

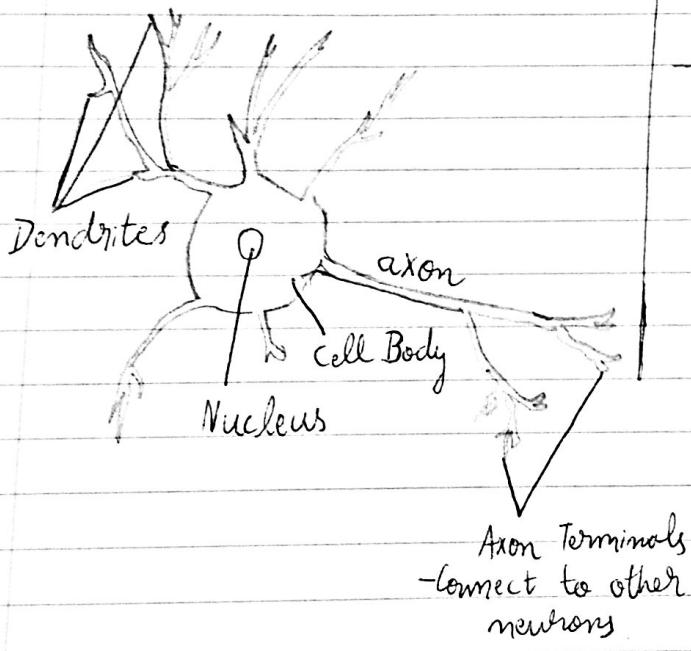
Exaggerated!  
Finally!  
Magic and superpowers!

# Neural Networks!

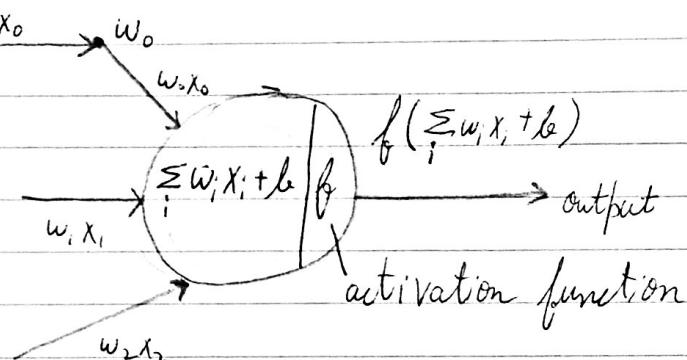
\* Spoiler: It's not actually a network, just a series of operations. False advertising.

\* Neuron: A neuron is a computational model of a function that accepts inputs and "fires" an output depending on some measure of strength of the inputs.

Biological Model



Computational Model



\* Activation Function: Computes the output / "strength" of the sum of inputs (weighted inputs). This is ~~the~~ the core element of Neural Nets. More on this later.

\* So essentially, Neuron (also called "Perception" or "Artificial Neuron") is a model of a function that accepts weighted inputs.

\* Depending on Activation function, a Neuron can "act" as other computational machinery (?)

For example, choosing it to be sigmoid results in:

$\sigma(\sum_j w_j x_j + b) \rightarrow \text{Binary Softmax classifier}$

$$= P(y_i = 1 | x_i; w)$$

Use hinge loss function (max-margin)  $\rightarrow \text{Binary SVM}$

\* Activation Functions in Neural Nets:

i) Sigmoid:  $\sigma(x) = \frac{1}{1 + e^{-x}}$   $\rightarrow$  Squashes values to between 0 & 1.  
Range: [0, 1]

Disadvantages:

i) Saturates & kills gradients: Gradient is 0 at both 0 and 1 tails of sigmoid (because value becomes constant at 0 and 1 for large values of input, so differential is 0.) So gradient signal is 0 and network learns nothing.

ii) Sigmoid output not zero-centered: Since sigmoid is always greater than 0, all gradients are either positive or all negative

2. Tanh: Hyperbolic Tangent. Range: [-1, +1]  
 $\tanh(x) = 2\sigma(2x) - 1$

3. ReLU (Rectified Linear Unit):  $f(x) = \max(0, x)$   
 → Simple thresholding, but shown to be 6x better for convergence of SGD.  
 → However, they are fragile. Can never be activated and can "die". ~~lot of network can die~~.

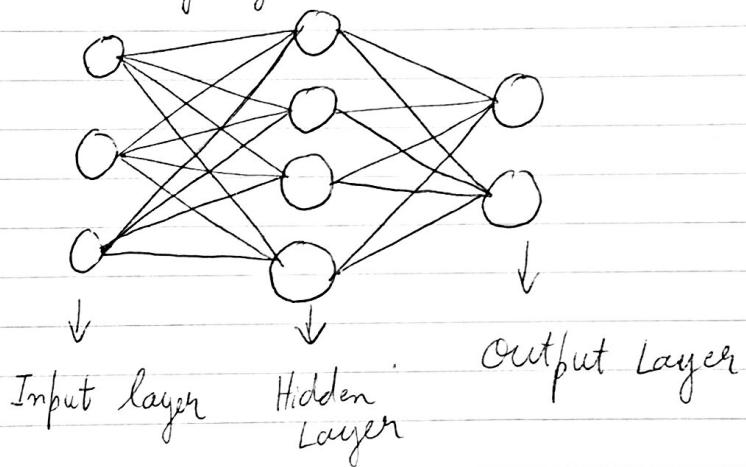
if learning rate is  
too high

→ As much as 40% of network can be dead, so it is important to set ~~the~~ learning rate.

ReLU is most widely used in these modern times.

## ★ Neural Network Architectures

### Layer-wise Organization:



- Input Layer: This is the vector of inputs.
- Output Layer: The output vector of the network.
- Hidden Layers: This is the layer that contains Neurons and does the computation. When we say  $n$ -layer neural net,  $n$  is the no. of Hidden layers. These are called "Hidden" because they detect the "hidden" features not provided / visible in the training data.

## \* Feed-Forward Computation:

A Neural Network computation is Repeated matrix multiplications interwoven with activation function. Example:

Inputs are  $3 \times 1$  vectors (3 nos.), Outputs are  $k$ -vectors containing scores for each of the  $k$ -classes.

$$\begin{array}{l} \text{Input : } \\ \text{Layer } 0 \end{array} \quad \text{im} = \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix}$$

$$\begin{array}{l} \text{Hidden : } \\ \text{Layer 1} \end{array} \quad \text{W1} = \begin{bmatrix} -1 & 2 & 9 \\ 0 & 1 & 6 \\ 11 & 21 & 4 \\ 6 & 69 & 83 \end{bmatrix} \quad h_1 = f(W1, \text{im}) = \begin{bmatrix} * \\ * \\ * \\ * \end{bmatrix} \quad (4 \times 1)$$

$$\begin{array}{l} \text{Hidden : } \\ \text{Layer 2} \end{array} \quad \text{W2} = \begin{bmatrix} * & * & * & * \\ : & : & : & : \\ . & . & . & . \\ ; & ; & ; & ; \end{bmatrix} \quad h_2 = f(W2, h_1) = \begin{bmatrix} * \\ * \\ * \\ * \end{bmatrix} \quad (4 \times 1)$$

$$\begin{array}{l} \text{Output : } \\ \text{Layer } 3 \end{array} \quad \text{out} = W_3 = ((\text{no. of outputs}) \times 4) \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \quad W_3 \text{ is } (y \times 4) \text{ where } y \text{ is desired outputs, we can use later.}$$

$$\text{out} = W_3 \cdot h_2 = \begin{bmatrix} * \\ * \\ * \\ * \end{bmatrix} \quad (y \times 1)$$

- $W_1, W_2, W_3, b_1, b_2, b_3$  are all learnable parameters.
- Input neurons = ~~width~~ rows in input vector.  
Input vector could hold all (batch of training data, where each column is an input example.)

### Representational Power:

- Neural Networks with fully connected layers can be seen as representing a family of functions parameterized by the weights.
- ★ Neural Nets can model/represent these functions because a Neural Net with at least one hidden layer is a Universal Approximator capable of approximating any real valued function, given we choose an appropriate non-linearity as the activation function. (See paper by G. Cybenko, 1989, for details.)

This is why Neural Nets can learn a huge class of functions and approximate map new, hidden features.  
(Keep in mind much of ML deals with approximating/mining the function or the process that generated the data. We learn this function and use it for predicting new information.)

### About Layers, Neurons and Capacity:

- No. of layers = Depth      No. of Neurons = Width
- Although more layers does not mean better representation power. One hidden layer - necessary, but more is not better (in general) except in case of Convolutional Nets.
- Layers + Neurons = Capacity. High capacity = Overfitting, less noise in the data. Low capacity = Underfitting.

• However, to minimize Overfitting, use Regularization, do not reduce capacity. It is easier to train larger networks. Larger networks contain many more local minima and is easier to for loss function to converge to.

## Regularization:

1. L2 Regularization:  $\text{Loss}_{\text{penalized}} = \text{Loss} + \frac{1}{2} \lambda w^2$ , for every weight in the network. Like before, the squared-magnitude penalizes large weights. The  $\frac{1}{2}$  is there because when we take its differential in gradient, it ends up as  $\lambda w$ . ( $\lambda$  - regularization strength - hyperparameter)

2. Dropout: While training, keep a neuron alive with probability  $p$ , or else set its output to 0.

$$H_1 = f(W_1 \cdot X + b)$$

if ( $p = \text{randomNumber}() < 0.5$ ):

$$H_1 = 0$$

$$H_2 = \dots$$

Note: In L2 reg, we add the regularization gradient to the gradient w.r.t to the weights as well.

# How To Train Your Neural Network:

(Coming to ~~Movies~~ a Cinema near you, Winter 2017.)

1. Eat Initialize weights, hidden layers using these weights, hyper parameters.

2. Gradient Descent start. Loop:

3. Evaluate class scores by performing feed-forward computation, as shown previously.

4. Compute class probabilities using softmax over scores  
(optional, but better)

5. Compute Loss, using loss function + Regularization.

6. Compute gradient on scores. For example:

$$L_i = -\log \left( \frac{e^{f_k}}{\sum_j e^{f_j}} \right) \quad \begin{array}{l} \text{[Cross-Entropy} \\ \text{Loss]} \end{array}$$

$$\frac{dL_i}{df_k} = \left( \frac{e^{f_k}}{\sum_j e^{f_j}} \right) - 1(y_i = k) \quad \begin{array}{l} [-1 \text{ for all scores if} \\ \text{probs of correct class}] \end{array}$$

Python:

$$\text{probs} = [p_{y_1}, p_{y_2}, \dots, p_{y_k}]$$

$$d \text{ scores} = \text{probs}$$

$$d \text{ scores} [\text{range}(0, \text{num-examples}), y] = 1$$

7. Backpropagate gradient to weights, biases (if any)

8. Add regularization to gradient  $[dW_i = dW_i + \lambda w_i]$

9. Update parameters

$$\boxed{\begin{array}{l} w_i = w_i + dW_i \\ b_i = b_i + dB_i \end{array}}$$