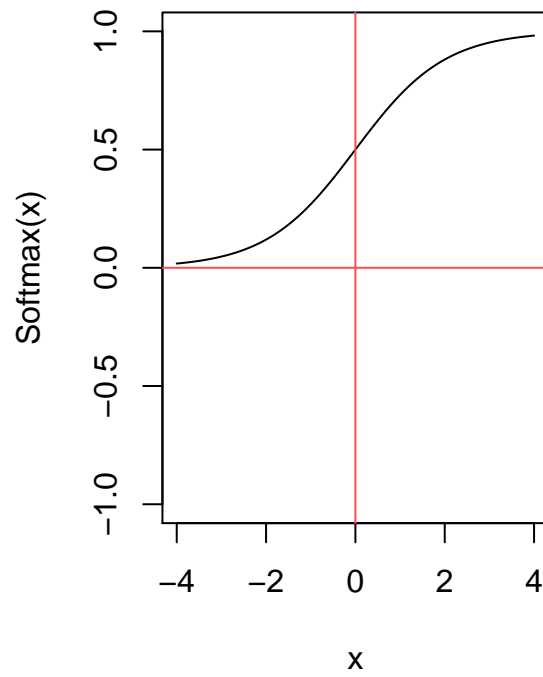
$$Z^{(6)} = Softmax(x) = \frac{e^{\alpha_i}}{\sum_{i=1}^n e^{\alpha_i}} \quad (7)$$

<sup>9</sup>Josh Patterson, Adam Gibson, Deep Learning; A Practitioner's Approach, 2017, O'Reilly

### ReLU Profile



### Softmax Profile



The benefit of these activation functions is that they are now differentiable. This fact becomes important for *Back-Propagation*, which is discussed later.

## The Two Neuron System

Building up in complexity, let us could consider our first Neural Network by using *only* two neurons. In two neuron systems, let us first generalize a bit more by adding that  $X$  is an array of all the inputs as is  $W_1$  and  $W_2$  is also an array of weights for each neuron. See figure #6.5.

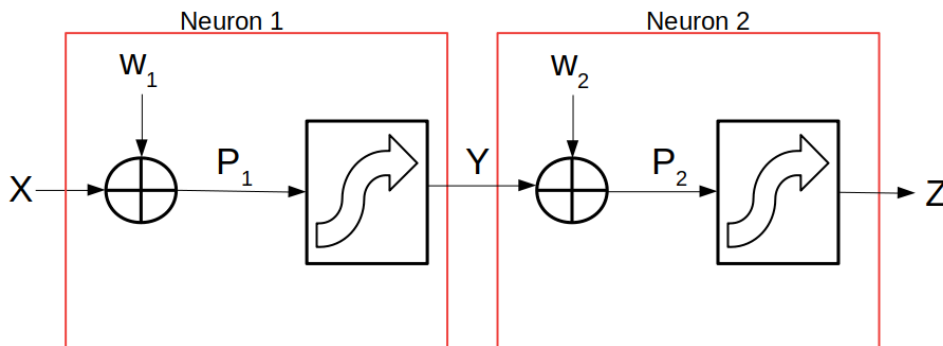


Figure 5: A Two Neuron System

### Feed-Forward In A Two Neuron Network

In our two neuron network, we can now write out the mathematics for each step as it progresses in a “forward” (left to right) direction.

Step #1: To move from  $X$  to  $P_1$

$$f^1(\vec{x}, \vec{w}) \equiv P_1 = (X^T \cdot W_1 - T) \quad (8)$$

Step #2:  $P_1$  feeds forward to  $Y$

$$f^2(P_1) \equiv \hat{Y} = \left( \frac{1}{1 + e^{-\alpha}} \right) : \text{ where } \alpha = P_1 \quad (9)$$

Step #3:  $Y$  feeds forward to  $P_2$

$$f^3(\vec{y}, \vec{w}) \equiv P_2 = (Y^T \cdot W_2 - T) \quad (10)$$

Step #4:  $P_2$  feeds forward to  $Z$

$$f^4(P_2) \equiv \hat{Z} = \left( \frac{1}{1 + e^{-\alpha}} \right) : \text{ where } \alpha = P_2 \quad (11)$$

Step #5: Our complicated function is simply a matter of chaining one result so that it may be used in the next step.

$$\hat{Z} = f^4(f^3(f^2(f^1(X, W)))) \quad (12)$$

In our **Feed-Forward Propagation**, we can now take the values from any numerical system and produce zeros, ones, or probabilities. Remember, in this set of experiments, we are using the concentrations of the 20 amino acids to provide a categorical or binary output, belongs to a) Myoglobin protein family, or b) does not.



## Error Back-propagation

Now that we have learned to calculate the output of our neurons using the Feed-Forward process, what if our final answer is incorrect? Can we build a feed back system to determine the weights needed to obtain our desired value of  $\hat{z}$ ? The short answer is yes. The process for determining the weights is known as Error Back-Propagation. Error Back-Propagation, also known as Back-Propagation, is crucial to understanding and tuning a neural network.

Simply stated Back-Propagation is an optimization routine which iteratively calculates the errors that occur at each stage of a neural network. Starting from randomly seeded values for the initial weights, Back-Propagation uses the partial derivatives of the feed forward functions. The chain rule and gradient descent are also used to determine the weights ( $W_1$  and  $W_2$ ) which are propagated through the network to find weights used in the summation step of a neuron.<sup>10</sup>

This thumbnail sketch gives the building blocks to calculate  $W$  which can be run until we reach a value that we desire. However the first time the back-propagation is carried out all the weights are chosen randomly. If the weights were set to the same number there would be no change throughout the system.

In the two neuron system, our first step is to generate an error or performance (Perf) function to minimize. If we call  $d$  our desired value, we can minimize the square error, a common choice.<sup>11</sup>

Step #1: Performance (Perf)

$$\mathbf{Perf} = c \cdot (d - \hat{z})^2 \quad (13)$$

Step #2:

$$\frac{dZ}{dx} = \frac{d\{f^4(f^3(f^2(f^1(X, W))))\}}{dx} \quad (14)$$

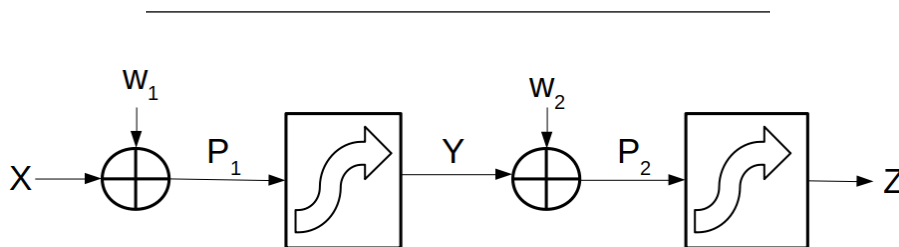


Figure 6: A Two Neuron System

Using the chain-rule, and figure 6.6, *Two Neuron System* as a guide, we can backwards to derive the formulas for error back-propagation. We find:

Step #3: Neuron 2  $\Rightarrow$  1

$$\frac{\delta Perf}{\delta w_1} = \frac{\delta Perf}{\delta z} \cdot \frac{\delta z}{\delta P_2} \cdot \frac{\delta P_2}{\delta y} \cdot \frac{\delta y}{\delta P_1} \cdot \frac{\delta P_2}{\delta w_1} \quad (15)$$

Step #4: Performance

$$\frac{\delta Perf}{\delta z} = \frac{\delta \left\{ \frac{1}{2} \| \vec{d} - \vec{z} \|^2 \right\}}{\delta z} = d - z \quad (16)$$

Step #5: Substitute  $P_2 = \alpha$

$$\frac{\delta z}{\delta P_2} = \frac{\delta ((1 + e^{-\alpha})^{-1})}{\delta \alpha} = e^{-\alpha} \cdot (1 + e^{-\alpha})^{-2} \quad (17)$$

<sup>10</sup>David Rumelhart, Geoffrey Hinton, & Ronald Williams, Learning Representations By Back-Propagating Errors, Nature, 323, 533-536, Oct. 9, 1986

<sup>11</sup>Ivan N. da Silva, Danilo H. Spatti, Rogerio A. Flauzino, Luisa H. B. Liboni, Silas F. dos Reis Alves, Artificial Neural Networks: A Practical Course, DOI 10.1007/978-3-319-43162-8, 2017

Step #6: Rearrange the right expression

$$\frac{e^{-\alpha}}{(1+e^{-\alpha})^{-2}} = \frac{e^{-\alpha}}{1+e^{-\alpha}} \cdot \frac{1}{1+e^{-\alpha}} \quad (18)$$

Step #7: Add 1 and subtract 1

$$= \frac{(1+e^{-\alpha})-1}{1+e^{-\alpha}} \cdot \frac{1}{1+e^{-\alpha}} \quad (19)$$

Step #8: Rearrange to find

$$= \left( \frac{1+e^{-\alpha}}{1+e^{-\alpha}} - \frac{1}{1+e^{-\alpha}} \right) \left( \frac{1}{1+e^{-\alpha}} \right) = \left( 1 - \frac{1}{1+e^{-\alpha}} \right) \left( \frac{1}{1+e^{-\alpha}} \right) \quad (20)$$

Step #9: Therefore we find

$$\frac{\delta z}{\delta \alpha} = \frac{\delta ((1+e^{-\alpha})^{-1})}{\delta \alpha} = \left( 1 - \frac{1}{1+e^{-\alpha}} \right) \left( \frac{1}{1+e^{-\alpha}} \right) \quad (21)$$

Nevertheless, we need one more part to ascertain the weights. As the error back-propagation is computed this process does not reveal how much the weights need to be adjusted/changed to compute the next round of weights given their current errors. For this we require one last equation or concept.

Once we compute the weights from our chain rule set of equations we must change the values in the direction proportional to the change in error. This is performed by using gradient descent.

Step #10: Learning Rate

$$\Delta W : W_{i+1} = W_i - \eta \cdot \frac{\delta Perf}{\delta W} \quad (22)$$

where  $\eta$  is the learning rate for the system. The key to the learning rate is that it must be sought and its range mapped for optimum efficiency. However smaller rates have the advantage of not overshooting the desired minimum/maximum. If the learning rate is too large the values of  $W$  may jump wildly and not settle into a max/min. There is a fine balance that must be considered such that the weights are not trapped in a local minimum and wildly oscillate unable to converge.

The last step of *error back-propagation* is simply setting up the derivatives mechanically and is not shown for brevity.

## Neural Network Experiment For Binary Classification

```
# Load Libraries
Libraries <- c("dplyr", "knitr", "readr", "caret", "MASS", "nnet", "purrr", "doMC")
for (p in Libraries) {
  library(p, character.only = TRUE)
}

# Load Data
c_m_TRANSFORMED <- read_csv("../00-data/02-aac_dpc_values/c_m_TRANSFORMED.csv",
                             col_types = cols(Class = col_factor(levels = c("0", "1")),
                                                PID = col_skip(),
                                                TotalAA = col_skip()))
```

```

# Create Training Data
set.seed(1000)
# Stratified sampling
TrainingDataIndex <- createDataPartition(c_m_TRANSFORMED$Class, p = 0.8, list = FALSE)

# Create Training Data
trainingData <- c_m_TRANSFORMED[ TrainingDataIndex, ]
testData      <- c_m_TRANSFORMED[-TrainingDataIndex, ]

TrainingParameters <- trainControl(method = "repeatedcv",
                                   number = 10,
                                   repeats = 5,
                                   savePredictions = "final") # Saves predictions

TuneSizeDecay <- expand.grid(size = c(16, 18, 20),
                             decay = c(1, 0.1, 0.01))

```

## Train model with neural networks

```
end_time - start_time          # Display time
```

```
## Time difference of 5.439274 mins
```

## Confusion Matrix and Statistics

```

NNPredictions <- predict(NNModel, testData)

# Create confusion matrix
cmNN <- confusionMatrix(NNPredictions, testData$Class)
print(cmNN)

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 239    9
##           1    4 215
##
##           Accuracy : 0.9722
##           95% CI : (0.9529, 0.9851)
##           No Information Rate : 0.5203
##           P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.9442
##
##           Mcnemar's Test P-Value : 0.2673
##
##           Sensitivity : 0.9835
##           Specificity : 0.9598

```

```

##          Pos Pred Value : 0.9637
##          Neg Pred Value : 0.9817
##          Prevalence     : 0.5203
##          Detection Rate : 0.5118
##          Detection Prevalence : 0.5310
##          Balanced Accuracy : 0.9717
##
##          'Positive' Class : 0
##

```

```

NNModel

```

```

## Neural Network
##
## 1873 samples
## 20 predictor
## 2 classes: '0', '1'
##
## Pre-processing: scaled (20), centered (20)
## Resampling: Cross-Validated (10 fold, repeated 5 times)
## Summary of sample sizes: 1685, 1686, 1686, 1686, 1685, ...
## Resampling results across tuning parameters:
##
##  size  decay  Accuracy  Kappa
##  16    0.01  0.9675458  0.9349866
##  16    0.10  0.9724570  0.9448152
##  16    1.00  0.9609233  0.9216202
##  18    0.01  0.9703226  0.9405495
##  18    0.10  0.9708545  0.9416022
##  18    1.00  0.9618830  0.9235428
##  20    0.01  0.9703157  0.9405366
##  20    0.10  0.9716003  0.9430995
##  20    1.00  0.9612419  0.9222614
##
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were size = 16 and decay = 0.1.

```

### Obtain List of False Positives & False Negatives

```

fp_fn_NNModel <- NNModel %>% pluck("pred") %>% dplyr::filter(obs != pred)

# Write/save .csv
write.table(fp_fn_NNModel,
            file = "./00-data/03-ml_results/fp_fn_NN.csv",
            row.names = FALSE,
            na = "",
            col.names = TRUE,
            sep = ",")

nrow(fp_fn_NNModel) ## NOTE: NOT UNIQUE NOR SORTED

```

```

## [1] 258

```

## False Positive & False Negative Neural Network set

```
keep <- "rowIndex"

fp_fn_NN <- read_csv("./00-data/03-ml_results/fp_fn_NN.csv")

NN_fp_fn_nums <- sort(unique(unlist(fp_fn_NN[, keep], use.names = FALSE)))

length(NN_fp_fn_nums)
```

```
## [1] 81
```

```
NN_fp_fn_nums
```

```
## [1] 4 6 15 16 46 57 94 97 100 114 115 116 130 136 149
## [16] 150 170 179 182 183 185 249 445 449 453 503 518 522 526 530
## [31] 531 532 534 546 547 566 570 580 592 655 910 913 980 1033 1034
## [46] 1035 1093 1094 1100 1101 1117 1121 1130 1190 1219 1226 1233 1264 1300 1471
## [61] 1510 1522 1575 1576 1579 1585 1587 1594 1608 1618 1621 1693 1697 1734 1771
## [76] 1773 1780 1789 1831 1833 1873
```

```
write_csv(x = as.data.frame(NN_fp_fn_nums),
          path = "./00-data/04-sort_unique_outliers/NN_nums.csv")
```

## Neural Network Conclusion

The Neural Network set included a total of 79 unique observations containing both FP and FN.

Accuracy was the primary criteria used to select the optimal model. There were 20 models tested. The neural network was configured with 20 inputs (one for each amino acid), one hidden layer and one output. The hidden layer was tested with either 10, 12, 14, 16, 18, 20 neurons and a array of different decays (1, 0.1, 0.01, 0.001). The caret software has a tuning parameter named `tuneGrid` that allows users to **expand** a set of arrays to a matrix of combinations to be tested. Therefore 20 models were tested with the training data set and the best values were size = 20 and decay = 0.1.

At this time, the author is not aware of any heuristic that gives the proper number of hidden layers and the proper number of neurons in each layer therefore one must search the experimental space for an optimized configuration. If a more thorough search of the experiemntal space was carried out using two or three hidden layers would be investigated. The poor showing of the Neural Network suggests that the data may have some additional decision boundary that is not yet represented by only 20 neurons.