

# 1.1 Softmax and Cross Entropy

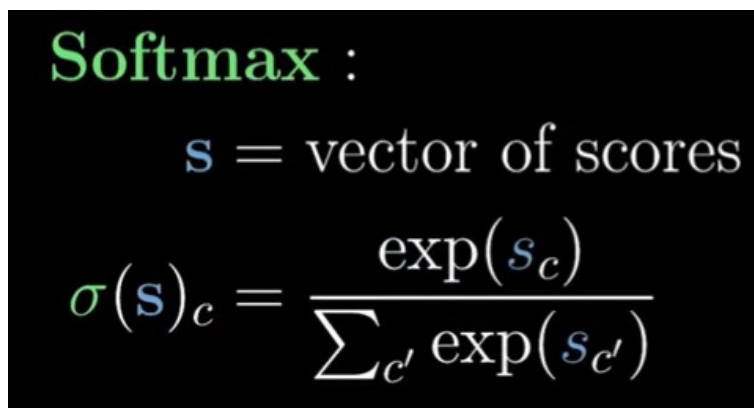
## Summary

Why artificial neural network works?.....	1
Artificial neural network structure.....	1
How neural network predicts: forward propagation.....	5
How neural network learns: back propagation.....	7
Cost/error function.....	8
Gradient vector.....	9
Math behind network learning.....	11
Using the matrix's power to calculate the gradient vector.....	16
Bibliography.....	18

## What is Softmax?

Softmax is an special activation function only used at the output layer of a neural network, which allows classification neural networks to assign a “probability” value to each class, a decimal value between 0 and 1. The sum of all values in output layer is always equal to 1.

This is the formula:



**Softmax :**

$s$  = vector of scores

$$\sigma(s)_c = \frac{\exp(s_c)}{\sum_{c'} \exp(s_{c'})}$$

S is the output layer, c is the index of the node in that output layer, and exp(x) equals  $f(x) = e^x$ . For each value in every neuron in output layer, we calculate the exponentiation of that value divided by the sum of the exponentiation of all values in the output layer.

Why exponentiation is used instead of another power? It has many advantages like guarantee the sum of errors will never be below 0. Here is more explained: <https://www.youtube.com/watch?v=p-6wUOXaVqs>

# What is Cross Entropy?

Cross Entropy is the used cost and loss function when training a more likely classification neural network. With softmax function, it can be simplified.

**Both Softmax and this cross entropy only works in multi-class only-1-labeled problems. It cannot be used in multi-label problems because output layer “probability” values are not independent each other.**

The cost function expression is basically (where not-hat-y is the true expected value, and hat-y is the predicted output value):

$$L = - \sum_i y_i \log(\hat{y}_i)$$

The derivative of the cost respect the output value of each output neural network it's basically:

$$\hat{y}_i - y_i$$

Notice always the true expected value is 1 or 0, while the predicted output value is a “probability” value

## Bibliography

**Page with cost function:** <https://levelup.gitconnected.com/killer-combo-softmax-and-cross-entropy-5907442f60ba>

**Where derivatives come from(15:45) + crossentropy+softmax in action(17:00):**  
<https://www.youtube.com/watch?v=xBEh66V9gZo>

**Deeper derivatives for cross entropy+softmax:** <https://www.youtube.com/watch?v=M59JElEPgIg&pp=ygUdc3RhdHF1ZXN0IHNVZnRtYXggZGVyaXZhdGl2ZXNM%3D>