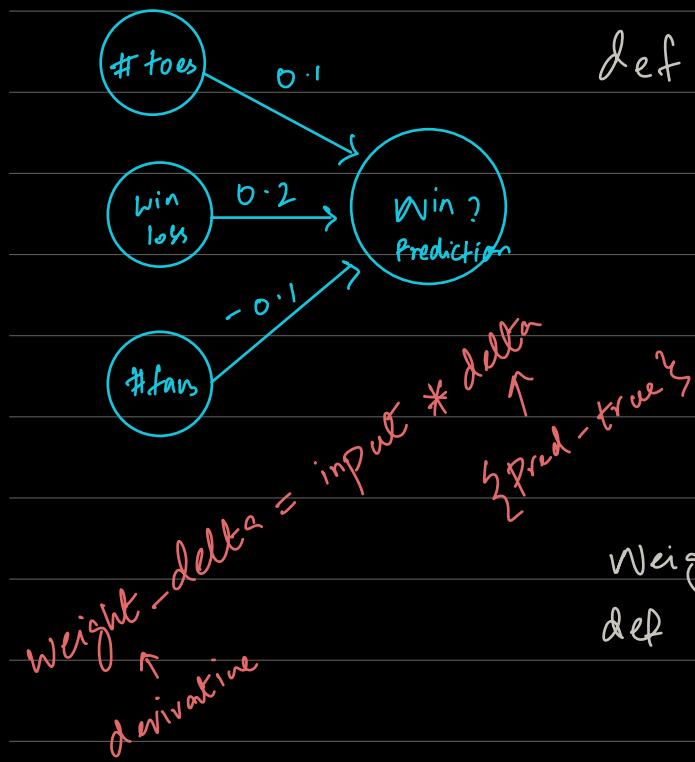


Gradient descent learning with multiple inputs:

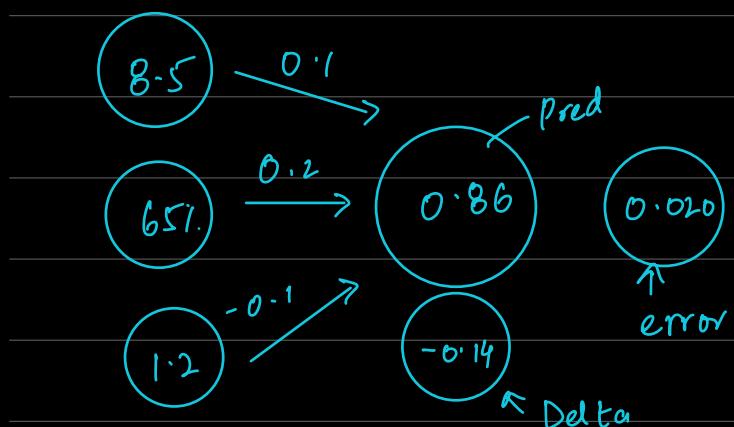
- a) Using gradient descent to update a weight.
- b) Using same technique to update a network containing multiple weights.



```
def w_sum(a, b):
    assert (len(a) == len(b))
    out put = 0
    for i in range (len(a)):
        out put += (a[i] * b[i])
    return out put
```

$$\text{Weights} = [0.1, 0.2, -0.1]$$

```
def neural_network(input, weights):
    pred = w_sum(input, weights)
    return pred
```



$$\text{toes} = [8.5, 9.5, 9.9, 9.0]$$

$$\text{wires} = [0.65, 0.8, 0.8, 0.9]$$

$$\text{nfans} = [1.2, 1.3, 0.5, 1.0]$$

win_or_lose_Binary = [1, 1, 0, 1]

true = win_or_lose_Binary[0]

input = [toes[0], nrec[0], nflw[0]]

pred = neural_network(input, weights)

error = (pred - true) ** 2

delta = pred - true.

-1.19 ↓ error goes down ↑
↑ if weight goes up
↓ unit

def ele_mul(number, vector):

output = [0, 0, 0]

assert len(output) == len(vector)

for i in range(len(vector)):

output[i] = number * vector[i]

return output

* Each weight-delta

pred = neural_network(input, weight)

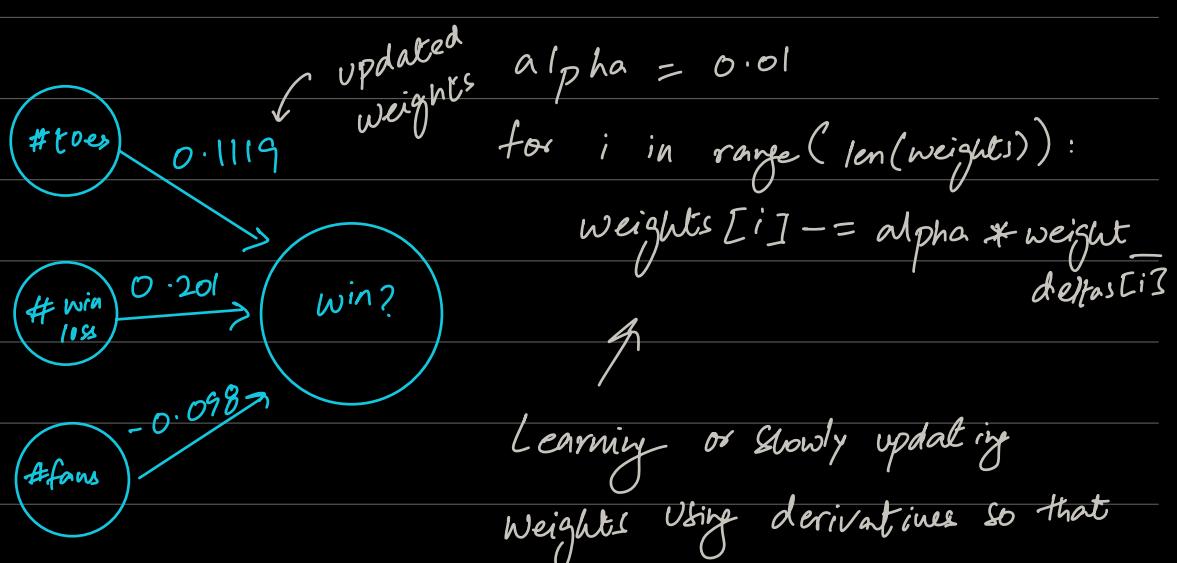
is calculated by taking output error = (pred - true) ** 2

delta and multiplying it delta = pred - true

by its input.

weight_deltas = elem_mul(delta, input)

Learn



Old weights weight -
 \downarrow \downarrow error reduces to 0.

$$0.1 - (-0.119 * 0.01) = 0.1119 = \text{weights}[0]$$

$$0.2 - (-0.091 * 0.01) = 0.2009 = \text{weights}[1]$$

$$-0.1 - (-0.168 * 0.01) = -0.098 = \text{weights}[2]$$

Pred - True

*) Delta: A measure of how much higher or lower you want a node's value to be, to predict perfectly given the current training example.

input * Delta - to calculate for each input

*) Weight - delta: is an estimate of the direction

and amount you should move a weight

to reduce node - delta, inferred by derivative.

→ make it close to zero.

A derivative-based estimate of the direction and amount you should move a weight to reduce the node-delta, accounting for scaling, negative reversal, stopping.

Weight-delta \rightarrow derivative based estimate of $\frac{dy}{dx}$ \rightarrow you should move a weight

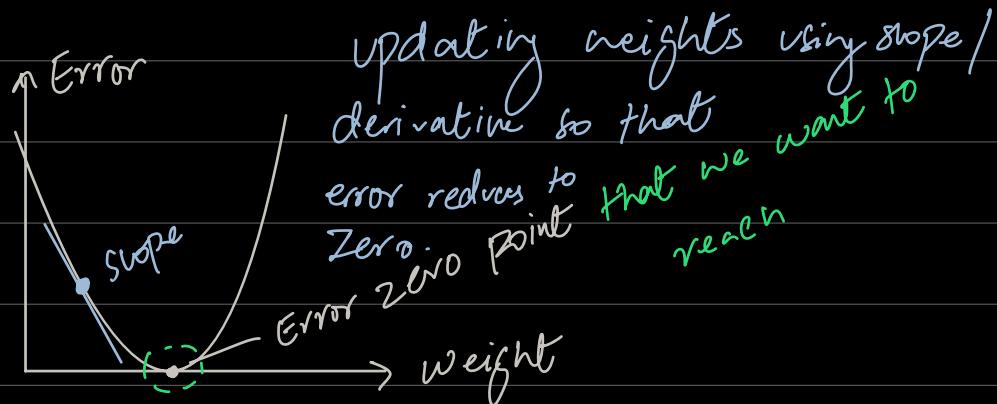
$\left\{ \begin{array}{l} \text{direction and amount} \end{array} \right\}$

\rightarrow to reduce $\frac{dy}{dx} = \frac{\text{change in node_delta}}{\text{change in weight}}$ - goal.

weight-delta or $\frac{dy}{dx} = \frac{1}{2} \frac{\text{change in } y \text{ increasing in } x}{\text{change in } x - \text{increased } 2 \text{ in } x}$

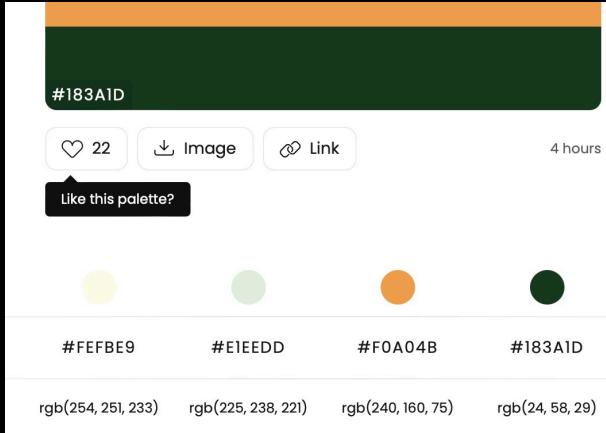
* Alpha Controls how quickly the network learns.

* With neural networks trying to find the lowest point on this big error plane, Where the lowest point refers to the lowest error.



- Gradient descent with multiple inputs and outputs:

o) generalizes to arbitrarily large networks.



$$256 \times 256 \times 256$$

Each pixel for Black & white image consists of 8 bits (1 byte) typically. Whereas each colored pixel consists of 24 bits (3 bytes)

$$\text{rgb} = [8 \text{ bits}, 8 \text{ bits}, 8 \text{ bits}]$$

$2^8 = 256$ different permutations
exist of a single Byte
Up to 8 permutations
so each Byte can represent a number from 0 - 256

$$2^8 \cdot 2^8 = 2^{16} = 65536$$

around and different combinations lead to different color.

$$2^2 = 4 \quad 2^{24} = 16,777,216 \text{ different colors}$$

possible outcomes
↑
2 coin

$$2^{[1, 0]}$$

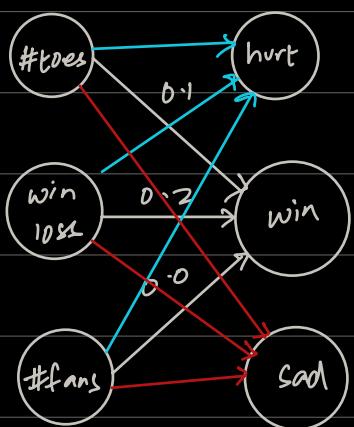
varyiations.

$$2^{24} \leftarrow \text{blocks or bits}$$

↑ possible outcome 1 or 0

inputs

Prediction



#toes %win #fans

$$\text{Weights} = \begin{bmatrix} [0.1, 0.1, -0.3] \\ [0.1, 0.2, 0.0] \\ [0.0, 1.3, 0.1] \end{bmatrix}$$

#hurt
#win
#sad

def vect_mat_mul(vec, matrix)

assert len(vec) == len(matrix)

def w_sum(a, b):

assert len(a) == len(b)

Output = 0

for i in range(len(a)):

Output += (a[i] * b[i])

return Output

Output = [0, 0, 0]

for i in range(len(vec)):

output[i] = w_sum(vec,

matrix[i])

def neural_network(input, weights):

pred = vect_mat_mul(input,
weight)

return Output

$$\text{toes} = [8.5, 9.5, 9.9, 9.0]$$

$$\text{wIrec} = [0.65, 0.8, 0.8, 0.9]$$

$$\text{nfans} = [1.2, 1.3, 0.5, 1.0]$$

$$\text{hurt} = [0.1, 0.0, 0.0, 0.1]$$

$$\text{win} = [1, 1, 0, 1]$$

$$\text{sad} = [0.1, 0.0, 0.1, 0.2]$$

$$\alpha = 0.01$$

$$\text{input} = [\text{toes}[0], \text{wIrec}[0], \text{nfans}[0]]$$

$$\text{true} = [\text{hurt}[0], \text{win}[0], \text{sad}[0]]$$

$\text{pred} = \text{neural_network}(\text{input},$
 $\text{weights})$

$$\text{error} = [0, 0, 0]$$

$$(\text{pred} - \text{true}) \leftarrow \delta = [0, 0, 0]$$

for i in range(len(true)):

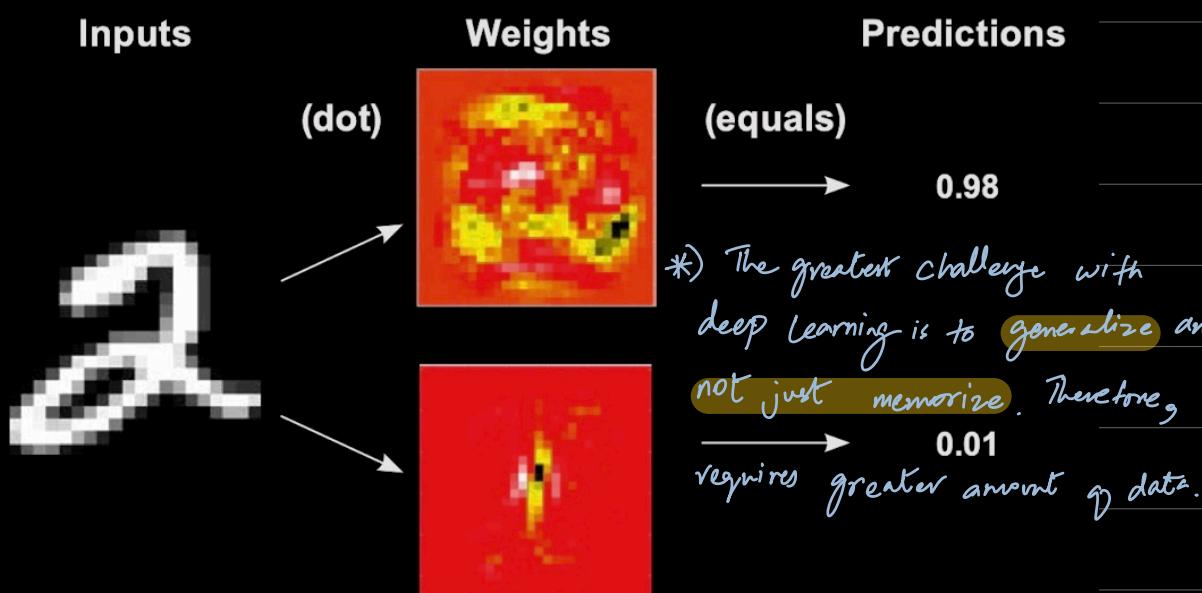
$$\text{error}[i] = (\text{pred}[i] - \text{true}[i])$$

$$\delta[i] = \text{pred}[i] - \text{true}[i]$$

* Calculating each weight_delta and putting it on each weight.

* A dot product is a loose measurement of similarity between two vectors.

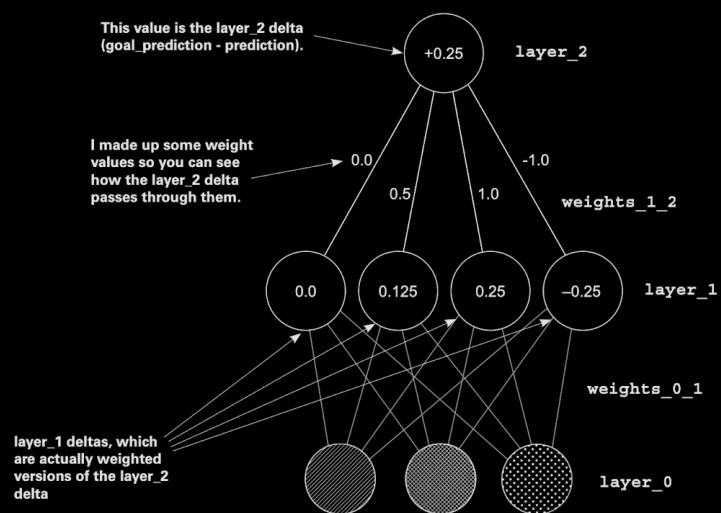
* If weight vector is similar to input vector for number = 2 image then dot product score / value is high and prediction is accurate



- The Prediction is a weighted sum of the inputs. The learning algorithm reward inputs that correlate with output with upward pressure (toward 1) on their weight while penalizing inputs with decorrelation with downward pressure. The weighted sum of inputs find perfect correlation b/w

Backpropagation: Long-distance error attribute inputs too. input and output by weighting decorrelated

- The prediction from layer_1 to layer_2 is weighted sum of the values at layer_1. If layer_2 is too high by x amount, then the higher weights contributed more.
- The ones with lower weights from layer_1 to layer_2 contributed less.
- Those weights also describe how much each layer_1 node contributes to the layer_2 prediction.
- The delta at layer_2 can be figured out by multiplying delta with each respective weights for layer_2. This process of moving delta signal is called backpropagation.



*) Regularization is advantageous because if a weight has equal pressure upward and downward then it is not good for anything. In essence, regularization aims to say that only weights with strong correlation can

- Linear vs non linear.

Stay on. Every else should be silenced
as contributing noise. (Trains

For any two multiplications, I can accomplish the faster).

Same thing using single multiplication. Which is bad.

\Rightarrow If trained the neural network as it is now,
it wouldn't converge.

Problem:

\Rightarrow For any two consecutive weighted sum of the input,
there exists a single weighted sum with exactly
identical behavior. Anything that the three-layer
network can do, the two layer network can also do.

- From a Correlation standpoint. Each node in the middle layer subscribes to a certain amount of correlation with each input node.
- If the weight from an input to the middle layer is 1.0, then it subscribes to exactly 100% of that nodes' movement. If that node goes up by 0.3, the middle node will follow.
- If the weight connecting two nodes is 0.5, each node in the middle layer subscribes to exactly 50% of the nodes' movement.

- Each hidden node subscribes to a little correlation from the input nodes.
- Since, the new data set has no correlation. The middle layer provides the ability to selectively correlate with the input.
- The middle layer should sometimes correlate with the input and sometimes not. This way middle layer not just x_i . Correlated to one input and y %. Correlated to another input. Instead, it can be x_i . Correlated to one input only when it wants to be, but other times not be correlated at all. Also known as Conditional Correlation or sometimes correlation.

\Rightarrow By turning off any middle node whenever it would be negative, you allow the network to sometimes subscribe to correlation from various inputs, which was impossible for two-layer neural networks.

\Rightarrow "If the node is negative, set it to 0" logic is non-linearity.

\Rightarrow Many kinds of non-linearities exist. Simplest one

that will be used is called ReLU

```
import numpy as np  
np.random.seed(1)
```

```
def relu(x): # returns x if x > 0;  
    return (x > 0) * x # returns 0 otherwise.
```

```
def relu2deriv(output): # returns 1 for input > 0;  
    return output > 0 # 0 otherwise.
```

```
streetlights = np.array([[1, 0, 1], [0, 1, 1],  
[0, 0, 1], [1, 1, 1]])
```

Transpose.

```
walk_vs_stop = np.array([[1, 1, 0, 0]]).T
```

$\alpha = 0.2$

$\text{hidden_size} = 4$

$\text{weights_0_1} = 2 * \text{np.random.random}((3, \text{hidden_size})) - 1$ 3 inputs \hookrightarrow connection
 \downarrow weight of each input

$\text{weights_1_2} = 2 * \text{np.random.random}((\text{hidden_size}, 1)) - 1$

for iteration in range(60):

$$\text{layer_2_error} = 0$$

for i in range(len(streetlights)):

$$\text{layer_0} = \text{streetlights}[i:i+1]$$

$$\# \text{function} \leftarrow \text{layer_1} = \text{relu}(\text{np}. \text{dot}(\text{layer_0}, \text{weights_0_1}))$$

$$\text{taking multiple arguments bcoz} \quad \text{layer_2} = \text{np}. \text{dot}(\text{layer_2}, \text{weights_1_2})$$

of numpy array

$$\text{layer_2_error} += \text{np}. \text{sum}((\text{layer_2} - \text{walk_vs_stop}[i:i+1])^2)$$

and processing it

$$\text{at the same time.} \quad \text{layer_2_delta} = (\text{layer_2} - \text{walk_vs_stop}[i:i+1])$$

time.

$$\text{layer_1_delta} = \text{layer_2_delta} \cdot \text{dot}(\text{weights_1_2})$$

$$\begin{aligned} & \text{weights_1_2} += \alpha * \text{layer_1_T} \cdot \text{dot}(\text{layer_2_delta}) \\ & \text{weights_0_1} += \alpha * \text{layer_0_T} \cdot \text{dot}(\text{layer_1_delta}) \end{aligned}$$

Back propagation

$$\text{if (iteration \% 10 == 9):}$$

print("Error : " + str(layer_2_error)) derivation

Purpose, as you'll see in a moment.

$$(3,1) \leftarrow \text{layer_0_T} \cdot \text{dot}(\text{layer_1_delta}) \Rightarrow (3,4)$$

* Every input of layer_0 gets

multipled with every delta of lay

- Neural networks seek to find direct and indirect correlation b/w an input layer and an output layer, which are determined by the input and output datasets.
- Local Correlation: Any given set of weights optimizes to learn how to correlate its input layer to with what the output layer says it should be.
- Global Correlation: What an earlier layer says it should be can be determined by taking what a later layer says it should be and multiplying it by weights in b/w them. This way later layers can tell earlier layers what kind of signal they need, to ultimately find the correlation with the output. This cross communication is called back propagation.

- In a neural network lines are weight-matrices
and input nodes are slice vectors.

*) Good neural architectures channel signal

so that correlation is easy to discover. Great
architectures also filter noise to help prevent
overfitting.

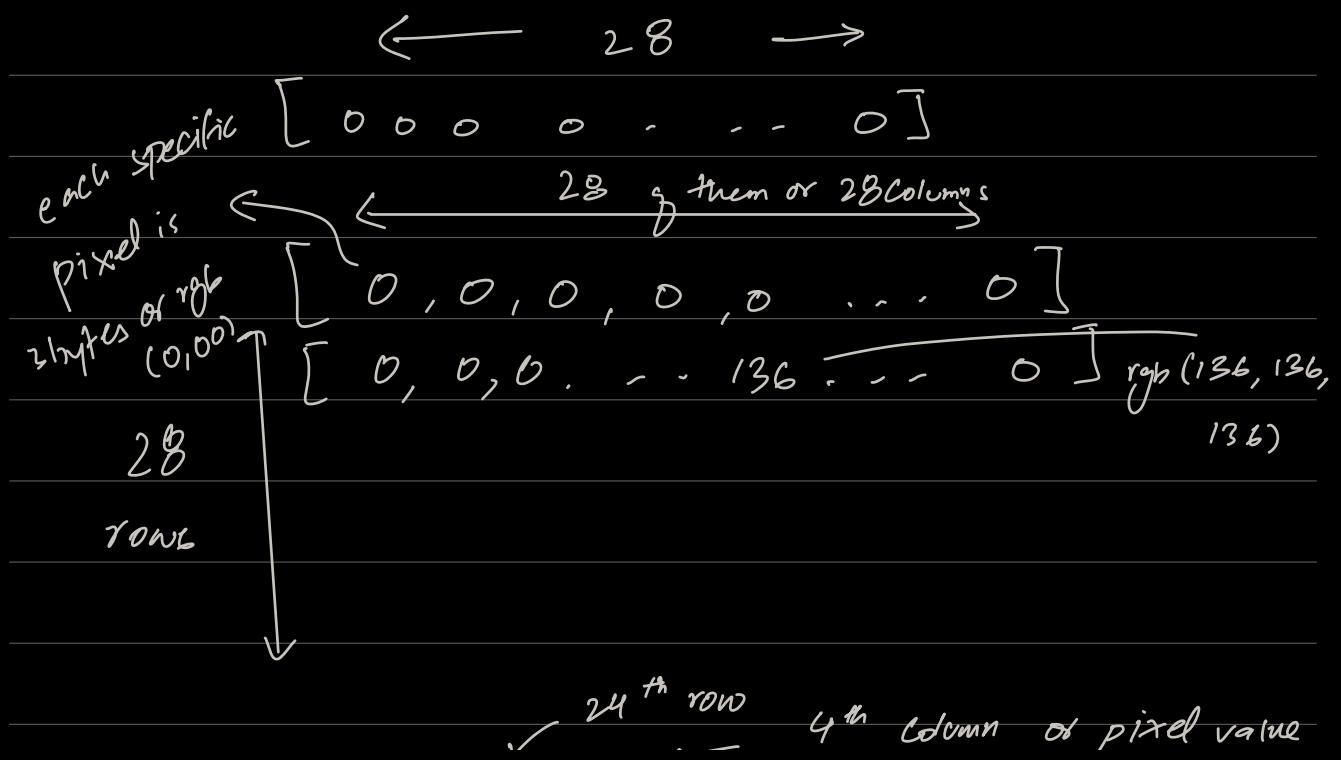
↳ incorrect correlation.

↳ learned to accurately predict but
somehow forgot to include a valuable
input.

*) The correlation is located wherever the weights
were set to high numbers.

28×28

$$\begin{array}{c} \cdots 28 \cdots \\ \hline | [0, 0, 0, 0, 0, 0, 0, 0] \cdots [0, 0, 0, 0, 0, 0, 0, 0] \\ | \\ \vdots \\ 28 \\ | \\ \vdots \\ | \end{array}$$



$x_train[0][24][4]$
 ↓
 first image
 ← Training dataset ← testing dataset
 (x_train, y_train) , (x_test, y_test)
 ↑ ↑
 input pixels output / result
 based on specific
 image.
 numpy.reshape (4, 3)
 ↑ ↑
 4 rows 3 columns

$x_train[0]$, $y_train[0]$
 ← First image ← output
 ← 1000 images ← columns
 ↓
 $(1000, 784)$
 ← 784 columns or pixels or one image

↑
 rows
 1000
 them
 or
 1000 images

→ 0.5333

0/255, 136/255

*) every pixel divided by 255.

but the color represented by

each pixel does not change.

* One-hot-labels - numpy array
 $\text{np.zeros}((\text{len(labels)}, 10))$

any particular image has 0 and 0
1000 rows and 10 columns
[0.0 . . . 0.0]
1000

Labels = [0 , 0 , 0, 0, 0, 0, 0, 0, 0, 0]

for i, l in enumerate(labels): # going through
one-hot-labels[i][l] = 1 1000 labels with
or indices with each

labels = one-hot-labels
index representing
actual output

with each index with a
number b/w 0

value of 1 representing the
actual true value / output or
and 0, set
it to 1 so

what number does the image
represent.
that we know
it is True
value.

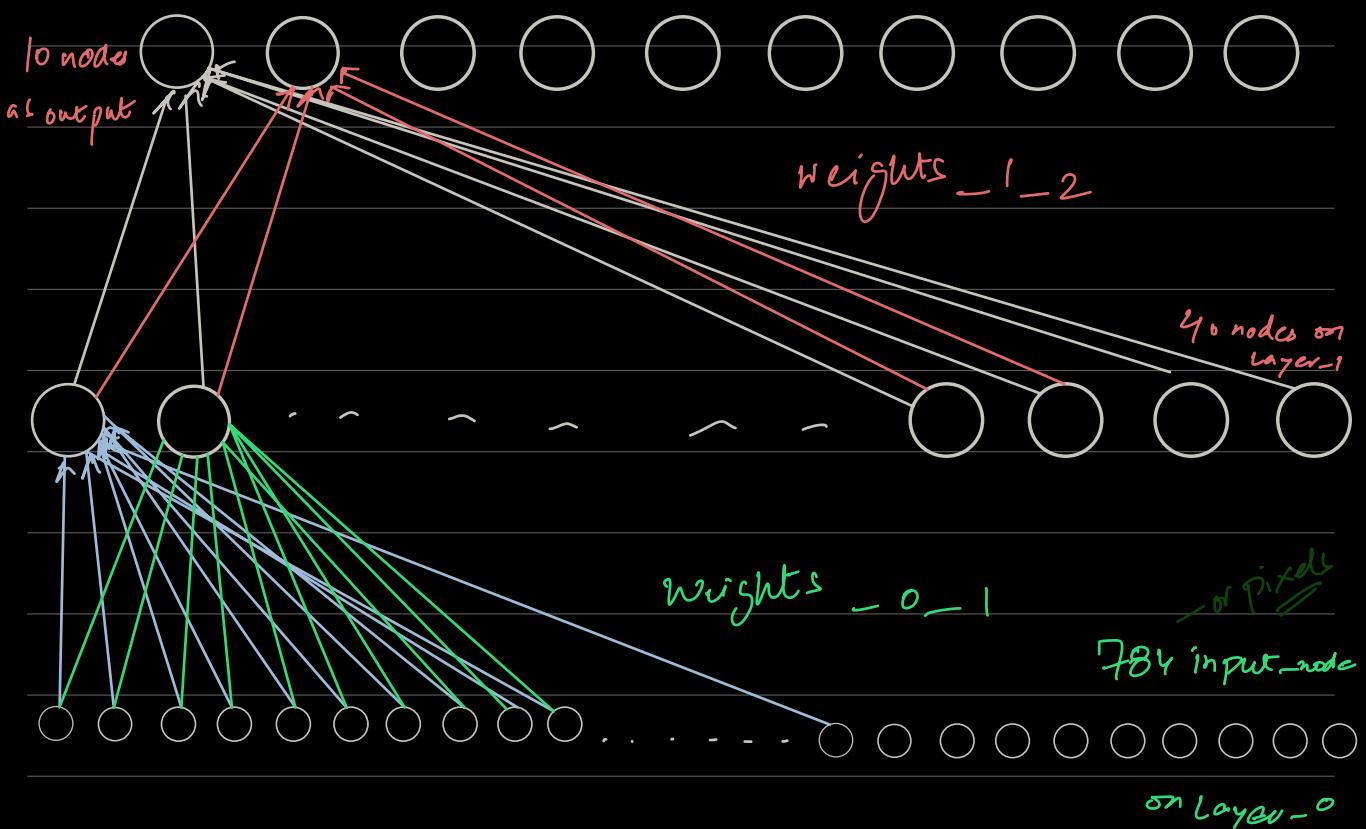
784 rows, 40 columns

layer-1 40 input nodes.

Layer-0 784 input nodes

$$\begin{bmatrix} 1 \times 784 \\ \text{in columns} \end{bmatrix} \begin{bmatrix} 784 \times 40 \end{bmatrix}$$

input nodes 11×40



```

import sys, numpy as np
from keras.datasets import mnist

(x_train, y_train), (x_test, y_test) = mnist.load_data()

images, labels = (x_train[0:1000].reshape(1000,28*28) \
                  255, y_train[0:1000])
one_hot_labels = np.zeros((len(labels),10))

for i,l in enumerate(labels):
    one_hot_labels[i][l] = 1
labels = one_hot_labels

test_images = x_test.reshape(len(x_test),28*28) / 255
test_labels = np.zeros((len(y_test),10))
for i,l in enumerate(y_test):
    test_labels[i][l] = 1

np.random.seed(1)
relu = lambda x:(x>=0) * x
relu2deriv = lambda x: x>=0
alpha, iterations, hidden_size, pixels_per_image, num_labels = \
(0.005, 350, 40, 784, 10)
weights_0_1 = 0.2*np.random.random((pixels_per_image,hidden_size)) - 0.1
weights_1_2 = 0.2*np.random.random((hidden_size,num_labels)) - 0.1

for j in range(iterations):
    error, correct_cnt = (0.0, 0)
    for i in range(len(images)):
        layer_0 = images[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)
        error += np.sum((labels[i:i+1] - layer_2) ** 2)
        correct_cnt += int(np.argmax(layer_2) == \
                           np.argmax(labels[i:i+1]))
        layer_2_delta = (labels[i:i+1] - layer_2)
        layer_1_delta = layer_2_delta.dot(weights_1_2.T) \
                      * relu2deriv(layer_1)
        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

    sys.stdout.write("\r"+ \
                    " I:"+str(j)+ \
                    " Error:" + str(error/float(len(images)))[0:5] + \
                    " Correct:" + str(correct_cnt/float(len(images))))

```

...
I:349 Error:0.108 Correct:1.0

*) The neural network learns to take a dataset of 1000 images and learn to correlate each input image with the correct label.

70.7% test accuracy - accuracy of the neural network on the data network was not trained on.

*) One way to view weights of a neural network is as a high-dimensional shape. As you train, this shape molds around the shape of data, learning to distinguish one pattern from another. The images in testing slightly different from patterns in the training dataset.

*) A neural network that overfits has learned the noise in the dataset instead of making decisions based on the true signal.

↳ learned noise instead of making decisions of true signal

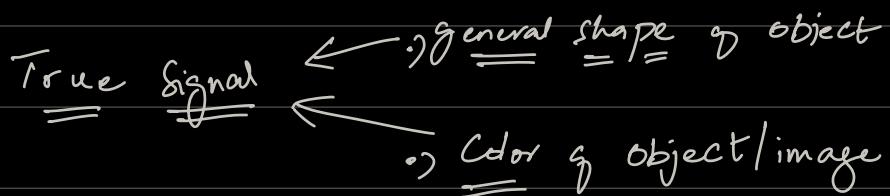
*) The mold got worse at the training dataset as you imprinted more forces because it learned more-detailed information about the

training dataset.

* How to get neural networks to train only the "signal" (the essence of dog) and ignore the noise. One way is early stopping.

* Large amount of Noise comes from fine-grained detail of image

* Most of the signal (for objects) is found in the general shape and in color of the image.



* Simplest regularization: Early Stopping:

- Stop training the network when it starts getting worse.
 - ↳ Don't let the network train long enough.
 - ↳ Early Stopping - cheapest form of regularization.
 - ↳ Subset of methods that help the neural network to learn the signal and ignore the noise.

Encourage generalization — often by making it difficult for a model to learn the fine-grained details of training data.

* Industry Standard Regularization: (Dropout)

The method: Randomly turn off neurons (set them to 0) during training.

↳ During training randomly setting neurons in network to 0.

↳ This way neural networks train exclusively using random subsections of the neural network.

↳ Dropout makes a big network act like a little one by randomly training little subsections of the network at a time, and little networks don't overfit.

↳ Larger strokes can capture more nuanced detail

↳ Form of training a bunch of networks and averaging them.

↳ neural networks learn by trial and error, and that every neural network learns on

little differently

↪ Randomly turning off

50% of the input

layer nodes or neurons.

Image[0 : 1] - image[0]
with 784 columns
pixels
row / First row
columns
excluding

1000 rows
columns
hidden

layer 0 = $\begin{pmatrix} 1 & 2 & \dots & 784 \end{pmatrix}$
weights₀ = $\begin{pmatrix} 1 & 40 \end{pmatrix}$
layer 1 = $\begin{pmatrix} 1 & 40 \end{pmatrix}$
40 columns

⇒ Modeling probabilities and non-linearities:

Activation Functions:

- Activation function
- Standard Activation functions \leq sigmoid
- Standard output Activation functions - Softmax
- Activation function installation instructions.

- What is an activation function:

A function applied to the neurons in a layer during prediction. Have used an activation relu in the three - layer neural network. The relu function had the effect of turning all negatives numbers to 0.

*) Constraints for Activation function:

- Constraint 1: The function must be continuous and differentiable in Domain, which means that there should be output number for any input.

- Constraint 2: Good activation functions are monotonic, never changing direction, which means the function should be 1:1, and must be either always increasing or decreasing

- Constraint 3: Good Activation functions are non-linear (they squiggle or turn).

To Create Some times Correlation you had to allow the neurons to selectively correlate to input

neurons such that a very -ve signal from one input into a neuron could reduce how much it correlated to any input (relu), Facilitated by any function that curves.

- Given a neuron with an activation function, you want one incoming signal to be able to increase or decrease how correlated the neuron is to all the other incoming signals.

Constraint 4: Good Activation functions (and their derivatives) should be efficiently computable.

* You'll often want to make multiple binary probabilities in one neural network. 98

* When a neural network has hidden layers, predicting multiple things at once can be beneficial. Often the network will learn something when predicting one label that will be useful to one of the other labels.

e.g.: learning / $^{O \times shape} \rightarrow ^Q$

*) Sigmoid Activation function can be of help in this instance because it models individual probabilities separately for each output node.

*) The Core Issue: Input have similarity
- Different Numbers Share Characteristics. The
Network Should Know about this.

*) The average 2 shares quite a bit in common with average 3.

*) General rule, similar inputs create similar outputs.

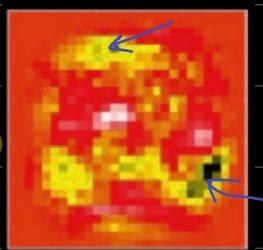
*) Side effect? Most images share lots of pixels in the middle of images, so the network will start trying to focus on the edges.

*) Consider the 2-detector node weights?

The heaviest weights are the end points

of the 2 toward the edge of the

image. It is best individual indicator of 2,



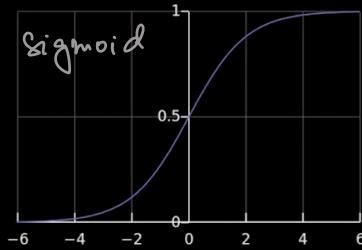
but the best overall is network that sees the entire shape for what it is.

*> Therefore, want an activation function that won't penalize labels that are similar.

*> Activation function that can pay attention to all the information that can be indicative of any potential input.

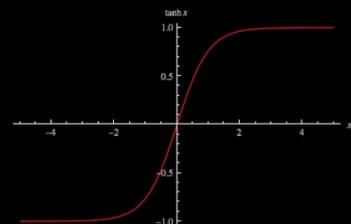
*> softmax Computation:

- Softmax raises each input value exponentially and then divides by the layer's sum.



(Image: Wikipedia)

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

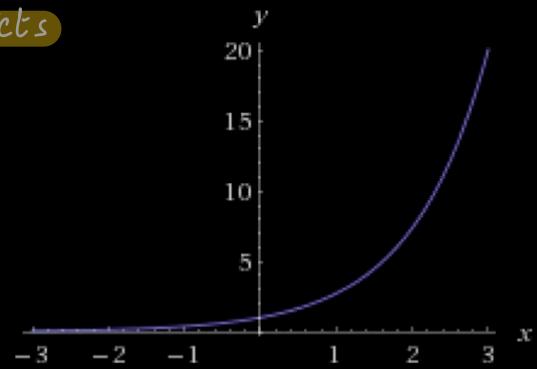


(Image: Wolfram Alpha)

*> For each value x , compute e^x , and then divide each value by the sum of all nodes in the layer.

* The higher softmax predicts one value, the lower it predicts all the others.

softmax



* It encourages the network to predict one output with very high probability.

* Adding an Activation function to forward propagation is straight forward. In the case of relu (rectified linear activation unit function), the function is applied to all the input nodes, where negative values are turned to 0's.

* However, properly compensating for activation function in backpropagation is bit nuanced.

* Whenever, relu had forced a layer_1 value to be 0, we also multiplied the delta by 0.

Reason:

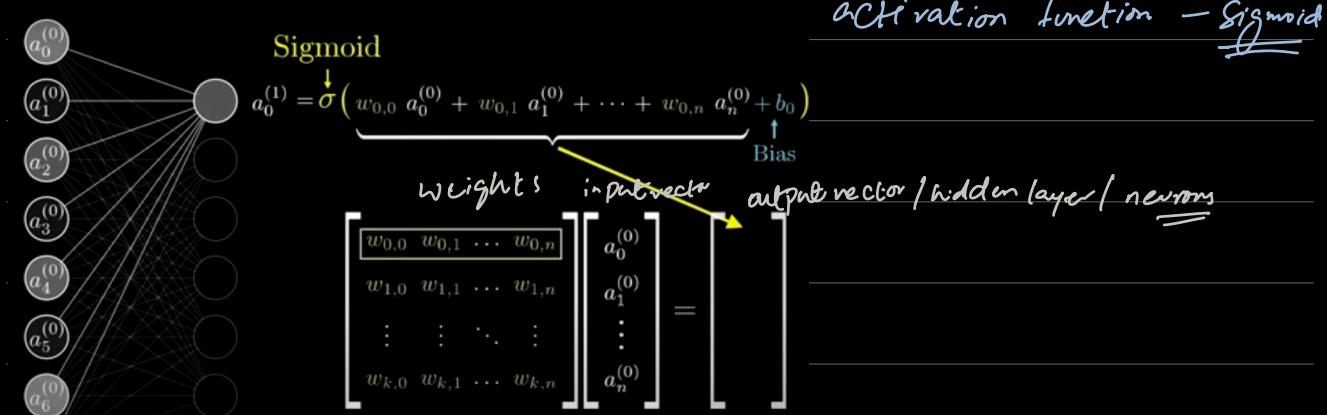
Because a Layer_1 value ≥ 0 had no effect on the output prediction, it shouldn't

have an impact on the weight update either.

It was not responsible for the error.

* The purpose of delta at this point is to tell earlier layers "make my input higher or lower next time". It modifies the delta propagated from the following layer to take into account whether this node contributed to error.

* When you backpropagate multiply backpropagated delta from Layer 2 (layer_2_delta . dot (weights_1_2 . T)) by the slope of relu at the point predicted in forward propagation. — Rule applies to other



What that means is that taking the weighted sum of the activations in the first layer according to these weights?

Following is a small table for the functions you've seen so far, paired with their derivatives. `input` is a NumPy vector (corresponding to the input to a layer). `output` is the prediction of the layer. `deriv` is the derivative of the vector of activation derivatives corresponding to the derivative of the activation at each node. `true` is the vector of true values (typically 1 for the correct label position, 0 everywhere else).

Function	Forward prop	Backprop delta
relu	<code>ones_and_zeros = (input > 0)</code> <code>output = input*ones_and_zeros</code>	<code>mask = output > 0</code> <code>deriv = output * mask</code>
sigmoid	<code>output = 1/(1 + np.exp(-input))</code>	<code>deriv = output*(1-output)</code>
tanh	<code>output = np.tanh(input)</code>	<code>deriv = 1 - (output**2)</code>
softmax	<code>temp = np.exp(input)</code> <code>output /= np.sum(temp)</code>	<code>temp = (output - true)</code> <code>output = temp/len(true)</code>

Note that the `delta` computation for `softmax` is special because it's used only for the last layer. There's a bit more going on (theoretically) than we have time to discuss here. For now, let's install some better activation functions in the MNIST classification network.

⇒ Neural Networks that understand Language :

King - man + Woman == ?

NLP (Natural language processing): Field that is exclusively dedicated to the automated understanding of human language. Deep learning's approach to this field.

matrix of numbers

↓

Raw text → Matrix of # → Supervised → What you
learning want to know

*) Converting text into numbers in such a way that the correlation b/w input and output is most obvious to the network. Resulting in faster training and better generalization.

*) Need to know what the input / output dataset looks like.

*) Topic classification

- Predicting whether people post positive or negative reviews.
- IMDB movie reviews dataset is a collection of review \rightarrow rating pairs (review(text), rating(number))

*) Input reviews are usually few sentences and output ratings are b/w 1 and 5 stars. — sentiment dataset.

*) Training neural network that can use the input text to make accurate predictions of the output score.

*) Adjusting the range of stars to be b/w 0 and 1 instead of 1 and 5, so that can use binary softmax to model output as probability.

*) Input data is text instead of numbers and it's variable-length text.

*) What about this data will have correlation with output.

*) Several words in this dataset would have a bit of correlation. Words such as terrible and unconvincing have negative correlation with the rating. ∵ As they increase in frequency in any input datapoint (review) \Rightarrow rating tend to decrease.

*) Words by themselves would have significant correlation with sentiment.

\Rightarrow Capturing word correlation in input data.

- Creating input matrix representing the vocabulary of review.
- Creating a matrix with each row(vector) representing

a review, and each column representing if a review contains particular word in the vocabulary

[$c_1 \ c_2 \dots$] - words in vocabulary
review

set - gets rid of duplicate elements.

tokens - list of sets, where every set refers to a review
vocab or words - length - [74074 and a list of words]

Word 2 index = ?

{ 'profession': index, 'some word': index }
1 74000

input - data set \Rightarrow every index or indices in

rows represent/refers some word.
reviews \Rightarrow every row/review consists of
list or vector of indices which are
values for keys or words appearing

target - data set if positive - 1 in that particular

negative - 0 review.

Topics: (We will use question/Answer dataset available on

- 1) Introduction to RNNs ✓ Pytorch.org)
- 2) RNNs ✓
- 3) Long short-term memory networks (LSTMs)
- 4) Implementation of RNN and LSTM.
- 5) Hyper parameters (might be used)
- 6) Seq2Seq (Encoder | Decoder)
- 7) Actual project Implementation

⇒ Chatbot: RNNs — LSTMs — Implementation —

Hyper parameters — seq2seq,

- Implementing different RNN solutions.
- RNNs modification to include a concept of memory.
- Advanced Implementation of LSTM - seq2seq
- Seq2seq uses Encoder/Decoder architecture, which led to transformers, BERT, GPT3.

- Requirements:

- 1) Pytorch, Numpy, Python
- 2) Build a simple neural network in pytorch.
- 3) Load CSV into a data from pandas.
- 4) Perform matrix multiplication with Numpy.

a) Object Oriented

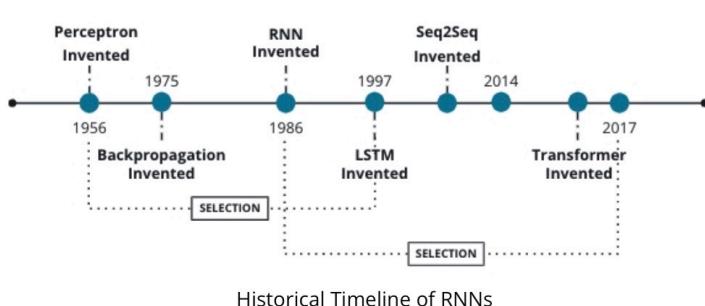
b) Visualizing data using plots in Jupyter notebook.

- matplotlib

- Jupyter notebook.

- Data processing and Object-Oriented architecture design patterns.

- Building new architectures for problems.



This chart shows the history of RNNs. Specifically, it describes how different techniques contributed to its development and how new ideas accelerated improvements through time. Our journey starts with the creation of the Perceptron in 1956 by Frank Rosenblatt. Frank Rosenblatt is considered by many to be the father of artificial intelligence. His creation, the Mark I Perceptron, is currently held in the Smithsonian Institute in Washington, D.C.

After the advent of the Perceptron, it took time before the other cornerstone of modern neural networks, the Backpropagation algorithm, was invented by Paul John Werbos in 1975 as a way to efficiently compute update gradients of the weights of a multi-layer network. David Rumelhart, Geoffrey Hinton, and Ronald Williams popularized this technique for neural networks.

Next, we progress to the RNN, which forms the foundation for this course. The RNN was invented in 1986 using a new technique to model sequential data called recurrence. This architecture was further modified by Jurgen Schmidhuber and Sepp Hochreiter from 1991 to 1997 into the LSTM to solve the vanishing gradient problem.

In 2014, the encoder-decoder architecture called Seq2Seq was introduced by Google to improve performance on sequence-to-sequence tasks. This architecture was carried forward by Google Brain into the advent of the Transformer in 2017. That is a lot of exciting ground to cover!

- RNN :

1) Modeling complex equations.

- to predict next value in a time series.
- model differential equations.

2) Natural language processing - understanding

complexities of word relationships in NLP

- generate text
- interactive chatbots.

Tools

- Python development environment.
- Pytorch - create neural networks.
- Pandas / Numpy
- Jupyter as IDE

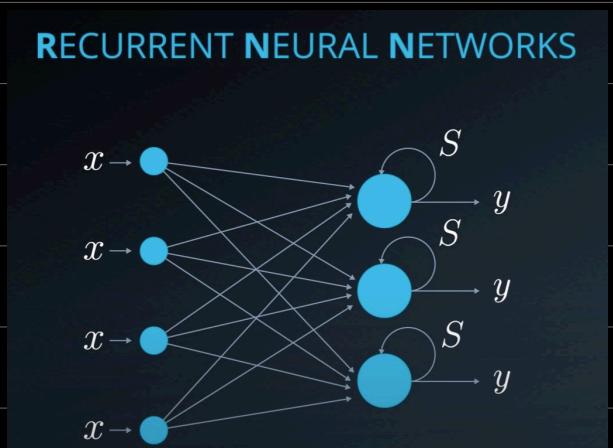
\Rightarrow Recurrent Neural Network:

1) Recurrence - occurring often

2) Memory - to produce output of the network.

- Backpropagation.

\Rightarrow Our current output does not belong on our current input but past input.



- * They can capture dependencies or dependencies over time
- * Perform same task for each element in the input sequence.

- Simple RNN ELMAN Network.

- Training RNNs

- LSTM

* RNNs give us a way to incorporate memory in our neural networks and are critical in analyzing sequential data.

*^U) Associated with text processing and text generation because of the way sentences are structured as a sequence of words.

*) RNNs and Text Generation:

TV Script generator — a model that takes in a series of words as input and outputs a likely next word, forming a text, one word at a time.

* Similarly, RNNs can be used to generate text given any other data. Giving an RNN a feature vector from an image can be used to generate a descriptive caption.

*) RNNs excel at sequence based tasks.

- Vanishing gradient problem:

In training our network, while using back propagation process, we adjust our weight-matrices using gradient/weight_delta/derivative/slope. In the process, gradients are calculated by continuous multiplication of derivatives. The value of derivatives

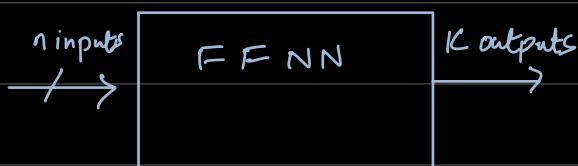
may be so small that these continuous multiplication may cause the gradient to practically "vanish".
LSTM - overcome vanishing gradient problem.

RNN and LSTM

- 1) Speech recognition
- 2) Time series
- 3) movie selection (Netflix)
- 4) Stock market
- 5) NLP (chatbots)

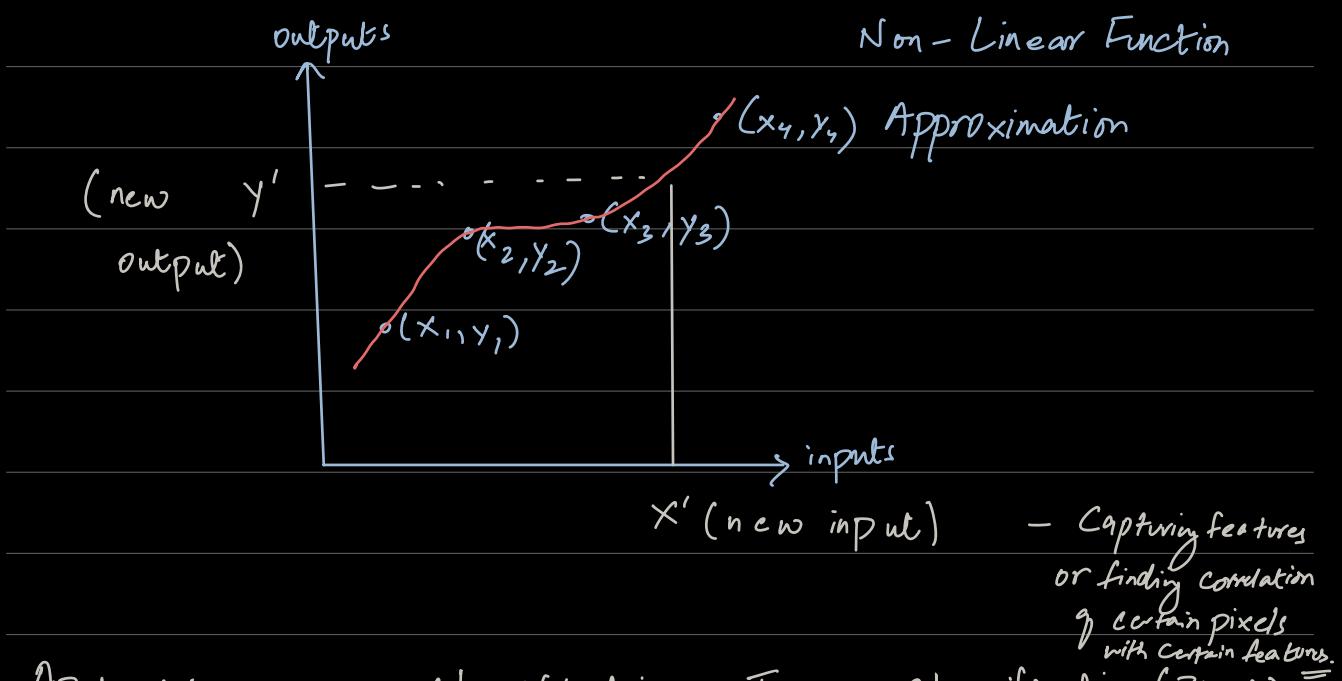
⇒ Feed Forward Neural Networks:

- Can have many hidden layers b/w the inputs and the outputs.
- Single hidden layer
- Simulate a artificial neural network by using a non-linear function approximation.



$$\text{output vector } \vec{y} = f(\vec{x}, w) \quad \text{input vector}$$

\vec{x} — weights

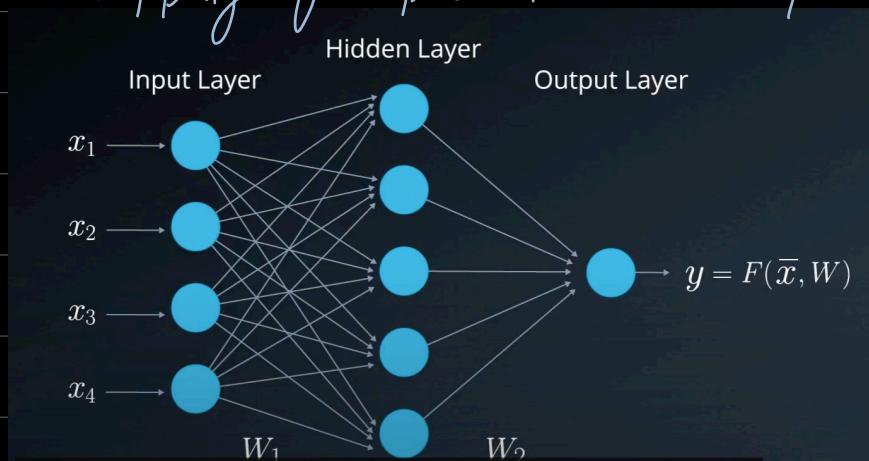


Applications :

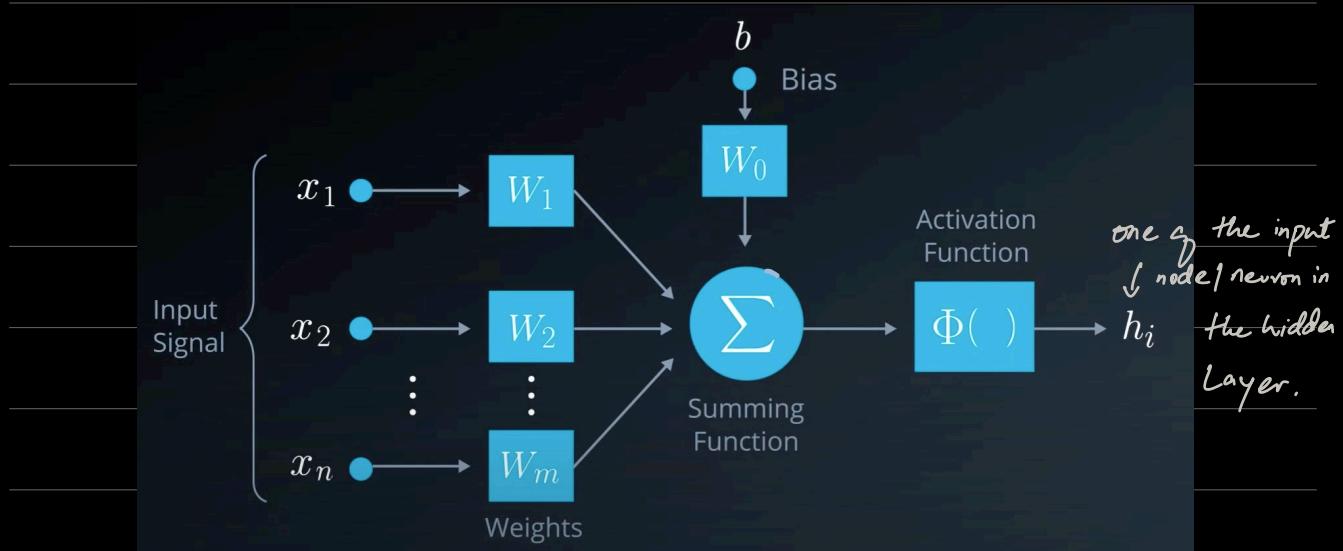
- Classification: Image Classification (Pixels) =
- Regression: Time-Series Forecasting
(predicting stock prices)

- Finding best set of weights that yield a good output.

- Static mapping of inputs to the output .



- Training: Training set or data, where neural network is trained to generalize beyond the training set or data.
- Evaluation: new inputs \longrightarrow Desired outputs.



\Rightarrow During the training phase, we take dataset (also called training set), which includes many pairs of inputs and their corresponding targets (outputs). Main goal is to find set of weights that would map the inputs to desired outputs.

In evaluation phase, we apply our new inputs and expect to obtain the desired outputs.

\Rightarrow Training phase:-

Feed forward and Backpropagation - steps repeated until we decide that our system has reached the best set of weights - error reduced to near zero.

i) W_K - weight matrix K

ii) W_{ij}^K - is the ij element of weight matrix K .
row column

\Rightarrow Calculate output value of perceptron

$$\begin{bmatrix} \text{(weights)}^T & \text{- Transpose} \\ 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} \text{input vector} \\ 1 \end{bmatrix} = 8 + 0 = 8$$

$$\begin{bmatrix} \text{input} \\ 1 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} \text{weights} \\ 8 \end{bmatrix}$$

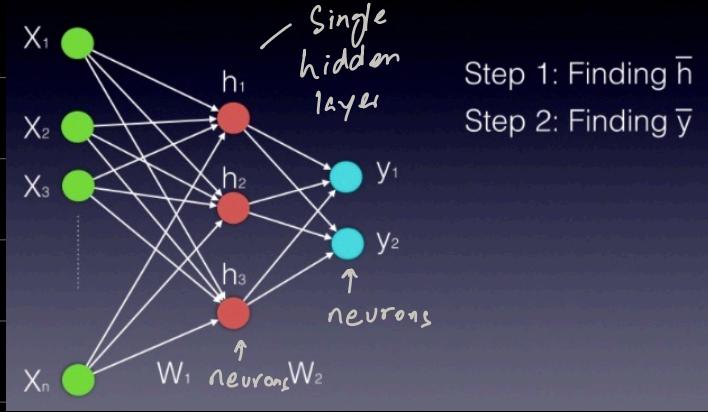
$$\text{Activation Function} \Rightarrow y = \frac{1}{(1 + e^{-z})}$$

$$y = \frac{1}{(1 + e^{-8})} = 1.0003$$

Activation

function model output

as probability



\Rightarrow Softmax encourages the network to predict one output with very high probability.

\Rightarrow Softmax is a useful activation function when we want all output values to be between 0 and 1, and their sum to be 1.

Error Functions:

\Rightarrow Mean squared Error function = $(\text{pred} - \text{true})^2$

\Rightarrow Cross Entropy

$$\bar{x} = [x_1 \ x_2 \ x_3 \dots x_n] \cdot \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ \vdots & \vdots & \vdots \\ w_{n1} & w_{n2} & w_{n3} \end{bmatrix}$$

$$\bar{h}' = [h_1' \ h_2' \ h_3']$$

↖ matrix multiplication gives us

hidden neurons

$$h'_i = x_1 \cdot w_{1i} + x_2 \cdot w_{2i} + \dots x_n \cdot w_{ni}$$

$f(x) = \tanh(x)$ - Activation function which makes sure

that values of neurons are b/w -1 and 1

$\sigma(x) = \frac{1}{1+e^{-x}}$ - Sigmoid Activation function makes sure that our values are b/w 0 and 1.

Activation function.

$$\boxed{\bar{h} = \phi(\bar{h}')} \quad \text{or} \quad \boxed{\bar{h} = \phi(\bar{x} \cdot \bar{w})}$$

hidden vector of neurons input

\Rightarrow Given \bar{h} \rightarrow Find \bar{y} (Apply step 1 again or matrix multiplication)

$$\boxed{h_i = \phi \left(\sum_i^n x_i w_{i1} \right)}$$

One hidden neuron

$$f(x, y) = 7x^2 - x^3y^4 + 5x^4y^3$$

partial derivative of function with respect to x .

- Every other variable except x - treat as constant

$$f_x = 7(2x) - (3x^2)y^4 + 5(4x^3)y^3$$

$$f_x = \boxed{14x - 3x^2y^4 + 20x^3y^3} \text{ first partial derivative.}$$

$$f_y = 0 - x^3(4y^3) + 5x^4(3y^2)$$

$$f_y = -4x^3y^3 + 15x^4y^2$$

$$z = f(x, y)$$

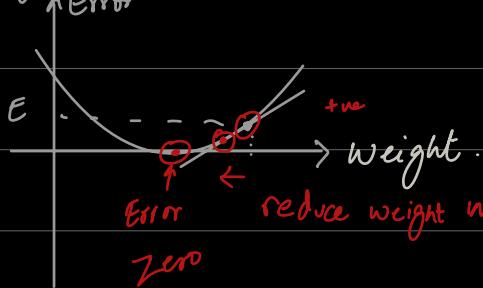
$$\frac{dz}{dx} = f_x, \quad \frac{dz}{dy} = f_y$$

$$\downarrow - \frac{1}{2} \leftarrow \text{cntr}$$

$\uparrow 2$ weight

slope / increase weight

\Rightarrow Back propagation:



$\uparrow 1$ slope /
 $\downarrow 2$ reduce weight

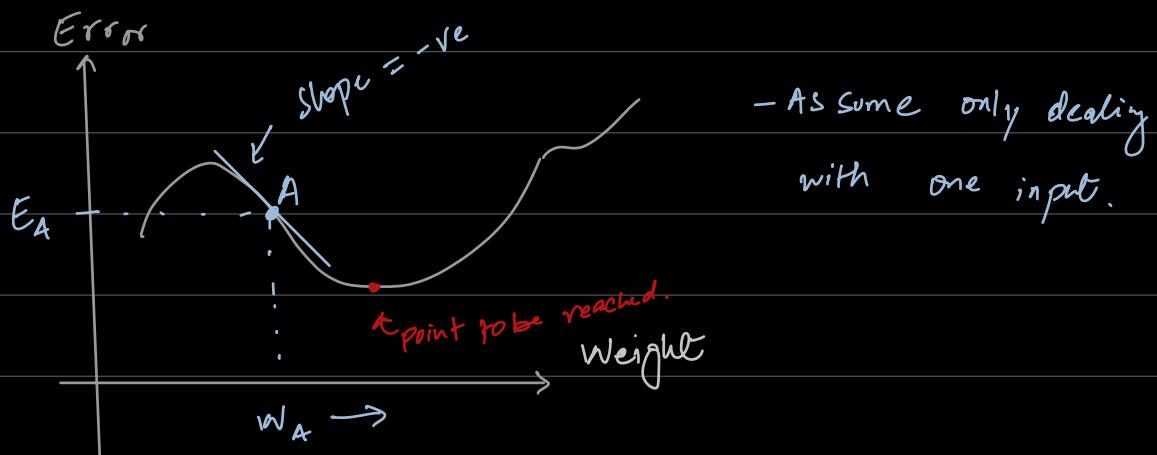
Reduce weight when gradient / slope is positive.

Increase weight when gradient / slope is negative.

- Partial derivative with respect to each weight.
- Error function consists of many weights / variables.

\Rightarrow Back propagation: Stochastic gradient Descent with chain rule.

- Finding set of weights that minimizes the network error.



Since, slope = -ve we have to change weight in opposite direction \therefore Have to increase weight so that the error goes down (vice versa)

$$W_{\text{new}} = W_{\text{previous}} + \alpha \left(- \frac{\partial E}{\partial w} \right)$$

\uparrow
step size

- The weight changes in the opposite direction of the partial derivative of the Error with respect to w.
- Partial derivative lets us measure how the error is impacted by each weight separately.

- Vector of partial-derivatives of the network error each with respect to different weight

$$W_{\text{new}} = W_{\text{previous}} + \alpha \underbrace{\left(\nabla_w (-E) \right)}_{\text{weight-deltas}} = \text{input} \times \text{delta}$$

Vector of partial derivatives.

①

$$\Delta w_{ij}^k = -\alpha \left(\frac{\partial E}{\partial w_{ij}^k} \right)$$

②

$$W_{\text{new}} = W_{\text{previous}} + \Delta w_{ij}^k$$

$$\text{loss function} \Rightarrow E = \frac{(\bar{d} - \bar{y})^2}{2} \quad \left. \begin{array}{l} \text{To model} \\ \text{error} \end{array} \right\}$$

↑
error function

\bar{d} - desired output

\bar{y} - calculated output.

$$\begin{matrix} \text{input} \\ \left[x_1 \ x_2 \right] \end{matrix} \quad \begin{matrix} \text{weights} \\ \left[\begin{matrix} a & b & c \\ d & e & f \end{matrix} \right] \end{matrix}$$

$$\begin{bmatrix} h_1 & h_2 & h_3 \end{bmatrix} \quad \begin{bmatrix} a \\ b \\ c \end{bmatrix} = \underline{y} \text{ or } \underline{\text{output}}$$

$$E = \frac{(d - y)^2}{2}$$

$$E = \frac{(\text{true} - \text{pred})^2}{2}$$

$$\Delta w_{ij} = -\alpha \left(\frac{\partial E}{\partial w_{ij}} \right)$$

Partial derivative of loss function E

with respect to the weight.

$$E = \frac{1}{2} (d - y)^2$$

constant
↓ Output/pred

$$F(x) = \frac{(x - y)^2}{2} = g(f(x))$$

$$F_x = \frac{1}{2} (x) (x - y) = x - y = g'(f(x))$$

$$F_y = -\frac{1}{2} (x) (x - y) = y - x$$

with respect to

$$y = + \frac{\alpha}{2} \cdot \left(\frac{\partial (d-y)^2}{\partial w_{ij}} \right) \quad E = (d-y)^2$$

$$\frac{\partial E}{\partial y} = +2(d-y)$$

$$\frac{dE}{dw} = \frac{dE}{dy} \cdot \frac{dy}{dw}$$

$(d-y)$ or delta $\frac{dE}{dw} = ?$

$$y = \sum_{i=1}^n (h_i w_i)$$

$\gamma = h_i w_i$ \downarrow
 $\frac{dy}{dw} = h_i$ * \triangleright ~~ell~~ $= \frac{d\gamma}{dw}$
 neuron

$$f(g(x)) = \sin(x^2)$$

$$f(y) = \sin(y) \quad y = g(x) = x^2$$

$$f'(y) = \cos y \quad g'(x) = 2x$$

$$f'(g(x)) \cdot g'(x)$$

$$f'(y) \cdot g'(x)$$

$$\cos y \cdot 2x - \frac{\alpha}{2}$$

$$f(g(a)) = \frac{(a-y)^2}{w_{ij}}$$

$$f(g(x)) = f'(g(x)) \circ g(x)$$

$$c = (\overset{\text{pred}}{a} - \overset{\text{true}}{y})^2$$

$$\frac{dc}{da} = 2(a-y) \quad f(g(x)) = g'(x) \cdot$$

$$\frac{d}{da} f'(g(x))$$

- constant =

$$\frac{dc}{dw} = \underbrace{\left(\frac{da}{dw} \right)}_{\text{Change in output over change in weight}} \circ \frac{dc}{da}$$

Change in output

over change in weight

$$a = i \cdot w$$

$$\frac{da}{dw} = i$$

$$E = \frac{(d-y)^2}{2} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{loss function}$$

$$*) \Delta w_{ij}^k = \alpha \left(- \frac{dE}{dw} \right) \quad \left. \begin{array}{l} \text{Partial derivative of loss function} \\ y = \sum_{i=1}^z (h_i w_i) \end{array} \right.$$

$$\Delta w_{ij}^k = \alpha \left(- \frac{dE}{dw} \right) = - \frac{\alpha}{2} \left(\frac{d \overbrace{(d-y)^2}^E}{d \overbrace{w_{ij}}^{\text{desired output}}} \right) \frac{dE}{dw_{ij}}$$

$$E = (d-y)^2$$

$$\frac{dE}{dw_{ij}} = -2(d-y)$$

$$dy$$

$$\frac{dE}{dw_{ij}} = \frac{dE}{dy} \cdot \frac{dy}{dw_{ij}} \quad \left. \begin{array}{l} \text{Chain rule} \\ \text{out} \end{array} \right\}$$

$$\frac{dE}{dw_{ij}} = -2(d-y) \cdot \frac{dy}{dw_{ij}}$$

$$\Delta w_{ij}^k = -\alpha \left(-\delta(d-y) \cdot \frac{dy}{dw_{ij}} \right)$$

$$\Delta w_{ij}^k = \alpha \left[(d-y) \cdot \begin{bmatrix} \frac{dy}{dw_{ij}} \\ \vdots \\ \frac{dy}{dw_{ij}} \end{bmatrix} \right]$$

gradient
or
 δ_{ij}

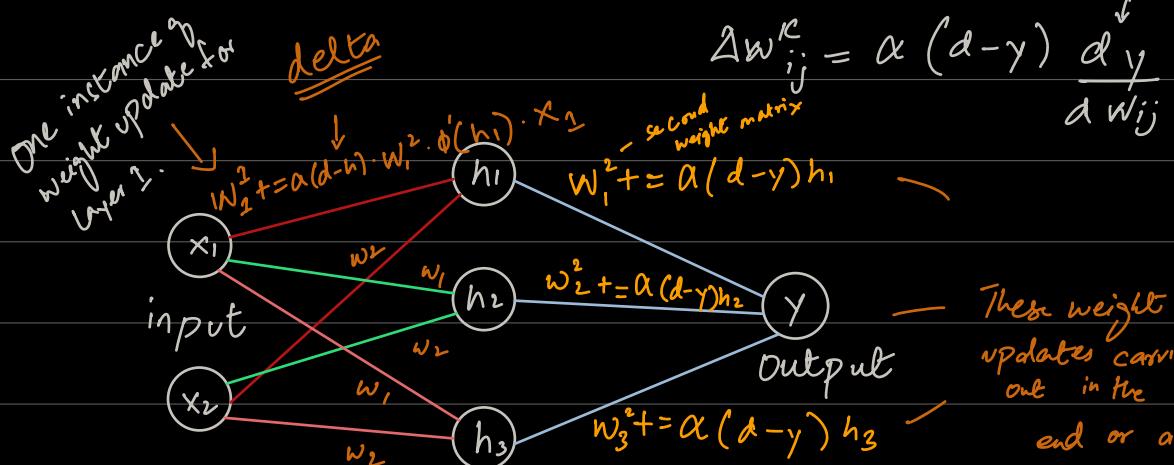
$$y = \sum_{i=1}^3 (h_i w_i)$$

$$\delta_{ij} = \frac{dy}{dw_{ij}}$$

$$\delta \text{elta} * \text{weight} = \text{node_delta}$$

$$\text{node_delta} * \text{input}(x_i) = \underline{\text{weight_delta}}$$

gradient



These weight updates carried out in the end or after Layer 1 update.

$$\text{output} \Rightarrow y = \sum_{i=1}^3 (h_i w_i) \quad x_1 w_1 + x_2 w_2$$

w^2 - 2nd weight matrix

w' - 1st weight matrix

$$h_1 = \sum_{i=1}^2 (x_i w_i)$$

$$h_1 = x_1 w_1 + x_2 w_2, \quad h_2 = x_1 w_3 + x_2 w_4$$

$$h_1 = \phi(x_1 w_1 + x_2 w_2) \quad h_3 = x_1 w_5 + x_2 w_6$$

function function

activation function applied.

$$h_j = \phi_j \left(\sum_{i=1}^2 x_i w_i \right)$$

hidden function

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} w_1 & w_3 & w_5 \\ w_2 & w_4 & w_6 \end{bmatrix}$$

$$w^2 = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix}$$

$$h_1 = x_1 w_1 + x_2 w_2$$

$$h_2 = x_1 w_3 + x_2 w_4$$

$$h_3 = x_1 w_5 + x_2 w_6$$

$$y = \sum_{i=1}^3 (h_i w_i)$$

$$y = h w$$

$$\frac{dy}{dw} = h$$

hidden neurons vector.

$$\frac{dy}{dw_{ij}} = h_i^T$$

Partial derivative
with respect to
 w is everything
else treated as
constant.

\Rightarrow Step 1 of Backpropagation:

$$\Delta w_{ij} = \alpha(d - y) h_i \quad \left. \begin{array}{l} \text{first} \\ \text{weight} \end{array} \right\}$$

update

(Backwards)

\Rightarrow Step 2 of Backpropagation.

Activation

$$h_j = \phi_j \left(\sum_{i=1}^2 x_i w_{ij} \right)$$

gradient

$$\frac{dy}{dw_{ij}} = \frac{dy}{h_i} \cdot \frac{dh}{w_{ij}} \cdot \frac{dh}{y} \cdot \frac{dy}{x_i \cdot \phi_j}$$

chain rule.

$$\frac{dh}{dw_{ij}}$$

$$\Delta w_{ij} = \alpha(d - y) \frac{dy}{w_{ij}}$$

$$w_{new} = w_{prev} + \alpha(d - y) \cdot w_j^2 \cdot \phi'_j \cdot x_i$$

$$y = \sum_{i=1}^3 (h_i w_i^2)$$

$$\frac{dy}{dh} = w_i^2 \quad \# \text{ weight - 2 matrix or vector.}$$

↑

Partial derivative
with respect to h .

Activation - second function. $f(g(x)) \Rightarrow \text{derivation} = f'(g(x)) \cdot g'(x)$

one function

$$h_j = \phi_j \left(\sum_{i=1}^2 x_i w_i \right)$$

$$h = x \cdot w$$

$$\frac{dh}{dw} = x$$

input passed to the activation
function derivative - a neuron

$$\frac{dh_j}{dw_i} = \phi'_j \cdot \left(\sum_{i=2}^2 (x_i w_i) \right) \circ x_i$$

Each neuron j in the hidden layer will have its own value for ϕ - activation function and

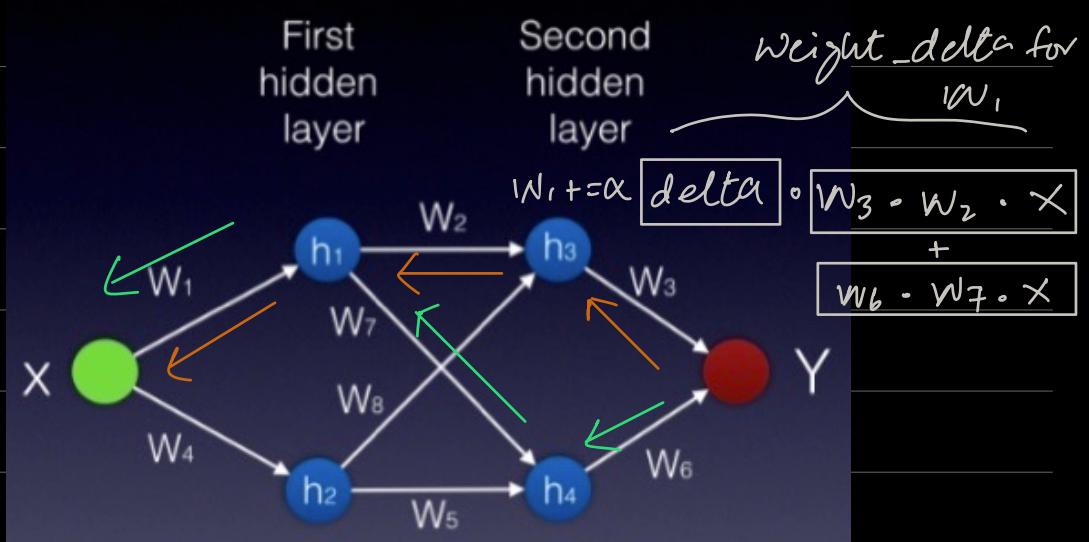
for ϕ' - activation function derivative

$$\therefore \frac{dh_j}{dw_{ij}^{-1}} = \phi'_j \cdot x_i \quad \begin{matrix} \text{input vector.} \\ \uparrow \\ \text{derivative of activation} \\ \text{function.} \end{matrix}$$

$$\therefore \frac{dy}{dw_{ij}^{-1}} = \frac{dy}{dh_j} \cdot \frac{dh_j}{dw_{ij}^{-1}} = w_j^2 \cdot \phi'_j \cdot x_i \quad \begin{matrix} \text{gradient} \\ \downarrow \\ \text{gradient} \end{matrix}$$

$$\therefore \Delta w_{ij}^{-1} = \alpha(d - y) \cdot w_j^2 \cdot \phi'_j \cdot x_i \quad \begin{matrix} \text{weight update} \\ \downarrow \end{matrix}$$

$$\therefore w_{\text{new}}' = w_{\text{prev}}' + \alpha(d - y) \cdot w_j^2 \cdot \phi'_j \cdot x_i$$



Weight update rule of weight matrix W_1

or partial derivative of y with respect to w_1

$$y = h_3 w_3$$

$$y = h_3 w_3 + h_4 w_6$$

$$\frac{dy}{h_4} = w_6$$

$$\frac{dy}{h_3} = w_3$$

$$h_4 = h_1 w_7 + h_2 w_5$$

$$h_3 = h_1 w_2 + h_2 w_8$$

$$\frac{dh_4}{h_1} = w_7$$

$$\frac{dh_3}{h_1} = w_2$$

$$h_1 = X w_1$$

$$h_2 = X w_4$$

$$\frac{dh_1}{w_1} = X$$

$$\frac{dy}{w_1} = \left(\left[\frac{dh_1}{w_1} \cdot \frac{dh_3}{h_1} \cdot \frac{dy}{h_3} \right] + \right.$$

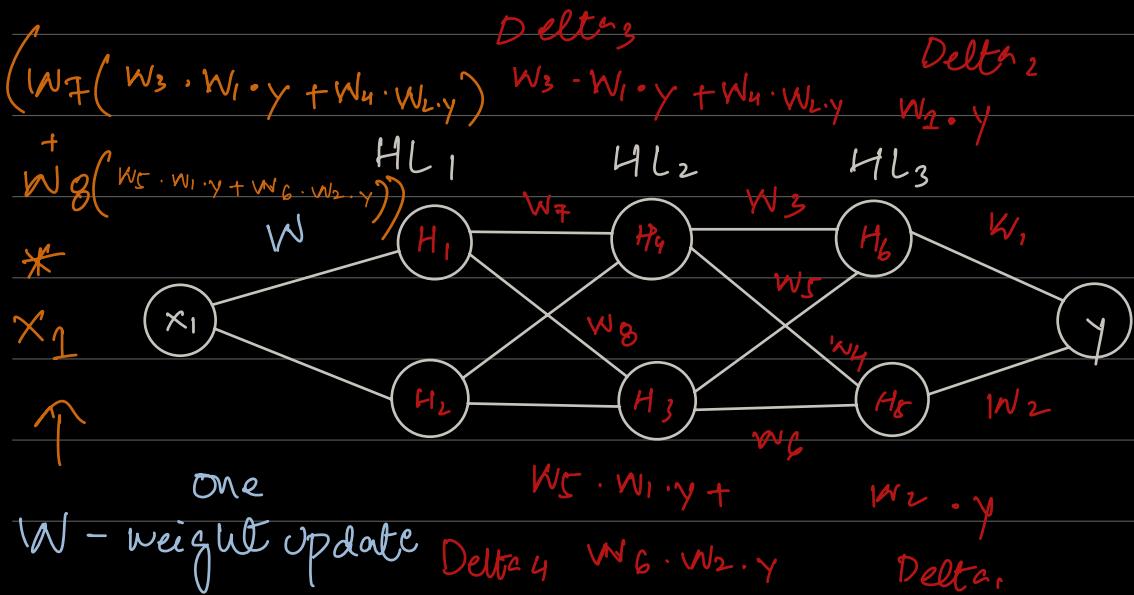
$$\left. \left[\frac{dh_1}{w_1} \cdot \frac{dh_4}{h_1} \cdot \frac{dy}{h_4} \right] \right)$$

W_1 - weight update

$$W_1 += \alpha (d - y) * \frac{\partial y}{\partial W_1}$$

↑ ↑ gradient
alpha delta

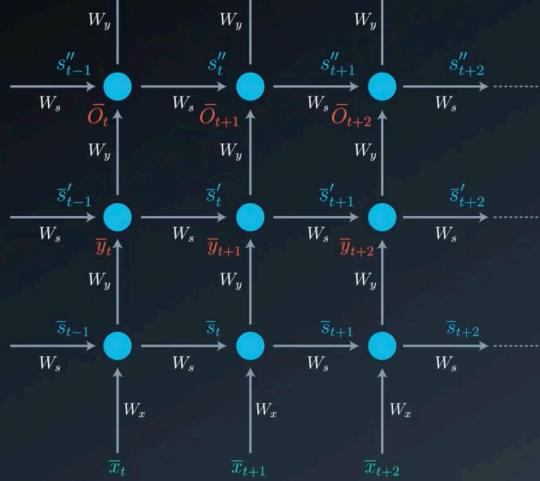
$$W_1 += \alpha (d - y) \left[(W_3 \circ W_2 \circ x) + (W_6 \cdot W_7 \cdot x) \right]$$



\Rightarrow Recurrent Neural Network :

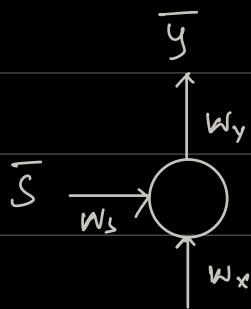
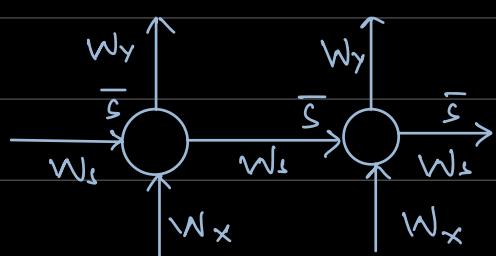
Recurrent - occurring repeatedly

RECURRENT NEURAL NETWORK



- Perform same task for each element in the input sequence
- Capturing information in previous inputs by maintaining internal memory elements also known as states.

Memory Elements



\bar{x}_t

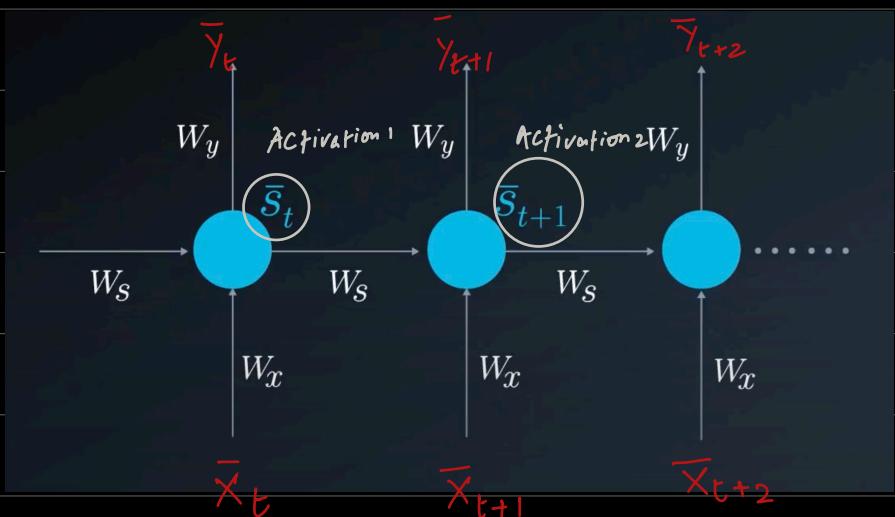
\bar{x}_{t+1}

\bar{x}

*) Temporal dependencies: Current output does not only depend on the current input, but also on a memory element, which takes into account the past inputs — use RNN

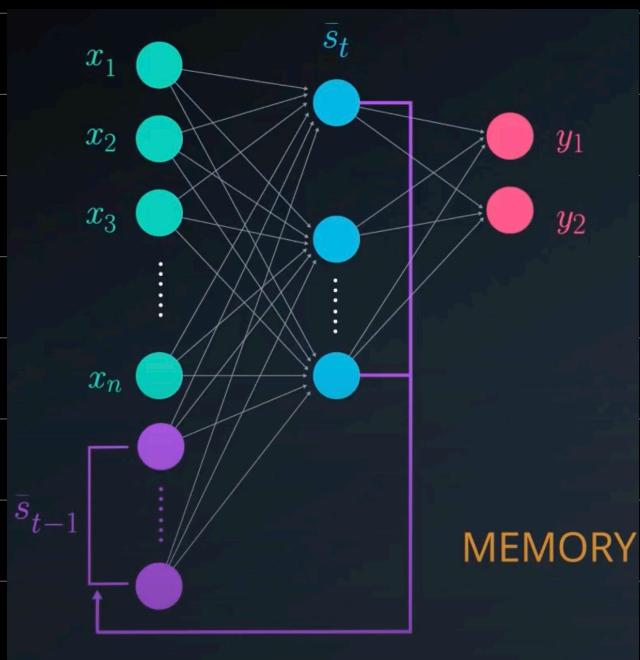
*) A good use of RNN — is predicting the next word in a sentence, which requires looking at the last few words rather than the current one only

Ex : To be or not to —



1) We train with Sequences because previous inputs matter.

2) The second difference, stems from memory elements that RNN's host. Current inputs, as well as activations of neurons serve as inputs to the next time step.



* The Feed forward scheme changes and includes feedback or memory elements — memory as output of the hidden layer, which will serve as additional input to the network at the following training step. This way the system learns to optimize the weight matrix.

RNN uses:

- sequences as inputs in the training process.
- memory elements.

- FFNN: The output at any time t is a function of the current input and the weights alone.

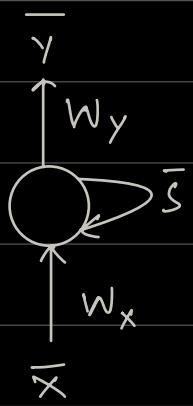
Time t

$$\bar{h}_t = F(\bar{x}_t, w)$$

- RNN: The output at time t does not only depend on current inputs and weights, but also on previous inputs.

Time t

$$\bar{y}_t = F(\bar{x}_t, \bar{x}_{t-1} \dots \bar{x}_{t-n}, w)$$



- w_x represents weight matrix connecting the inputs to the state layer.

- w_y represents weight matrix connecting the

State to output.

- W_s represents weight

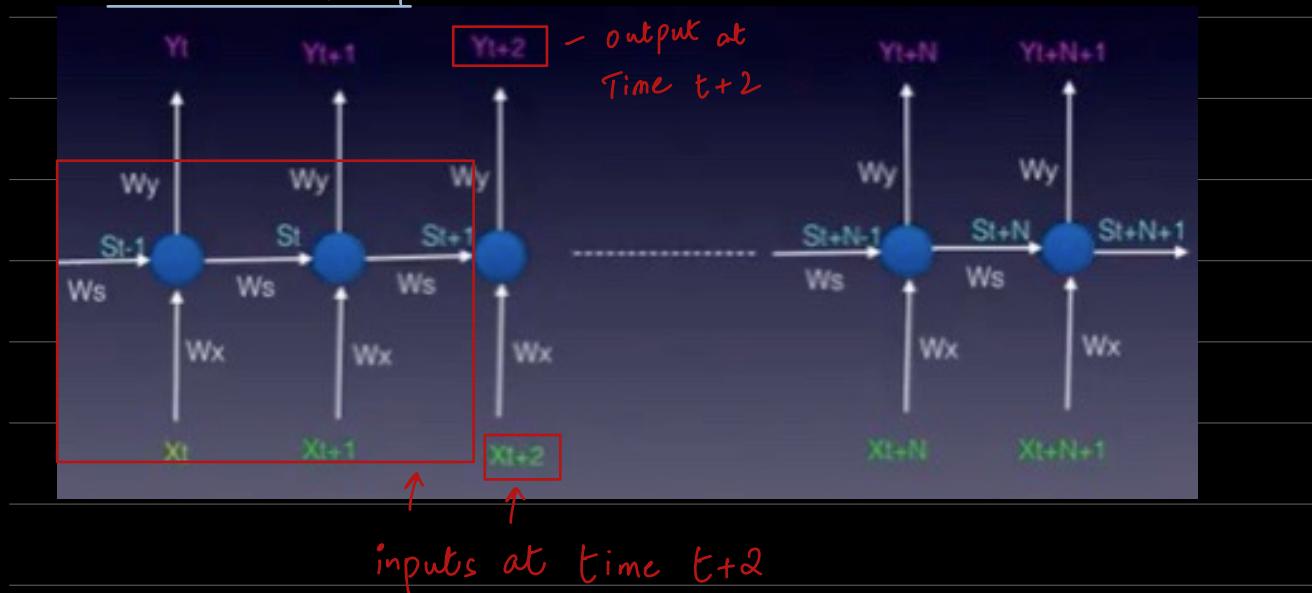
matrix connecting the

the state from previous

time step to the state

in the next time step.

\Rightarrow Unfolded RNN:



inputs at time $t+2$

- Trained on sequences as inputs.

- Memory elements.

FFNN

$$\bar{h} = \Phi(\bar{x}, W)$$

Activation

RNN

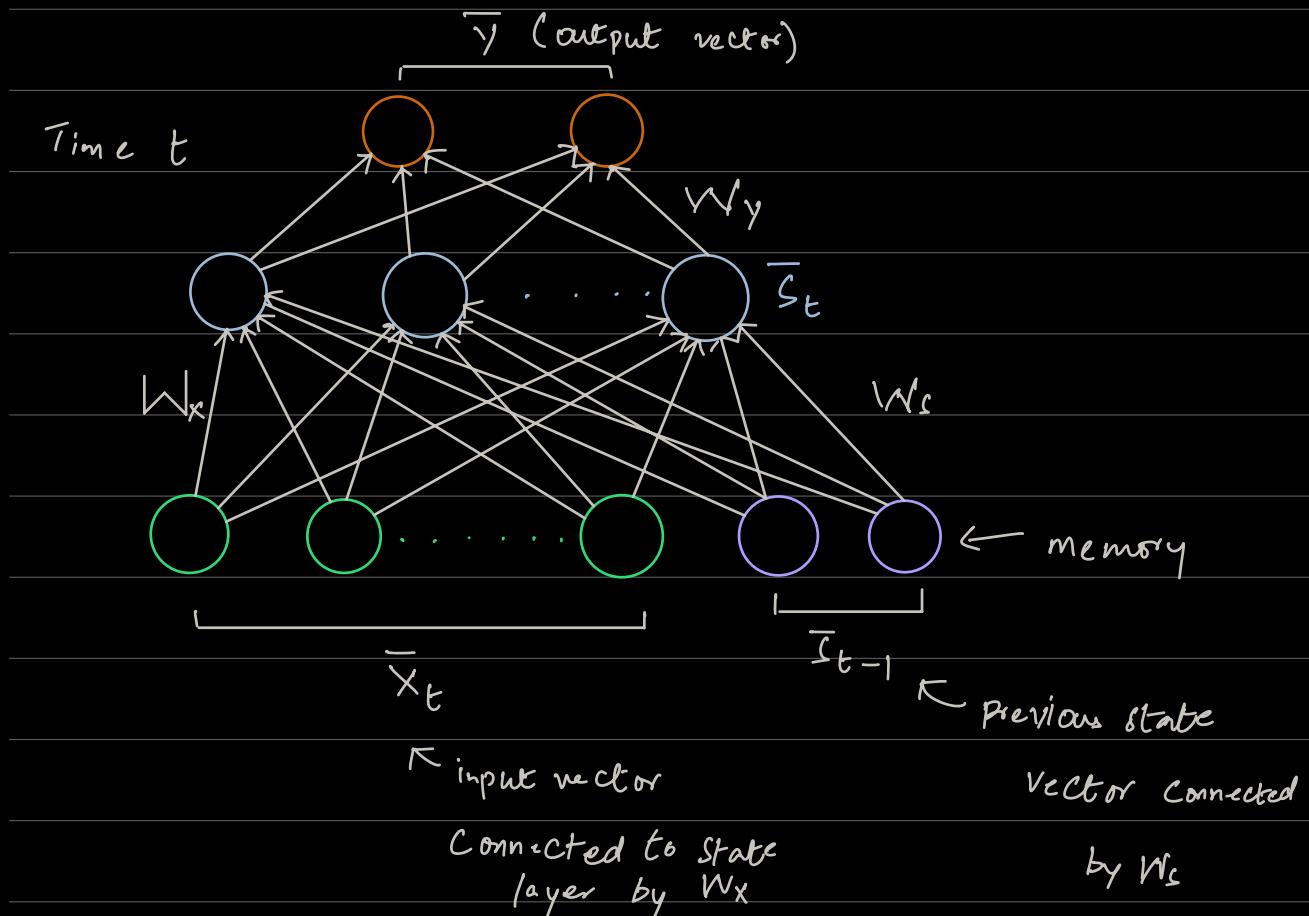
$$\bar{s}_t = \Phi(\bar{x}_t W_x + \overbrace{\bar{s}_{t-1} \cdot W_s}^{\text{Previous state}})$$

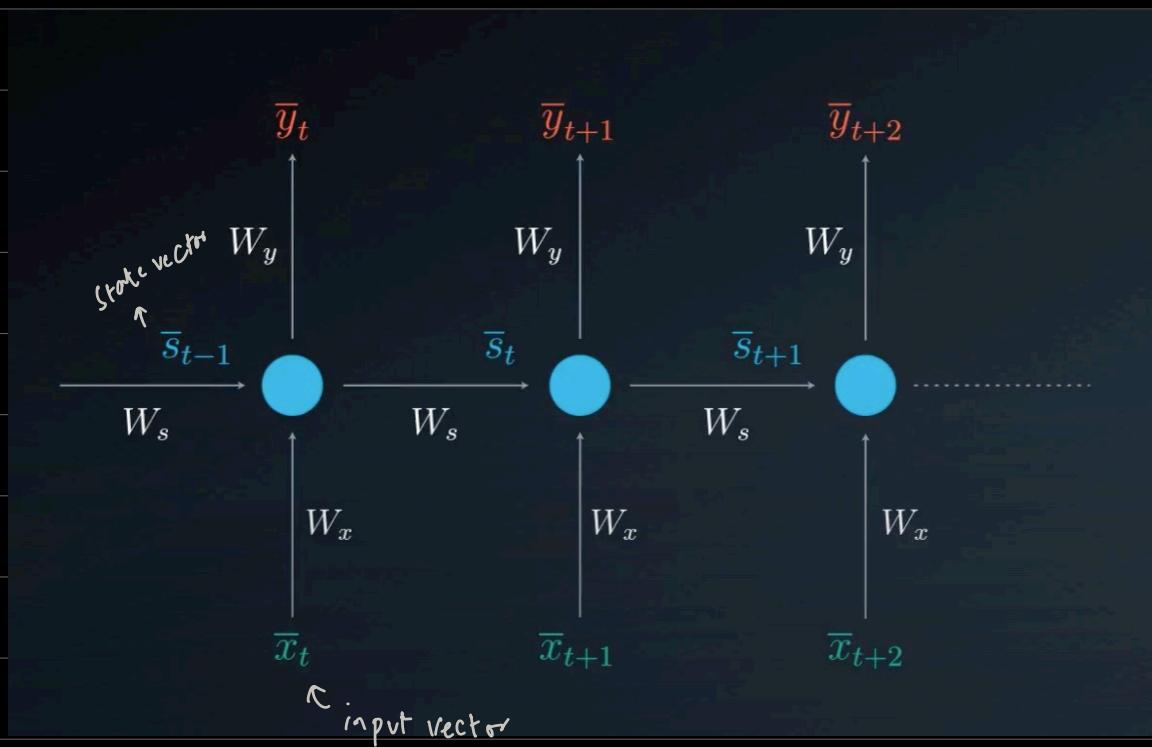
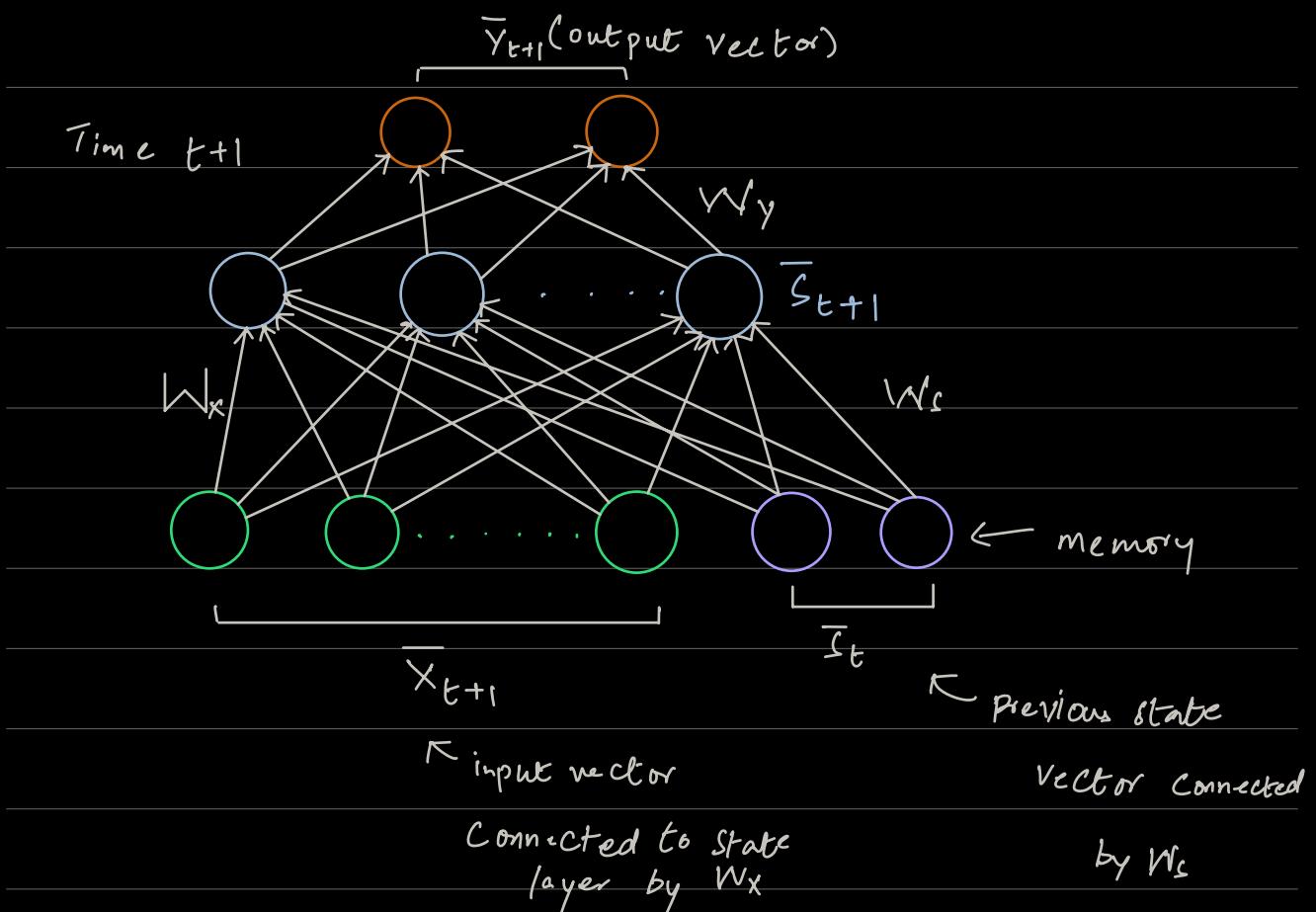
$$\bar{y}_t = \bar{s}_t \circ w_y$$

↑ ↑
Output or State input vector
output vector.

Weight matrix

\Rightarrow Elman Network:





- RNN example :

.) Sequence / Word detector :

Detecting word "happy"

Arranging letters of Word "happy" in ascending order
and creating one-hot vector encoding to be fed
as input to the neural network.

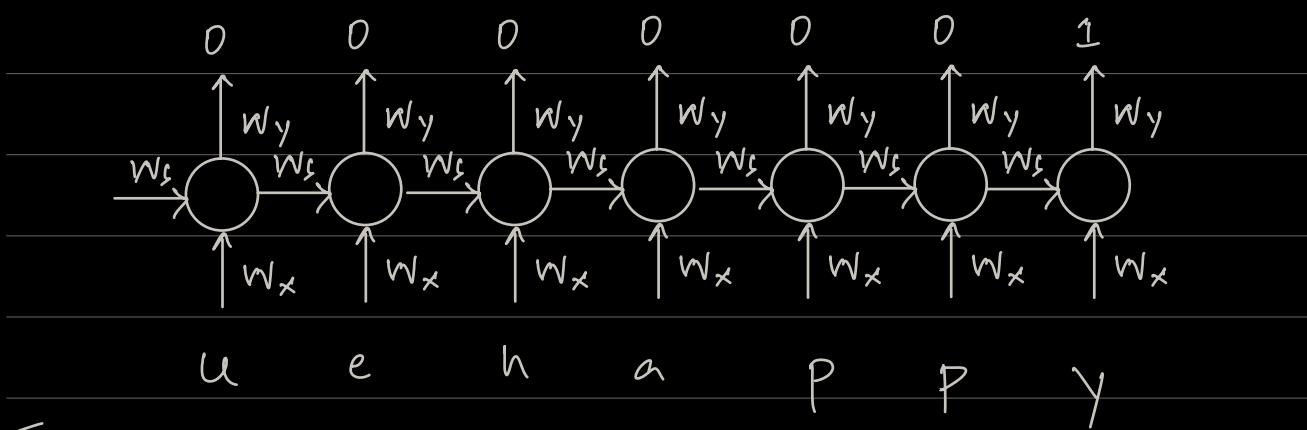
1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

a h p P y

Sequence detected Configuration:

0	1	0	0	0	0
1	0	0	0	0	0
0	0	1	0	0	0
0	0	0	1	0	0
0	0	0	0	1	0

h a p P y "e"
DNF ..

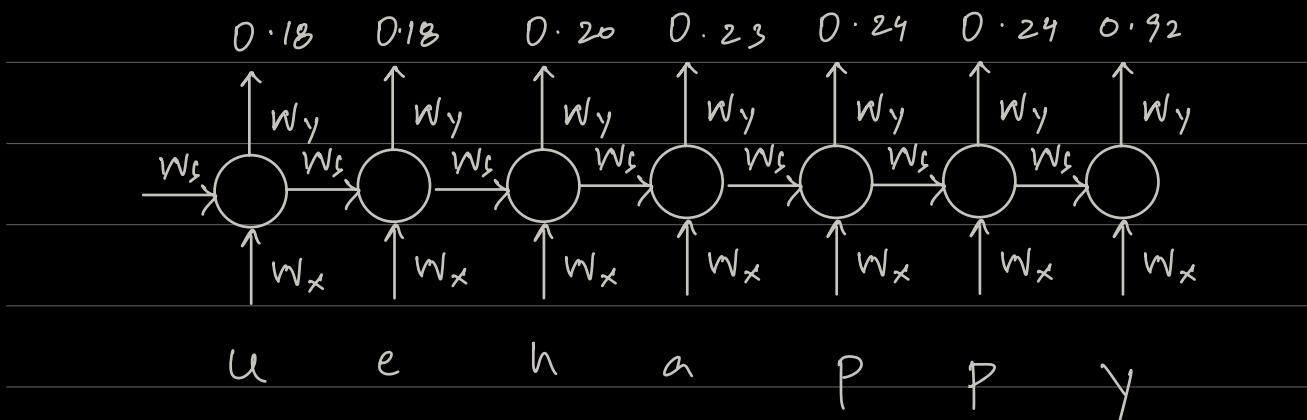


Training:

* Will keep feeding the network random letters and setting target value to 0 if sequence word not detected.

* When sequence / word detected set target to 1 as occasionally will feed word "happy".

Evaluation phase:



When threshold or probability > 0.9 — sequence is detected

\Rightarrow Back propagation through time:

- Find derivatives of loss function at time $t=3$ as a function of all of the weight matrices.

(FFNN), but have to consider previous time steps.

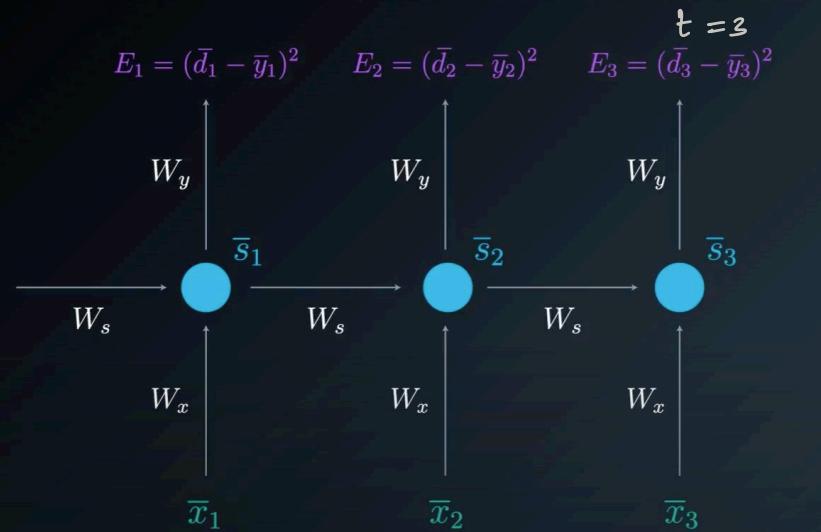
$\frac{dE}{dW_x} \left(\begin{matrix} \text{delta} \cdot w_x \\ \text{delta}^+ \cdot w_s \end{matrix} \right)$

$E = \bar{s} \cdot w_y$

$\frac{dE}{dW_y} = \bar{s}$

$\frac{dE}{dW_s} = ?$, $\frac{dE}{dW_x} = ?$

BACKPROPAGATION THROUGH TIME



- In this model we have to modify three weight matrices :

W_x - linking the network input to the state or hidden layer

W_s - connecting one state to the next.

W_y - connecting this state to the output.

- Adjusting Weight Matrix W_y :

$$E_3 = \frac{d_{\text{desired}}}{d_{\text{output}}} = \frac{(d_3 - y_3)^2}{\xi_2}$$
$$\frac{dE_3}{W_y} = \frac{\delta E}{d\bar{y}_3} \cdot \frac{dy_3}{dW_y}$$

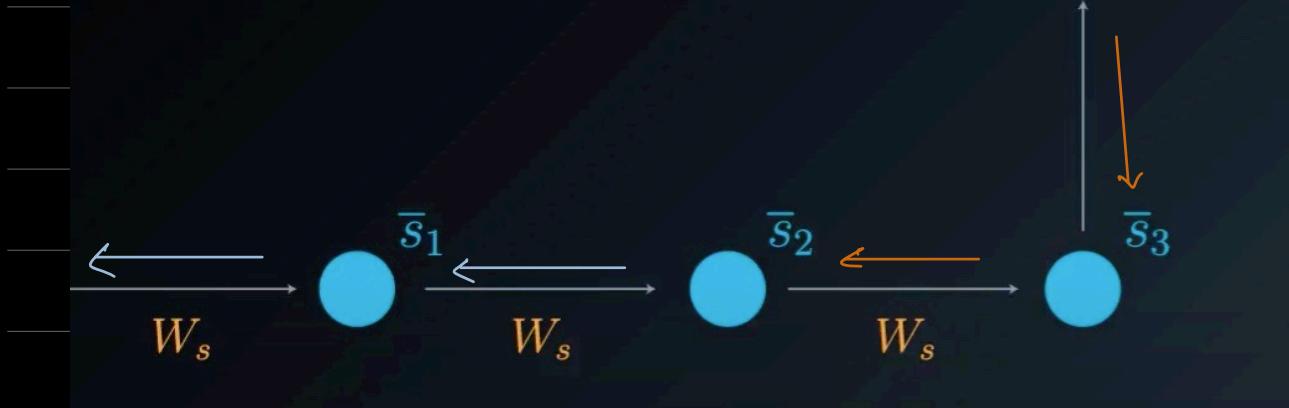
$$W_y = W_y + \alpha (d - y) \cdot \frac{\delta y_3}{dW_y}$$

- Adjusting weight matrix - W_s :

- weight matrix connecting one state to the next and removing anything that we don't need

for now.

$$E_3 = (\bar{d}_3 - \bar{y}_3)^2$$



$$s_3 = s_2 \cdot w_s \quad s_2 = s_1 \cdot w_s$$

at time $t = 3$

$$\frac{ds_3}{ds_2} = w_s \quad \frac{ds_2}{dw_s} = s_1$$

$$\frac{dE}{dw_s} = \frac{dE}{dy_3} \cdot \frac{dy}{ds_3} \cdot \frac{ds_3}{dw_s} \quad \left. \right\} \text{ till } \bar{s}_2$$

+

$$\frac{dE}{d\bar{y}_3} \cdot \frac{d\bar{y}}{ds_3} \cdot \frac{ds_3}{ds_2} \cdot \frac{ds_2}{dw_s} \quad \left. \right\} \text{ till } \bar{s}_1$$

+

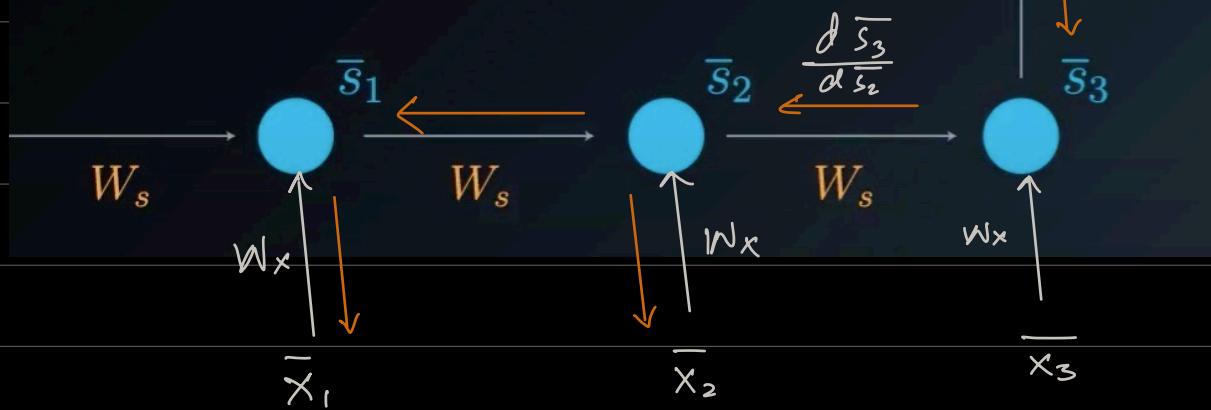
$$\frac{dE}{d\bar{y}_3} \cdot \frac{d\bar{y}}{ds_3} \cdot \frac{ds_3}{ds_2} \cdot \frac{ds_2}{ds_1} \cdot \frac{ds_1}{dw_s} \quad \left. \right\} \text{ till } \bar{s}$$

Weight update of w_s

$$\frac{dE_N}{dw_s} = \sum_{i=1}^n \frac{dE_N}{d\bar{y}_N} \cdot \frac{d\bar{y}_N}{d\bar{s}_i} \cdot \frac{ds_i}{dw_s}$$

— Adjusting w_x :

$$\begin{aligned} \bar{s}_3 &= \bar{x}_3 \cdot w_x \quad y = s_3 \cdot w_y \quad E = (\bar{d}_3 - y_3)^2 \\ \frac{d\bar{s}_3}{dw_x} &= \bar{x}_3 \quad \frac{dy}{dw_y} = s_3 \quad \frac{dE}{dy} = (d_3 - y_3) \quad E_3 = (\bar{d}_3 - \bar{y}_3)^2 \end{aligned}$$



at time $t = 3$:

$$\frac{dE}{dw_x} = \frac{dE}{d\bar{y}_3} \cdot \frac{d\bar{y}_3}{d\bar{s}_3} \cdot \frac{d\bar{s}_3}{dw_x}$$

+

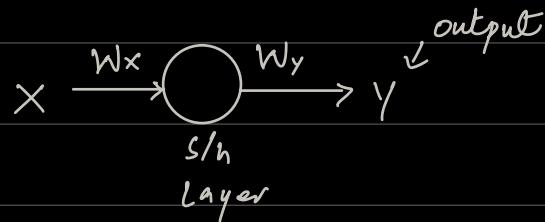
$$\frac{dE}{d\bar{y}_3} \cdot \frac{d\bar{y}_3}{d\bar{s}_3} \cdot \frac{d\bar{s}_3}{d\bar{s}_2} \cdot \frac{d\bar{s}_2}{dw_x}$$

$$\frac{dE_N}{dw_x} = \sum_{i=1}^n \frac{dE_i}{d\bar{y}_N} \cdot \frac{d\bar{y}_N}{d\bar{s}_i} \cdot \frac{d\bar{s}_i}{dw_x}$$

+

$$\frac{dE}{d\bar{y}_3} \cdot \frac{d\bar{y}_3}{d\bar{s}_3} \cdot \frac{d\bar{s}_3}{d\bar{s}_2} \cdot \frac{d\bar{s}_2}{d\bar{s}_1} \cdot \frac{d\bar{s}_1}{dw_x}$$

*) Simple Back propagation :



$$y = h \circ w_y \quad h = x \cdot w_x$$

Error function

$$E = \frac{(d - y)^2}{2}$$

$$\begin{aligned} \frac{dE}{dw_y} &= \frac{dE}{dy} \cdot \frac{dy}{dw_y} \\ &\stackrel{\text{update to weight}}{=} \end{aligned}$$

$$\frac{dE}{dy} = -(d - y)$$

$$\frac{dy}{dw_y} = - (d - y) \cdot h$$

$$\frac{dE}{dy} = -\alpha (d - y) \cdot h$$

$$\frac{\partial E}{\partial w_x} = \underbrace{\frac{\partial E}{\partial y} \circ \frac{\partial y}{\partial h} \circ \frac{\partial h}{\partial w_x}}_{\text{Chain rule}} = -(d-y) \cdot w_y \cdot x$$

update to weight matrix w_x

$$w_x += -\alpha (d-y) \cdot w_y \cdot x$$

Gradient Calculations

$$\frac{\partial E_N}{\partial W_y} = \frac{\partial E_N}{\partial \bar{y}_N} \cdot \frac{\partial \bar{y}_N}{\partial W_y}$$

$$\frac{\partial E_N}{\partial W_s} = \sum_{i=1}^N \frac{\partial E_N}{\partial \bar{y}_N} \cdot \frac{\partial \bar{y}_N}{\partial \bar{S}_i} \cdot \frac{\partial \bar{S}_i}{\partial W_s}$$

$$\frac{\partial E_N}{\partial W_x} = \sum_{i=1}^N \frac{\partial E_N}{\partial \bar{y}_N} \cdot \frac{\partial \bar{y}_N}{\partial \bar{S}_i} \cdot \frac{\partial \bar{S}_i}{\partial W_x}$$