

Introduction to Neural Networks

Introduction (From Tosh)

ANN are mlc learning techniques which simulate the mech of learning biological organisms.

The human nervous system contains cells, referred as neurons.

The Biological neuron

Dendrite \rightarrow input ; Axon - output

Neurons Pass electrochemical impulses across synapses from one cell to the next.

Impulse must be strong enough to activate the next neuron by release of chemicals across a synaptic bifurcation.

The neuron is made up of a nerve cell consisting of a soma that has many dendrites but only one axon.

Axons are special nerve fibres with a special cellular extension that comes from the cell body.

The single axon can branch hundreds of times, however.

Dendrites are thin structures that arise from the main cell body.

Synapses are connecting junction b/w neurons and dendrites. Send signals from the axon of a neuron to the dendrite of another neuron.

Dendrites have fibers branching out from the soma in a bushy network around the cell.

Dendrites allow the cell to receive signals from connected neighbouring neurons. It performs multiplication by that dendrite's weight value. Multiplication means increase or decrease in the ratio of synaptic neurotransmitters to signal chemicals introduced into the dendrites.

Axons are single, long fibers extending from the main soma. They stretch out longer distance than dendrites and measure generally 1m in length. They will branch 100 of times and connect to other dendrites. Neurons are able to send electrochemical signal that travels along the cell's axon and activates synaptic connections with other neurons.

Information flow across the biological neuron

Synapses that increase the potential and those that decrease the potential.

Neurons also have been shown to form new connections over time and even migrate. This combined mechanism of connection

Total connection in the human brain

Biological neural networks are composed of 86 billion neurons connected to many other neurons 500 billion connections between neurons in the human brain.

Activation function: Tells what should be the output
↳ Various functions are there: Sigmoid, ReLU (Rectified Linear Unit)

| e.g. Better understanding what cat is!
if $x=1 \rightarrow$ respond

DATE: Step function = Activation function

Humans versus computers: stretching the limits of AI (chart)

Success of neural network: increased data availability and computational power of modern computers has outgrown the limits of traditional machine learning.

Smaller data sets = Machine learning is better.

Reasons: more choice, greater ease of model interpretation, tendency to handcraft interpretable features that incorporate domain specific insights.

Neural Networks: (Josh)

NN are computational model that shares some properties with the animal brain in which many simple units are working in parallel with no centralized control unit. Behaviour is shaped by its network architecture (NA). NA can be defined by the following

(i) Number of neurons (ii) No. of layers (iii) Types of connections between layers.

From biological to artificial

Animal brain is fundamental components of the mind. Do not completely understand how this collection of decentralized functional units provides the foundation for thought and the seat of consciousness.

Modeling the neuron with beginning of perception.

Perception

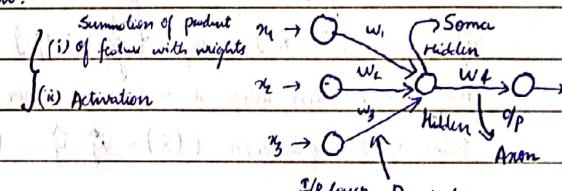
Linear model used for binary classification. Neural networks consider perception as artificial neuron.

Preprocessing: Multiplication of input with weight.

$$Y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}$$

$$Z = \text{Act}(Y)$$

$$= Z * W_f$$



Definition of the perceptron

The perceptron is a linear model binary classifier with a simple input-output relationship as, it shows we're summing no. of inputs times their associated weights and then sending this "net input" to a step function with a defined threshold. Typically with perceptrons, this is a step function with a threshold value of 0.5. This function will output a single binary value (0 or 1), depending on input.

$$f(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$$

Net Activation function = dot product of the input and connection weights.

If the bias value is negative, it forces the learned weights sum to be much greater value to get a 1 classification output. Bias term is learned through perceptron learning algorithm.

The perceptron learning algorithm

The perceptron learning algorithm takes web input word, as we can see in Fig., and computes the output classification to check against the actual classification label.

CLASSMATE

PAGE: [] []

Hidden layer - Rule
Output layer - Anti sigmoid

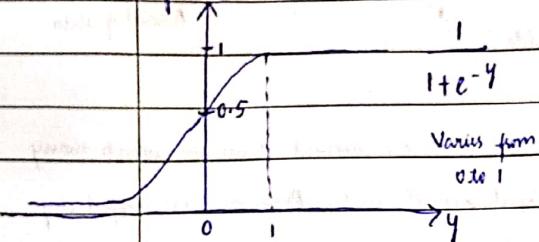
y = Summation of weight and inputs
 \hat{y} = Predicted

DATE

- To produce the classification, the columns are matched up to weights where m is the number of dimensions in both our input and output weights. Bias value = 1 because we don't affect the bias input.

Sigmoid function:

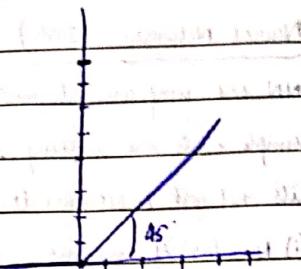
Sigmoid(y)



Rule: max($y, 0$)

$$\downarrow \text{Rectified } y = -ve \max(-ve, 0) = 0$$

$$\text{Linear unit } y = +ve \max(+ve, 0) = +ve$$



The Basic Architecture of Neural Networks (Chanc C. Aggarwal)

In any type of application, one would have historical cases in which the class variable is observed and other cases in which the class variable has not yet been observed but needs to be predicted.

The input layer does not perform any computation on its own.

The predicted \hat{y} is computed as follows.

$$\hat{y} = \text{sign}\{\sum_j w_j x_j\} = \text{sign}\{\sum_j w_j x_j\}$$

The sign function maps a real value either +1 or -1, which is appropriate for binary classification.

Error of the prediction $E(x) = y - \hat{y}$ (keep it minimum as possible).

If error value is good enough, the weights in the neural network need to be updated in the direction of the gradient.

Perception contains single computational layer; it is considered a single layer network.

Invariant part of the prediction (some entity that is not altered) which is referred to as bias.

$$\hat{y} = \text{sign}\{\sum_j w_j x_j + b\}$$

The bias can be incorporated as the weight of an edge by using a bias neuron.

Invariant: quantity that is unaltered.

Bias (Important): a neuron improves properties of the neuron. It allows moving the threshold of activation function.

Bias value allows shifted left or right to better fit data, it is always 1.

Weights represent the connection between units. Weights of NN can and will become positive or negative depending on training data.

Training: overall algo used to train NN to recognize a certain I/P and map it to O/P.

Involves finding a set of weights in the network that proves to be good at solving specific problem.

Learning: deals with individual weights and thresholds and decides how those would be manipulated.

$\eta \rightarrow$ learning rate
values from 0.01 - 0.0001,
normally we take 0.001

$O_{11} \rightarrow$ output

$W_{11}^1 \rightarrow$ hidden layer 1
 $O_{11} \rightarrow$ hidden layer node
 \downarrow
Input node

$f_{11} \rightarrow$ Neuron
 \downarrow
Mulen

$$\text{Cost} = \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

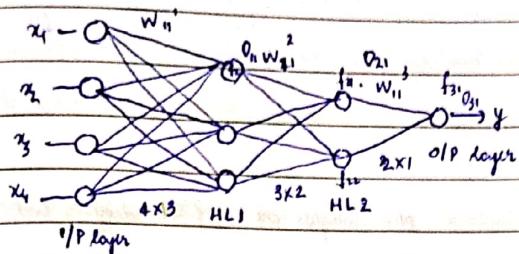
DATE

Loss function: $(y - \hat{y})^2$ to reduce that we use optimizers. Optimizers will undergo backpropagation.

$$\text{where } \Delta W_{111} = \eta \frac{\partial L}{\partial W_{111}}$$

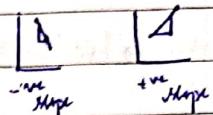
1 iteration = 1 forward propagation + 1 backward propagation
referred to as epoch.

Multilayer Neural Networks



Step 1: $\sum_{i=1}^n W_{ii} x_i + \text{bias}$

Step 2: $\text{Act}(z)$

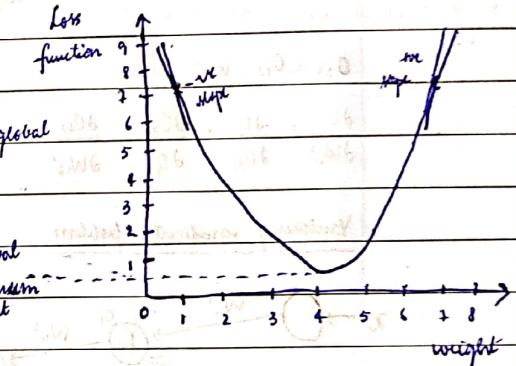


Matrix is formed w.r.t weight and next layer.

Gradient descent (optimizer) : Reduces the $(y - \hat{y})^2$; \rightarrow Gradient descent stochastic Gradient: picks one feature randomly and performs F/W and B/W propagation

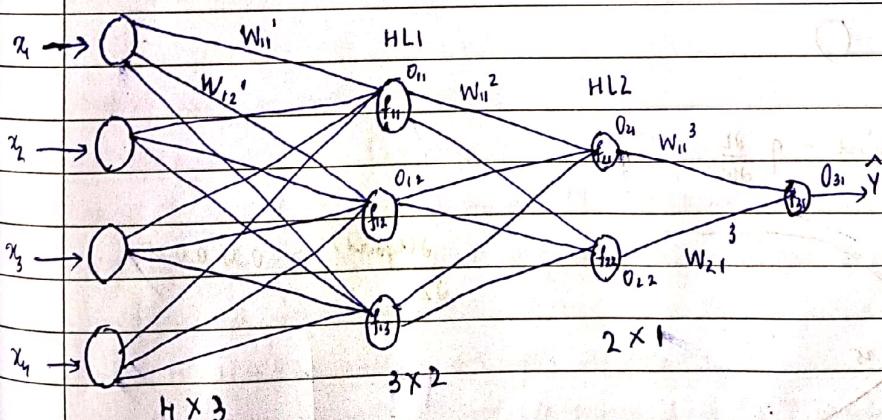
To update the weight such a way that the \hat{y} is to be similar to y . To find out the minimum efficient weights do we updated to get just \hat{y} .

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W} \quad \text{done by output node}$$



O_{31} is computed by
Act. function intern

inputted with weights.



$$W_{11}^{3 \text{ new}} = W_{11}^{3 \text{ old}} - \eta \frac{\partial L}{\partial W_{11}}$$

$$O_{31} = W_{11}^3 * O_{11} + W_{21}^3 * O_{21}$$

PAGE

Anaconda is a free and open-source programming Python distribution and collection of hundreds of packages related to data science, scientific programming, development and more.

$$\frac{\partial L}{\partial W_{ii}^3} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial W_{ii}^3}$$

From chain rule of derivative

$$W_{21}^3_{new} = W_{21}^3_{old} - \eta \frac{\frac{\partial L}{\partial O_{31}}}{\frac{\partial O_{31}}{\partial W_{ii}^3}}$$

$$\frac{\partial L}{\partial W_{ii}^2} = \frac{\partial L}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial W_{ii}^2}$$

$$W_{12}^2 \rightarrow W_{12}^2_{new} = W_{12}^2_{old} - \eta \frac{\frac{\partial L}{\partial O_{22}}}{\frac{\partial O_{22}}{\partial W_{ii}^2}}$$

New weights are computed during back propagation by using these computation

$$O_{21} = W_{ii}^2 \cdot O_{ii}$$

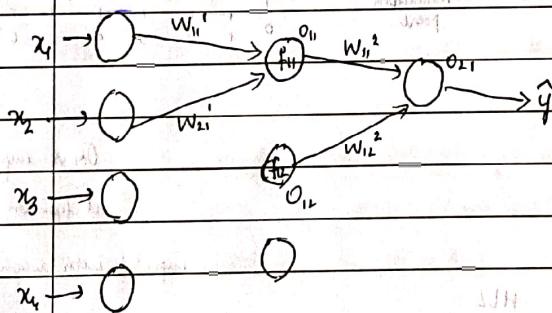
$$\frac{\partial L}{\partial W_{ii}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{ii}^2}$$

$$W_{12}^2 \rightarrow W_{12}^2_{new} = W_{12}^2_{old} - \eta \frac{\frac{\partial L}{\partial O_{31}}}{\frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{ii}^2}}$$

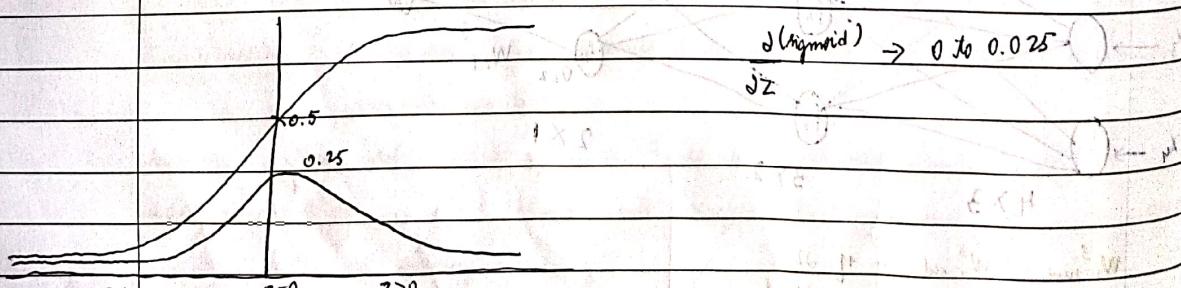
$$O_{22} = O_{12} \cdot W_{12}^2$$

$$\frac{\partial L}{\partial W_{ii}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial W_{ii}^2}$$

Vanishing gradient problem



$$W_{ii}^3_{new} = W_{ii}^3_{old} - \eta \frac{\frac{\partial L}{\partial O_{21}}}{\frac{\partial O_{21}}{\partial W_{ii}^3}}$$



loss function value is less implies y and \hat{y} are almost same and prediction is correct.

Vanishing gradient is because of sigmoid derivative of ReLU is 1, hence the gradient will not vanish. Hence we will not get any problem.

DATE []

If more number of layers are there both exploding and vanishing will occur together.

Only for sigmoid function loss is found.

$$W_{ii}^{new} = W_{ii}^{old} - \eta \frac{\partial L}{\partial W_{ii}}$$

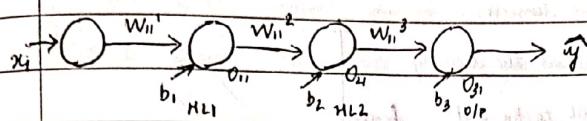
$$\text{loss} = (y - \hat{y})^2$$

$$\frac{\partial L}{\partial W_{ii}} = \frac{\partial L}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial W_{ii}} \stackrel{0.25}{=} 0.25$$

(Activation function) Derivative has caused the vanishing of gradient descent.

To avoid this vanishing gradient we came up with ReLU activation function.

Exploding gradient problem (Because of initialization of weight and sigmoid function) Here we will never go to global minima but we will be moving between positive and negative steps.



$$W_{ii}^{new} = W_{ii}^{old} - \eta \frac{\partial L}{\partial W_{ii}^{old}}$$

$$\frac{\partial L}{\partial W_{ii}} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial W_{ii}}$$

b is bias, and it is there for all the neurons.

$$O_{21} = \text{Act}(z)$$

$$z = W_{ii}^2 O_{11} + b_2$$

$$\frac{\partial O_{21}}{\partial O_{11}} = \frac{\partial [\text{Act}(z)]}{\partial O_{11}}$$

$$\frac{\partial O_{11}}{\partial O_{11}} = 1$$

$$\frac{\partial [\text{Act}(z)]}{\partial z} \times \frac{\partial z}{\partial O_{11}} \Rightarrow 0 \text{ to } 0.25 \times \frac{\partial [W_{ii}^2 O_{11} + b_2]}{\partial O_{11}}$$

$$0 \text{ to } 0.25 * [W_{ii}^2 + 0]$$

$$\Rightarrow 0 \text{ to } 0.25 * W_{ii}^2$$

Droput Layers in Multi Neural Network

Underfitting in NN:

One of the problems that occur during neural network training is underfitting. Underfitting is the case where the model has "not learned enough" from the training data, resulting in low generalization and unreliable predictions.

A model with too little capacity cannot learn the problem.

Too little learning will make the model to perform poorly on the training dataset and on new data.

An underfit model has high bias and low variance.

We developed some features in overfitting problem. We decide using dropout ratio, it is selection and distribution of neurons are random. In BP, updates only active neuron's weight. Weight is reduced to not being overfit.

DATE []

Overfitting in NN:

One of the problems that occur during the neural network's training is called overfitting. The error on the training set is driven to a very small value, but when new data is presented to the network the error is large.

The network has memorized the training examples, but it has not learned to generalize to new situations.

A model with too much capacity can learn it too well and overfit the training dataset. With too much learning, the model will perform well on the training dataset and poorly on new data.

An overfit model has low bias and high variance. The model learns the training data too well and performance varies widely with new unseen.

Good fit model: A model that suitably learns the training dataset and generalizes well to the old test dataset.

How to handle underfitting problem?

Underfitting can easily be addressed by increasing the capacity of the network.

How to handle overfitting problem?

(i) Reduce the network's capacity by removing layers or reducing the number of neuron elements in the hidden layers

(ii) Use dropout layers which will randomly remove certain features by setting them to zero.

Dropout layers are an easy and effective way to prevent overfitting in NN.

A dropout layer randomly drops some of the connections between layers. This helps to prevent overfitting. We reduce the weights of the model to small during learning.

These techniques not only reduce overfitting, but they can also lead to faster optimization of the model and better overall performance.

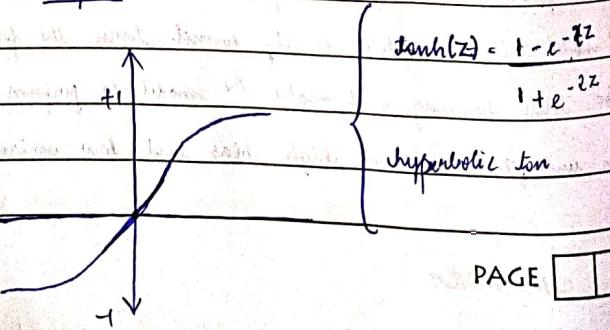
In deep NN, there are many weights and bias parameters. When we have large weights and bias parameters, the NN model tends to overfit the data set problem on a particular use case.

In deep NN, underfit will never happen. Usually underfitting happens with one layer NN.

To fix this problem technique used called dropout.

derivative will vary from 0 to 1.

CLASSMATE



hyperbolic tan

PAGE []

If we use sigmoid and tanh activation function, gradient descent (may or may not) converge and vanishing gradient problem occurs. (Disadvantages)

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial W}$$

We find negligible difference in W_{new} and W_{old} when sigmoid and tanh activation function are used.

Gradient descent may or may not converge when sigmoid and tanh AFs are used. It may converge when no. of iterations are performed for many times.

This problem does not arises in ReLU function, so it is used in deep neural network.

In ReLU case, gradient descent will not converge when $z=0$ i.e. $W_{new} = W_{old}$, so you need to ensure that when $z \leq 0$, you need to have some value great apart from zero.

This is lucky value when $z=0, z \neq 0$

How to get this component?

$$\begin{aligned} &= (\text{ReLU}) + 0.01 \\ &= -0.01 \end{aligned}$$

To avoid the use of dead neuron in the network, we make use of lucky neuron. Lucky neuron need not be applied when less number of neurons are used.

With lucky ReLU vanishing gradient problem never arises.

Weight Initialization techniques

Key points while initializing the weights

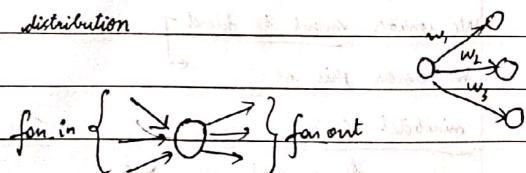
- (i) weights should be small (but must not be very very small)
- (ii) weights should not be same
- (iii) weight should have good variance

uniform distribution technique

weights are taken and sampled from uniform distribution

It should be within a range

will be
$$w = \frac{t - b}{\sqrt{fan_in}} + b$$



Sigmoid activation function works better for this.

In Minibatch SGD the error is somewhat smaller than GD.
w.r.t same data (first)
GD > Minibatch > SGD

DATE []

2. Xavier / Glorot distribution:

In this we have two cases:
(i) Xavier / Glorot uniform distribution
(select the smaller value)

$$W_i \text{ will range from } -\frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}} \text{ to } \frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}$$

(ii) Xavier / Glorot Normal:

We will range from Normal ($\mu = 0$ and $\sigma = \text{standard deviation}$)

$$\sigma = \sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$$

Sigmoid activation function works well.

3. He initialization

(i) He uniform:

$$W_i \text{ will be in range } -\frac{\sqrt{6}}{\sqrt{\text{fan-in}}} \text{ to } \frac{\sqrt{6}}{\sqrt{\text{fan-in}}}$$

(ii) He normal:

Normal ($\mu = 0$, std deviation = σ)

$$\sigma = \sqrt{\frac{2}{\text{fan-in}}}$$

ReLU activation function works well.

Stochastic Gradient:

Only one i/p feature is randomly

picked out of 'n' features, it is

forward propagated and then backward

propagation is performed.

Error = $(y - \hat{y})$; Loss = $(y - \hat{y})^2$

Cost = $(y - \hat{y})^2$

Gradient Descent

Error = $(y - \hat{y})$

Loss = $(y - \hat{y})^2$

Cost = $\sum_{i=1}^n (y - \hat{y})^2$

Accuracy is less than SGD.

Training is done more quickly

we consider full dataset

Benefit: Accuracy is more in this model

Training takes more time.

We consider subset of dataset

We consider this as

minibatch SGD.

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 \text{ where } n \geq k$$

For (GD):

To compute large data, we need to store all GD in the RAM. Computational power required is also large.

More amount of time is required to compute \hat{y} .

CLASSMATE

For (Minibatch):

To handle this less amount of memory is enough and also less computation is enough. Only the required row

activates. After sometime SGD overtakes GD \Rightarrow SGD minibatch is faster than GD.

PAGE []