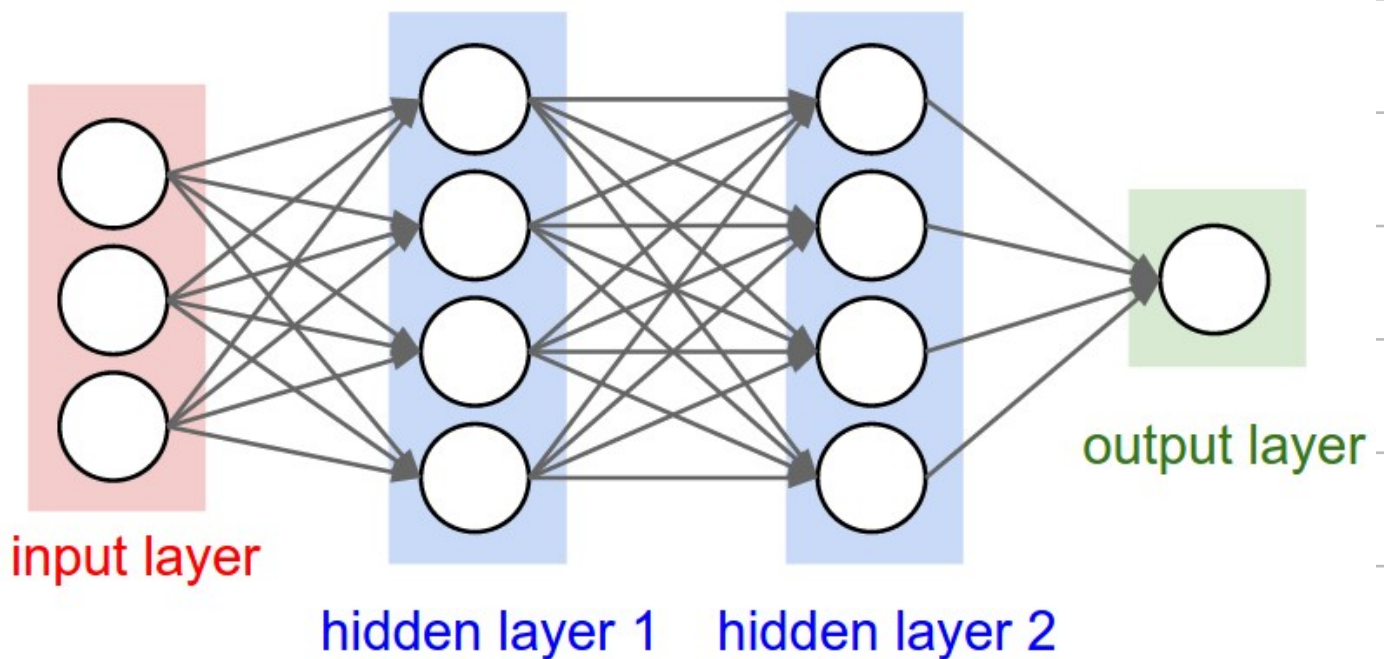


Neural Network :

- learns by taking inputs, mixing them with learned importance weights, and producing an output that improves with experience.
- many neurons connected together in layers.
- consists of nodes and connections between the nodes.



Activation Function :

- decides whether a neuron should activate

Without activation function

$$y = w_1 x_1 + w_2 x_2$$

no matter how many layers you add, the result is still linear

An activation function:

- takes the weighted sum from a neuron
- transforms it in a nonlinear way
- controls signal flowing forward

$$\text{output} = f(w \cdot x + b)$$

Weights decide how
imp. each input is
input value

Bias decides how
easy or hard it is
to activate

(handled internally)

Activation Functions

ReLU

$$f(x) = \max(0, x)$$

-ve \rightarrow ignored

+ve \rightarrow pass through

simple & fast

reduces vanishing

gradients

Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

Output: 0 - 1

Use case: probability

& binary classification

gradients vanish

for large

Tanh

$$f(x) = \tanh(x)$$

Output: -1 & 1



A gradient tells you how much model should change to reduce error.

$$\text{gradient} = \frac{\partial \text{loss}}{\partial \text{weight}}$$

Chain Rule:

raw weighted sum, $z = wx + b$ \rightarrow linear

activation function, $a = f(z)$ \rightarrow non-linear

loss associated, $L = \text{loss}(a)$

Chain, $w \rightarrow z \rightarrow a \rightarrow L$

Gradient, $\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}$

Question) Consider a single neuron NN defined as follows:

$$z = wx + b$$

$$a = \text{ReLU}(z)$$

where the loss for a single data point is given by the Mean Squared Error (MSE)

$$L = (a - y)^2$$

a) Derive the expression for the loss L in terms of w , b and x

$$z = wx + b$$

$$a = \text{ReLU}(z)$$

$$L = (a - y)^2$$

$$L = (\text{ReLU}(wx + b) - y)^2$$

(b) Compare the gradients w.r.t to weight and bias

$$\frac{\partial L}{\partial a} = 2(a - y)$$

$$\frac{\partial L}{\partial z} = 1$$

$$\frac{\partial z}{\partial w} = x$$

$$\frac{\partial z}{\partial b} = 1$$

gradient w.r.t weight,
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$
$$= 2(a - y) \cdot x$$

gradient w.r.t bias,
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b}$$
$$= 2(a - y)$$

(c) Explain how these gradients are used to update the parameters w and b during gradient descent.

$$w_{\text{new}} = w - \eta \frac{\partial L}{\partial w}$$

$$b_{\text{new}} = b - \eta \frac{\partial L}{\partial b}$$

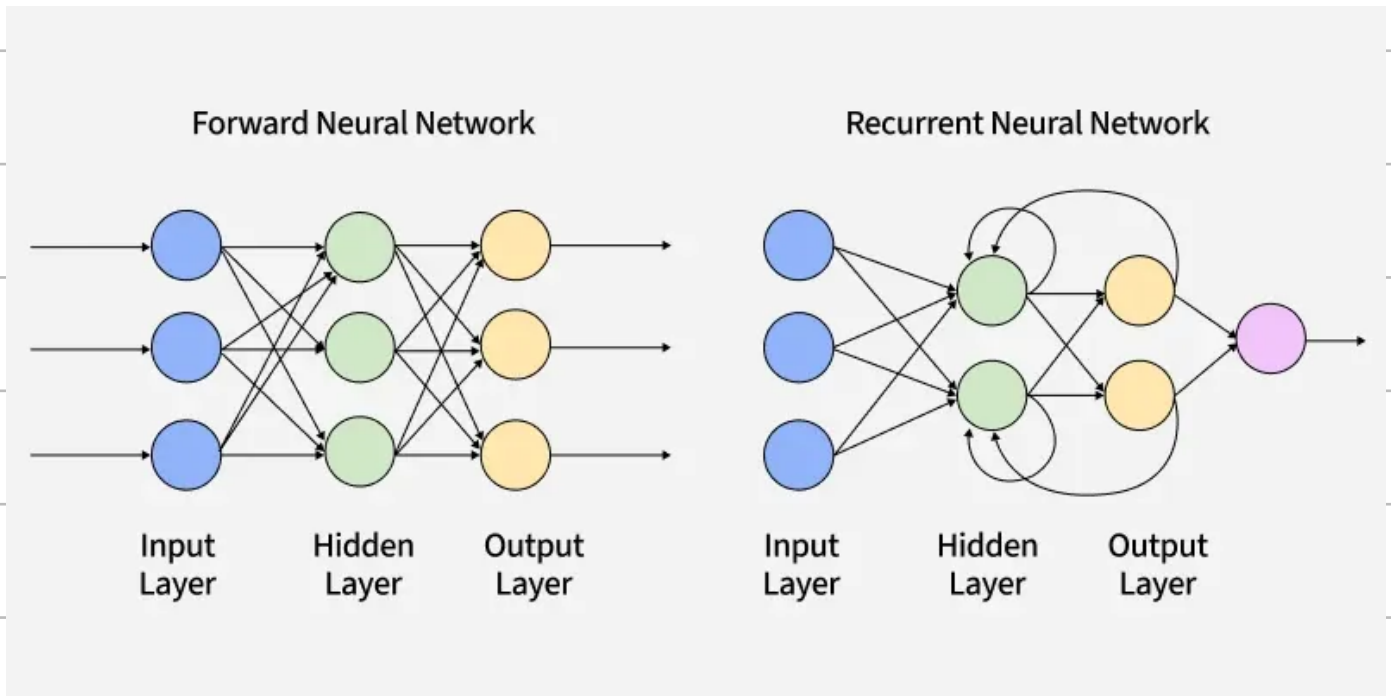
where η is the learning rate

Interpretation of a Gradient Value:

Consider a parameter θ with gradient:

$$\frac{\partial L}{\partial \theta}$$

Difference between Feed-Forward NN and Recurrent NN:



Sparse vs Dense Matrix:

Sparse

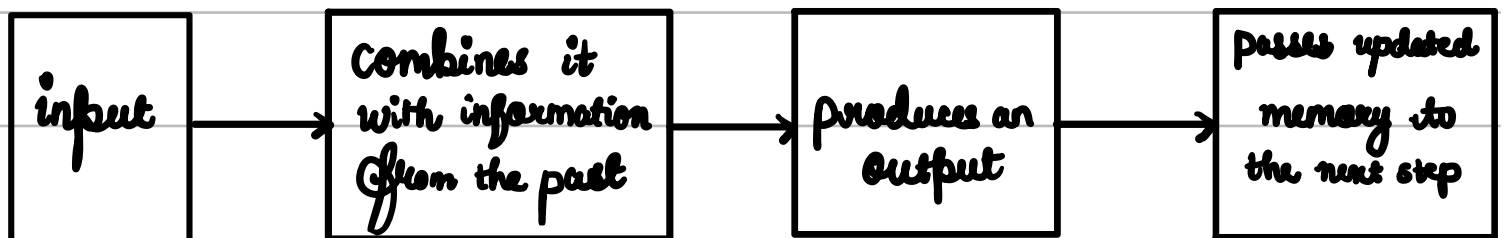
- Mostly zeros
- One hot encoded vectors
- low memory usage
- DL inputs
- no semantic meaning

Dense

- Packed with meaningful values
- NN weights
- High memory usage
- DL weights
- carries semantic meaning

Recurrent Neural Network :

- type of NN designed to work with sequences
- unlike standard neural network, it has memory.



Where RNNs are good at

- Capturing temporal patterns
- Learning trends & momentum
- Modeling sequences with short-term dependencies

Where RNNs struggle

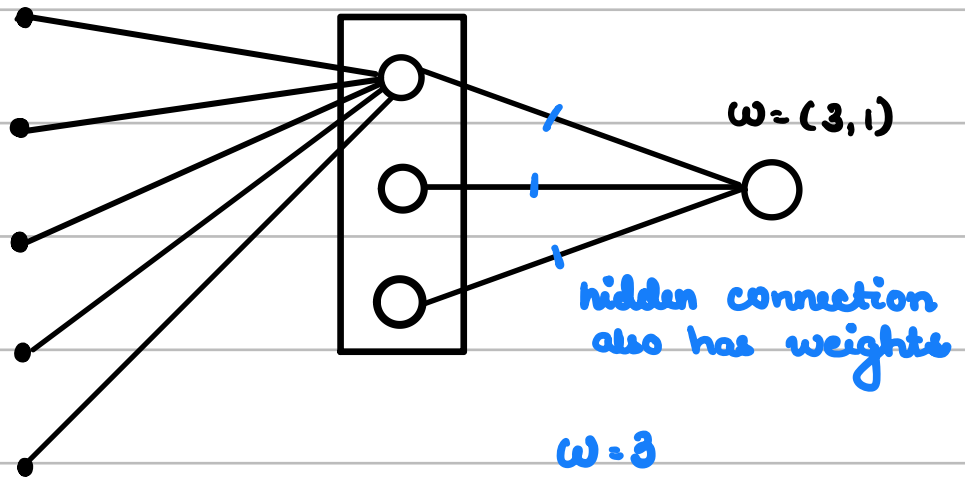
- Vanishing gradient (They forget long term dependencies)
- Over-smoothing (Predictions collapse toward averages)

RNN doesn't take the entire sequence at once, it takes one step at a time



Mechanism of vanilla RNN:

$$w = (5, 3)$$



weights

$$w = 5 \times 3$$

$$w = 3 \times 3$$

$$w = 3 \times 1$$

biases

$$b = 3$$

$$b = 1$$

Forward Propagation :

- the step where inputs flow through the neural network to generate output using current weights and biases

Suppose you have :

- input vector x
- weights w
- bias b
- activation function : $f(\cdot)$

For one layer, $z = Wx + b$ (linear transformation)

Output, $O_i = f(z)$

Output a becomes the input for next

Input			Output	Let x_{11} be 5 dimensional vector $[1\ 0\ 0\ 0\ 0]$ by One hot encoding (n features)
x_{11}	x_{12}	x_{13}	1	
x_{21}	x_{22}	x_{23}	0	
x_{31}	x_{32}	x_{33}	0	

Suppose you want 3 neurons. (n neurons)

Two weight matrices:

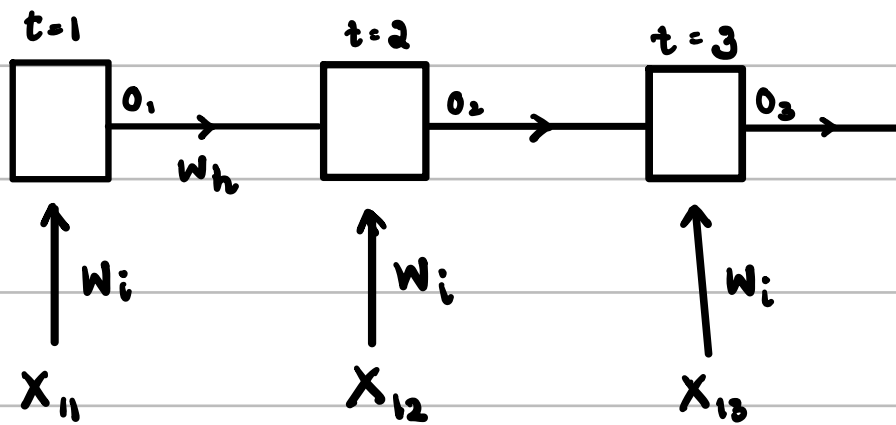
(1) W_i for input \rightarrow hidden

(2) W_h for hidden \rightarrow hidden

Weight dimensions are chosen to match your architectural choice

Activation
function

$$f(x_{11}w_n) \quad f(x_{12}w_i + O_1w_n) \quad f(x_{13}w_i + O_2w_n)$$



$$x_{11} \times w_i$$

$$x_{12} \times w_i + O_1 \times w_n$$

$$(1,5) \quad (5 \times 3)$$

$$(1,5) \quad (5,3) \quad (1,3) \quad (3 \times 3)$$

$$O_1 = (1 \times 3)$$

$$(1 \times 3) \quad (1 \times 3)$$

+
└──────────┘

$$O_2 = (1 \times 3)$$

To make the activation function consistent even at timestep 1, we need:

$$O_1 = f(x_1w_i + O_0w_n)$$

Common choices for O_0 :

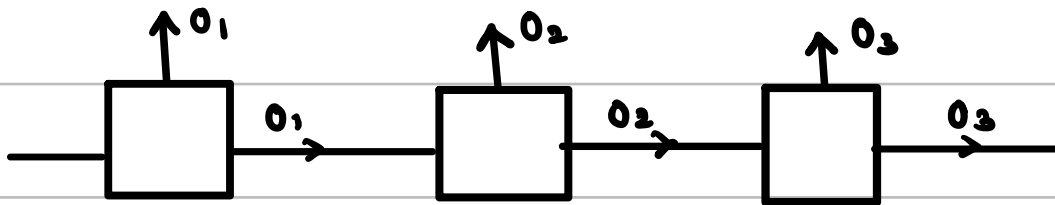
- Zero initialization (most common)
- Random (rare)
- Learned initial state

NOTE :

We use same weight matrix for each step.
So, at every step we are updating w_i and w_n

Q) Is there a way to check output at each node?

If you set `return_sequence = True`, you can also get output at each node.



Q) What if input size is not consistent?

We add extra values to data so everything has the same size. This technique is called

Padding

Loss Function :

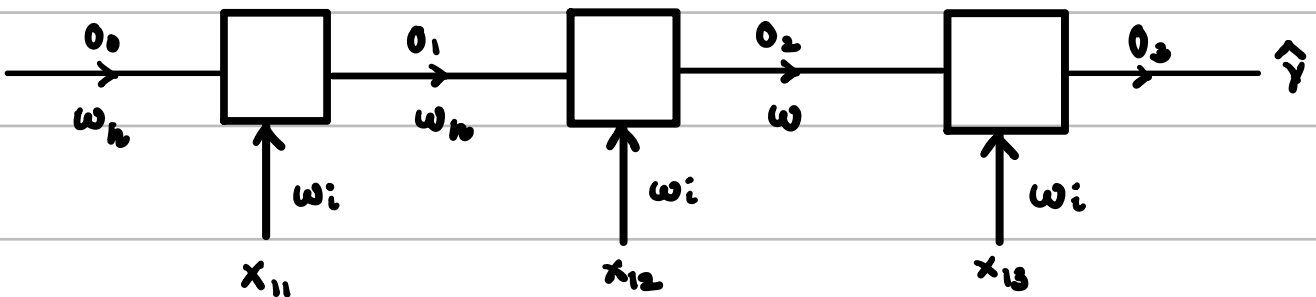
a mathematical function that measures how wrong a model's prediction is compared to the true value.

Prediction \rightarrow loss \rightarrow gradient \rightarrow Weight update

Name	Simple Definition	Formula	Where It Is Used
Mean Squared Error (MSE)	Average squared difference between actual and predicted values.	$L = \frac{1}{n} \sum (y - \hat{y})^2$	Regression
Mean Absolute Error (MAE)	Average absolute difference between actual and predicted values.	$L = \frac{1}{n} \sum y - \hat{y} $	Robust regression
Huber Loss	Combines MSE and MAE: Quadratic for small errors, linear for large errors.	$L = \frac{1}{2}(y - \hat{y})^2$ if $ y - \hat{y} \leq \delta$ else $\delta(y - \hat{y} - \frac{1}{2}\delta)$	Robust regression
Binary Cross Entropy (BCE)	Negative log likelihood for binary outcomes.	$L = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$	Binary classification
Categorical Cross Entropy	Negative log likelihood for multi-class outcomes.	$L = -\sum_k y_k \log(\hat{y}_k)$	Multi-class classification
Hinge Loss	Margin-based classification loss.	$L = \max(0, 1 - y\hat{y})$	Support Vector Machines
Negative Log Likelihood (NLL)	Log-based probabilistic loss.	$L = -\log P(y x)$	Probabilistic models

Backpropagation :

- A systematic way to compute gradients of the loss with respect to all weights using the chain rule



$$O_1 = f(x_{11} w_i + O_0 w_n)$$

$$O_2 = f(x_{12} w_i + O_1 w_n)$$

$$O_3 = f(x_{13} w_i + O_2 w_n)$$

$$\hat{y} = \sigma(O_3 w_0)$$

→ Prediction using sigmoid function

$$\text{loss} = \gamma - \hat{y}$$

$$= - \left[\gamma \log(\hat{y}) + (1 - \gamma) \log(1 - \hat{y}) \right]$$

Weights are calculated to minimize loss, so we use gradient descent.

$$w_i = w_i - \eta \frac{\partial L}{\partial w_i}$$

$$w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

$$w_0 = w_0 - \eta \frac{\partial L}{\partial w_0}$$

Since we are doing backpropagation,
we start with \hat{y}

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_0}$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{\partial}{\partial \hat{y}} \left[y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right] \quad \frac{\partial \hat{y}}{\partial w_0} = \frac{\partial \sigma(o_3 w_0)}{\partial w_0}$$

$$= -\left[\frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right]$$

$$= \frac{\partial}{\partial w_0} \left(\frac{1}{1 + e^{-o_3 w_0}} \right)$$

$$= -\frac{y - \cancel{y\hat{y}} - \hat{y} + \cancel{y\hat{y}}}{\hat{y}(1-\hat{y})}$$

$$= \frac{-1}{(1 + e^{-o_3 w_0})^2} \frac{\partial (1 + e^{-o_3 w_0})}{\partial w_0}$$

$$= \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$= \frac{-1}{(1 + e^{-0.3\omega_0})^2} \frac{\partial}{\partial \omega_0} (e^{-0.3\omega_0})$$

$$= \frac{-1}{(1 + e^{-0.3\omega_0})^2} (e^{-0.3\omega_0}) \frac{\partial}{\partial \omega_0} (-0.3\omega_0)$$

$$= \frac{e^{0.3\omega_0}}{(1 + e^{-0.3\omega_0})^2} \times (-0.3)$$

$$= \frac{-0.3 e^{0.3\omega_0}}{(1 + e^{-0.3\omega_0})^2}$$

$$= -0.3 \frac{1}{1 + e^{-0.3\omega_0}} \times \frac{e^{0.3\omega_0}}{(1 + e^{-0.3\omega_0})^2}$$

$$= -0.3 \hat{y} (1 - \hat{y})$$

$$\frac{\partial}{\partial \omega_0} \hat{y} = \frac{1}{1 + e^{-0.3\omega_0}}$$

$$1 - \hat{y} = 1 - \frac{1}{1 + e^{-0.3\omega_0}}$$

$$= \frac{1 + e^{-0.3\omega_0} - 1}{1 + e^{-0.3\omega_0}}$$

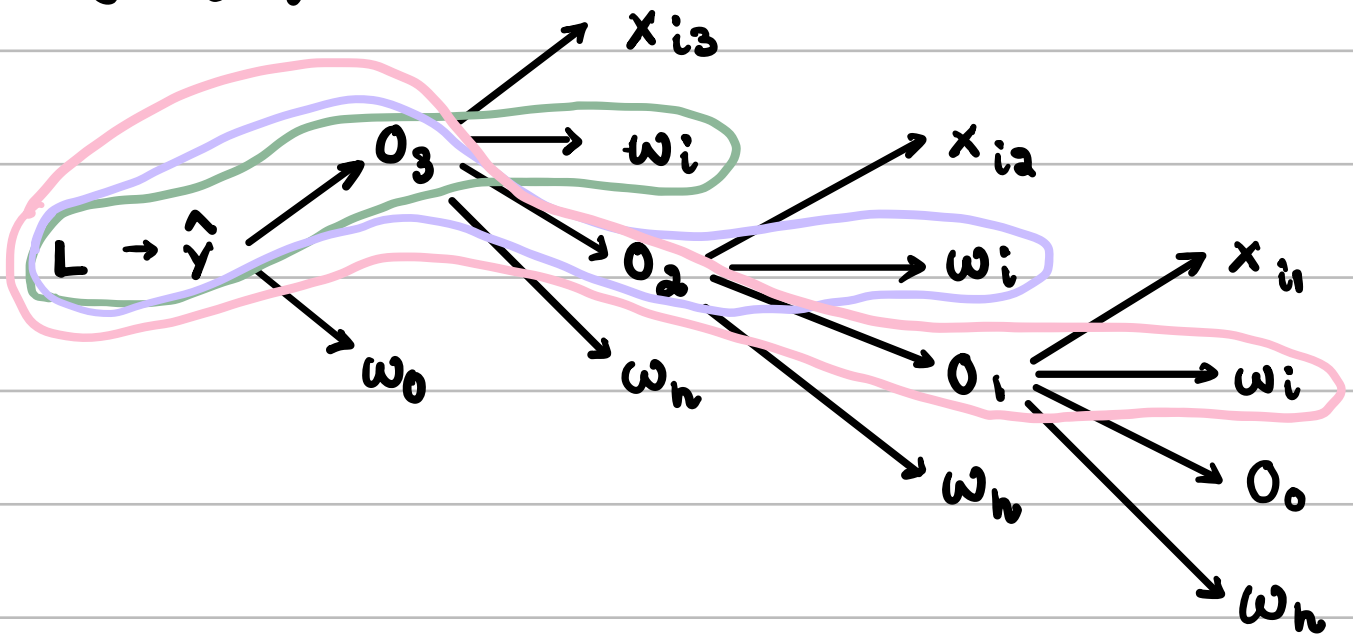
$$= \frac{e^{-0.3\omega_0}}{1 + e^{-0.3\omega_0}}$$

$$\frac{\partial L}{\partial \omega_0} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial \omega_0}$$

$$= \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \times -0.3 \hat{y} (1 - \hat{y})$$

$$= -o_3 (\hat{y} - y)$$

Finding $\frac{\partial L}{\partial \omega_i}$ by defining the independent paths



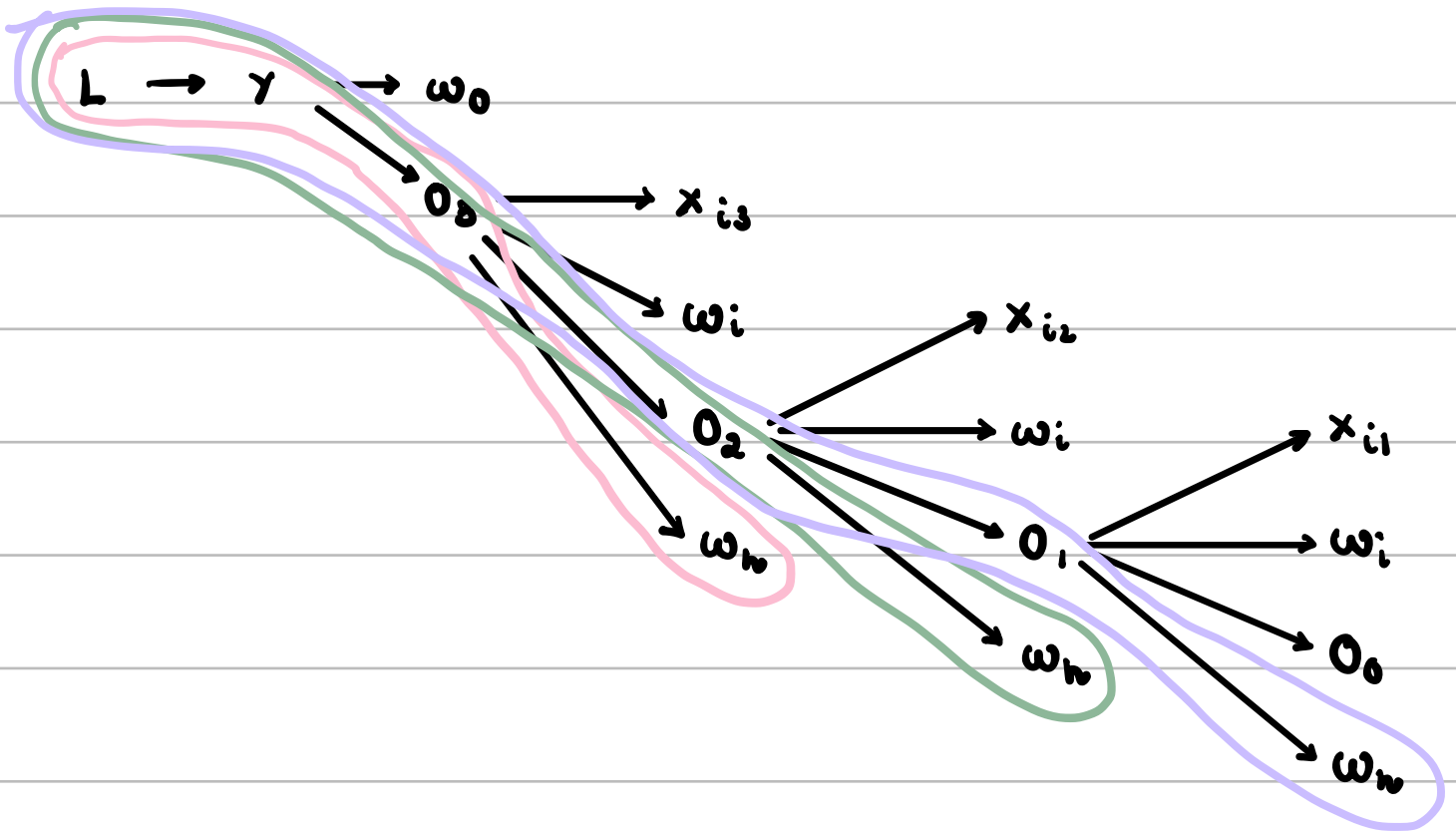
$$\frac{\partial L}{\partial \omega_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \omega_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial \omega_i} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial \omega_i}$$

For n timesteps

$$\frac{\partial L}{\partial \omega_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial o_j} \times \frac{\partial o_j}{\partial \omega_i}$$

Finding $\frac{\partial L}{\partial \omega_n}$,



$$\frac{\partial L}{\partial \omega_n} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial \omega_n} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial \omega_n} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial \omega_n}$$

For n steps,

$$\frac{\partial L}{\partial \omega_n} = \sum_{i=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial \omega_n}$$

RNN Summary :

Forward Propagation

Compute hidden states

$$O_1 = f(x_1 w_i + O_0 w_h)$$

$$O_2 = f(x_2 w_i + O_1 w_h)$$

$$O_3 = f(x_3 w_i + O_2 w_h)$$

Compute final prediction

$$\hat{y} = \sigma(O_3 w_o)$$

Compute loss

$$L = -[\gamma \log(\hat{y}) + (1-\gamma) \log(1-\hat{y})]$$

Update Weights

Choose a learning rate η

$$w_o \leftarrow w_o - \eta \frac{\partial L}{\partial w_o}$$

$$w_h \leftarrow w_h - \eta \frac{\partial L}{\partial w_h}$$

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i}$$

Back Propagation

Output weights : $\frac{\partial L}{\partial w_o}$

Hidden weights : $\frac{\partial L}{\partial w_h}$

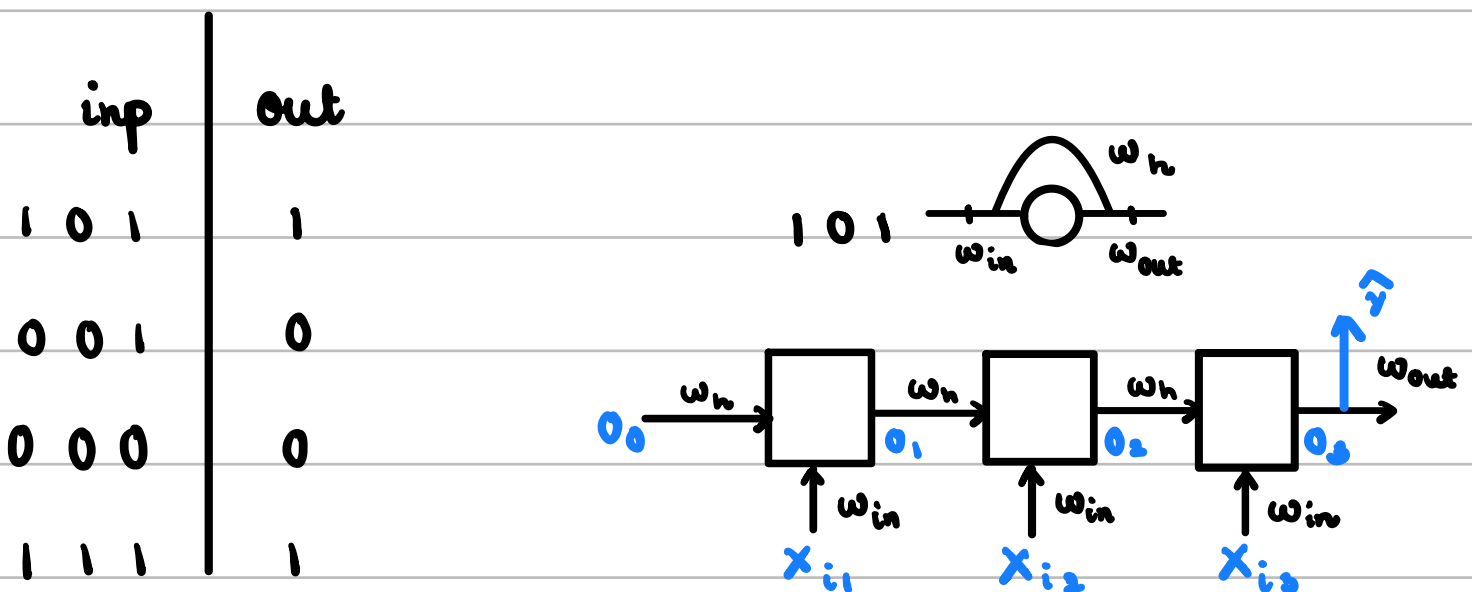
Input weights : $\frac{\partial L}{\partial w_i}$

Repeat

- Do this for many sequences
- Over many epochs
- Loss goes down
- Predictions improve

Drawback of RNNs :

(1) Vanishing and Exploding Gradient :



Since we are doing back propagation, we update weights to minimize the loss.

For 3 timesteps,

$$\frac{\partial L}{\partial w_{in}} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_{in}}} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial o_2} \frac{\partial o_2}{\partial w_{in}} +$$

short term memory

$$\boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w}}$$

long term memory

For n timesteps, long term memory will look like

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_n} \frac{\partial o_n}{\partial o_{n-1}} \dots \frac{\partial o_1}{\partial w_{in}}$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_n} \prod_{t=2}^n \left(\frac{\partial o_t}{\partial o_{t-1}} \right) \frac{\partial o_1}{\partial w_{in}}$$

Here $O_t = \tanh(x_{it} w_{in} + O_{t-1} w_h)$

$$\frac{\partial O_t}{\partial O_{t-1}} = \tanh' (x_{it} w_{in} + O_{t-1} w_h) w_h$$

Substituting,

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_n} \prod_{t=2}^n \underbrace{\tanh'(x_{it} w_{in} + o_{t-1} w_h)}_{0-1} \underbrace{w_h}_{0-1} \frac{\partial o_1}{\partial w_{in}}$$

≈ 0

As timesteps increase, repeated multiplication during backprop causes gradients to decay exponentially, preventing the model from learning long-term dependencies.

(2) Sequential Computation:

$$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots$$

each step depends on the previous one.

So,

- cannot fully parallelize
- Training is slow on GPU