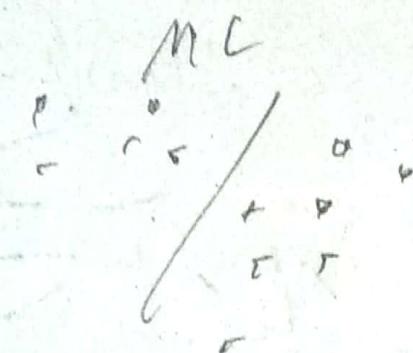


the people rather than freight all its hope on
of just one person.

Enter to NN

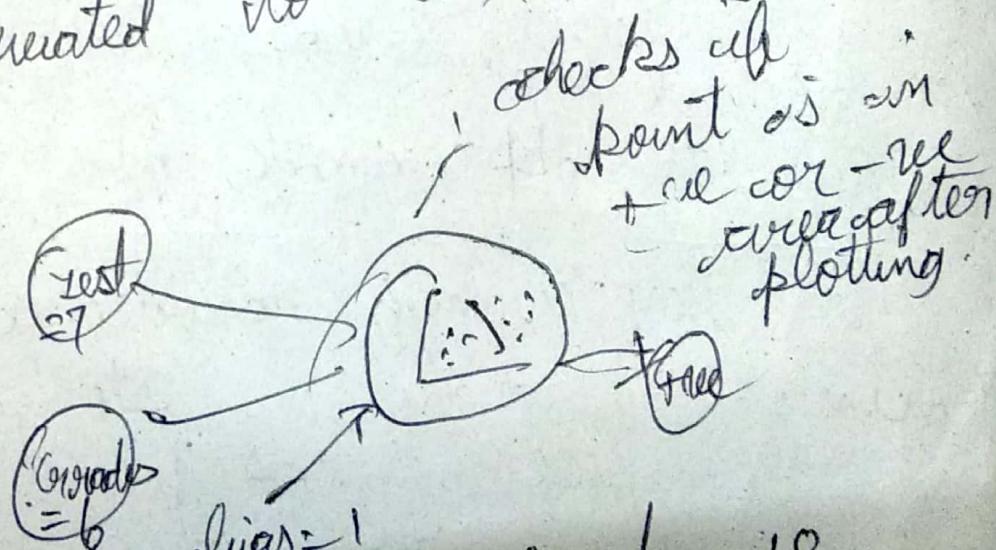
Neural network



n -dimensional space, boundary is a $n-1$ dim hyperplane which is a high dimensional line of line w/ plane in 2D & 3D respectively.

$w_1x_1 + w_2x_2 - w_nx_n + b = 0$ & at
 $w_1x_1 + w_2x_2 - w_nx_n + b = 0$.
can be abbreviated to $wx + b = 0$.

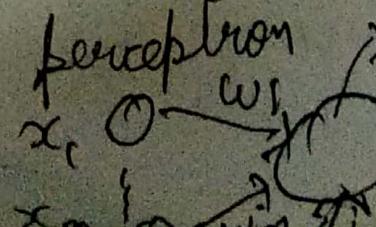
Perception :



$$2 \times \text{Test} + 1 \times \text{grades} - 18$$

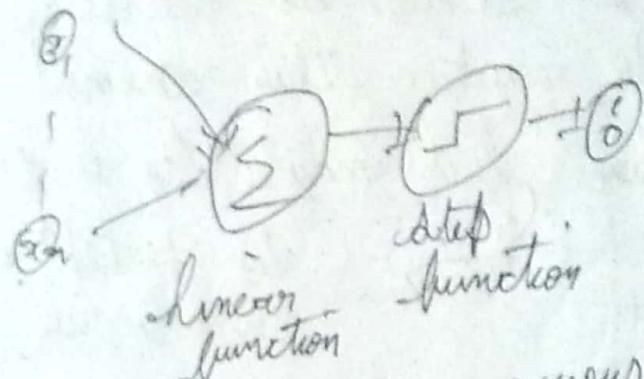
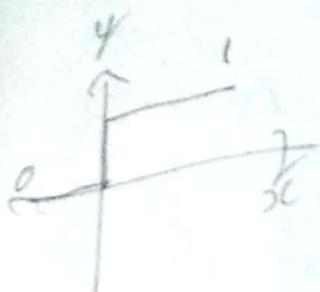
$$wx + b = \sum_{i=1}^n w_i x_i + b$$

Generic case of perceptron

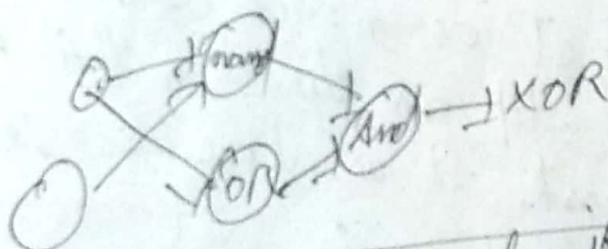


do here we use a step function i.e.

$$y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



nucleus → axon. Neurons takes nervous impulses as input through dendrites & output them using axon.



Perception algorithm

1) Start with random weights:

$$w_1, w_2, \dots, w_m, b$$

4 5 \rightarrow^{10}
2) For every misclassified point:

$$(x_1, \dots, x_n)$$

If prediction is 0

for $i=1 \dots n$

change $w_i + \alpha x_i$

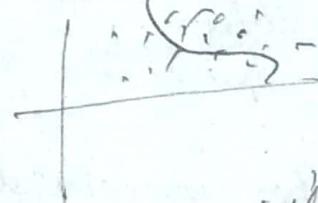
change $b + \alpha$

If prediction is 1:

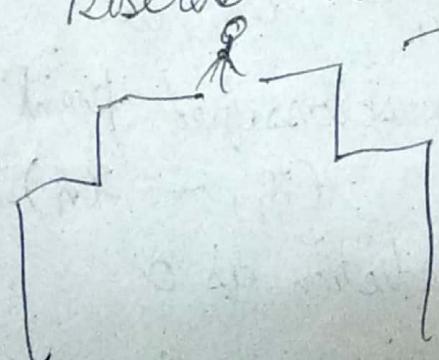
$$\text{for } i=1 \dots n$$

change $w_i - \alpha x_i$

change $b - \alpha$

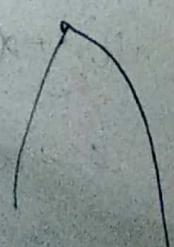
Non-linear regions -
What if grade is 7 but test is 9. This
would be classified as positive but we
do not want to accept a person with
low test scores. So we can classify
it to be negative. This cannot be done
by a line but maybe by a curve -

So perception won't
work & we have to
generalize so it'll work for other types
of curves.

Log-loss error function: Mountain
example - To get down we look at
direction & take step in direction that helps
us descend the most. It is possible to get
stuck at local minimum.
Discrete v/s continuous



here we always
see constant errors &
we are confused &
not sure what to
do.

But in a continuous scenario
one can detect small
variations in height & radar
in figure in what direction
it happened the most.

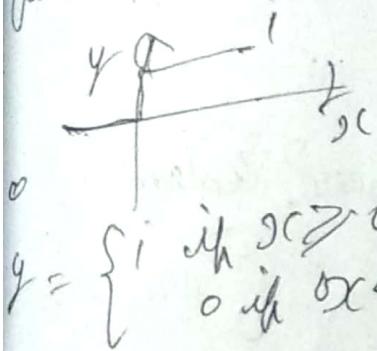


do our error function cannot be discrete but continuous (i.e. differentiable).

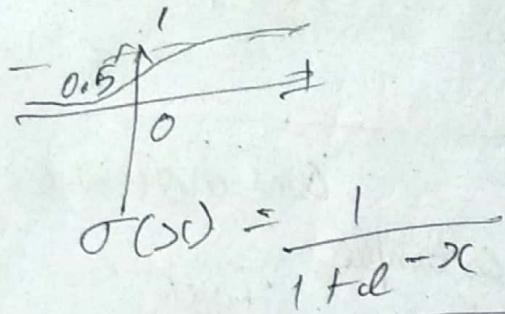
discrete vs continuous? e.g. discrete - yes or answer no

continuous - a number probably between 0 to 1. (continuous?).

A simple way to move from discrete to continuous predictions is to shift from step function to sigmoid activation function.



$$y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



Multi-class classification

It is equivalent to sigmoid function but for problems with 3 or more classes.

\exp (exponential) is a function that returns a real number for every input even for all numbers.

score

bear

2

duck

1

mechanus

0

number by sum of all numbers. e.g. $\frac{2}{2+1+0} = \frac{2}{3}$

or $\frac{1}{2+1+0}$

. But there's a probability

that scores well - & hence could lead to decisions by a error. So instead we use exp as it turns all numbers to +ve numbers. So

$$\frac{e^z}{e^z + e^{z_1} + \dots} - \text{This gives us the softmax function.}$$

For N classes & a linear model that
gives us scores.

$$P(\text{class}_i) = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

one hot encoding			
example	Duck	Beaver?	Walrus?
	1	0	0
Beaver	0	1	0
Walrus	0	0	1

Maximum likelihood: usually the best model is the one that assigns higher probabilities to events that happened to us. This method is called Maximum likelihood.

Maximizing probabilities using math
Maximizing the probability is equivalent to minimizing the error function.

For finding the best model we consider pixels of colors as independent events & multiply them. The product that yields

highest prob is the best one. However products of millions of data points are dead. Instead it can be converted to a sum problem using log.

Cross-entropy

$\log(1) = 0$ so $\log(\text{number between 0 to 1}) \approx -ve$.
so if we take $-ve$ of log of probabilities we get a $+ve$ number.

$$\log \frac{1}{6} = \ln(0.6) = -0.51 \text{ but } -\ln(0.6) = 0.61.$$

Sum up - we do log of numbers is called cross entropy. Good algo - fast low cross entropy.

bad algo - low high cross entropy.

Cross entropy is simple. If we have a bunch of events & probability, how likely is that the events happen based on the probabilities. If events happen based on the probabilities - its very likely we have low cross entropy. less likely? high entropy.

	Door 1	Door 2	Door 3	Actual?
P(gift)	0.8	0.7	0.1	= expog
P(wrong)	0.2	0.3	0.9	50%

Events with high prob - have low cross entropy & with low prob have high "

$$\begin{array}{ccc} \text{Door 1} & \text{Door 2} & \text{Door 3} \\ 0.8 & 0.7 & 0.9 \end{array} \quad \begin{array}{c} \text{right} \\ x - \text{no right} \end{array}$$

$$-\ln(0.8) - \ln(0.7) \\ -\ln(0.9)$$

$y_i = 1$ if present on door i

$$\text{Door 1 } y_1 = 1$$

$$p(\text{right}) =$$

$$p_1 = 0.8$$

$$\text{Door 2 } y_2 = 1$$

$$p_2 = 0.7$$

$$\text{Door 3 } y_3 = 0$$

$$p_3 = 1 - p_2 = 0.3$$

$$\text{So formula cross entropy} = - \sum_{i=1}^m y_i \cdot \ln(p_i) + (1-y_i) \cdot \ln(1-p_i)$$

Cross entropy tells us how similar are the vectors.

Multi-class cross entropy

Animal	Door 1	Door 2	Door 3	The numbers in columns should add to 1.
Duck	0.7	0.3	0.1	
Beaver	0.2	0.4	0.5	
Walrus	0.1	0.3	0.6	

nos. in rows need not add to 1.

Door 1 Door 2 Door 3

Duck Walrus Walrus

$$P = 0.7 + 0.3 + 0.4 = 1.4$$

$$CE = -\log(0.7) + -\ln(0.3) + -\ln(0.4) = 1.48$$

$y_{ij} = 1$ if Duck behind door j

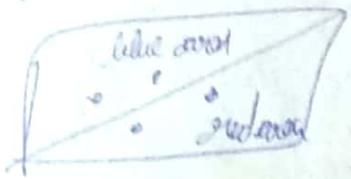
$y_{ij} = 1$ if Beamer

$y_{ij} = 1$ if weel

$$\text{Cross entropy} = -\sum_{i=1}^n \sum_{j=1}^m y_{ij} \ln(\hat{y}_{ij})$$

m - number of classes.

Logistic regression:



$$\begin{aligned} \text{If } y = 1 \\ P(\text{blue}) &= y^1 \\ \text{error} &= -\ln(y^1) \end{aligned}$$

$$\begin{aligned} \text{If } y = 0 \text{ (point is red)} \\ P(\text{red}) &= 1 - P(\text{blue}) \\ &= 1 - y^1. \end{aligned}$$

$$\text{error} = -\ln(1 - y^1)$$

$$\text{So error} = -(1 - y^1) \ln(1 - y^1) - y^1 \ln(y^1)$$

$$\text{error function} = \frac{-1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \hat{y}_i) + y_i \ln(\hat{y}_i)$$

now we considering average.

If activation is sigmoid then our error function should be in terms of w & b .

$$E(w, b) = \frac{-1}{m} \sum_{i=1}^m [-(y_i) \ln(1 - \sigma(w^T x^{(i)}) + b) + y_i \ln(\sigma(w^T x^{(i)}) + b)]$$

Above is for binary class. For multiclass it is

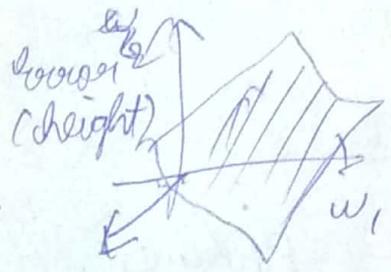
$$\frac{-1}{m} \sum_{i=1}^m \sum_{j=1}^k y_{ij} \ln(\hat{y}_{ij})$$

Minimizing error function

$$E(w, b) = -\frac{1}{m} \sum_{i=1}^m (y_i \ln(\sigma(wx^{(i)} + b)) + (1-y_i) \ln(1-\sigma(wx^{(i)} + b)))$$

we update weights & biases such that we get reduced error $E(w', b')$.

gradient descent



∇E = vector sum of gradient partial derivatives of E w.r.t w_1, w_2

∇E gives as the direction in which error increases more. So taking -ve of ∇E tells as in which direction we must move to decrease the error.

$$\begin{aligned} \frac{\partial E}{\partial w_2} \uparrow & \quad \nabla E \quad g = \sigma(w_1 x_1 + b) \\ \downarrow \frac{\partial E}{\partial w_1} & \quad g = \sigma(w_1 x_1 + w_2 x_2 + b) \\ -\nabla E & \quad \nabla E = (\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}) \end{aligned}$$

$\alpha = 0.1$ Taking a step is same as updating w & b .

$$w_i' \leftarrow w_i - \alpha \frac{\partial E}{\partial w_i}$$

$$b' \leftarrow b - \alpha \frac{\partial E}{\partial b}$$

$$g = \sigma(w_1 x_1 + w_2 x_2 + b) - \text{target}$$

$$\sigma^2(x) = \frac{\partial}{\partial x} \left(\frac{1}{1+d^{-x}} \right) = \frac{d^{-x}}{(1+d^{-x})^2}$$

$$= \frac{d^{-x}}{(1+d^{-x})^2} = \frac{1}{1+d^{-x}} \cdot \frac{1}{(1+d^{-x})^2}$$

$$= \frac{1}{(1+d^{-x})} \cdot \frac{d^{-x} + 1 - 1}{(1+d^{-x})} = \sigma(1-\sigma)$$

For m points
labelled $x^{(1)} \rightarrow x^{(m)}$

$$E = \sigma \cdot \left(1 - \frac{1}{1+d^{-x}} \right) = \sigma(1-\sigma)$$

$$E = \frac{1}{m} \sum_{i=1}^m y_i \ln(g_i) + (1-y_i) \ln(\bar{g}_i)$$

where $\bar{g}_i = \sigma(wx_i + b)$

$$\nabla E = \left(\frac{\partial}{\partial w_1} E, \dots, \frac{\partial}{\partial w_n} E, \frac{\partial}{\partial b} E \right)$$

To simplify

we think of error each point produces &
calculate derivative of error. The total error
is average of error at all these points.

Error at each point is: $E = y \ln g - (1-y) \ln(1-g)$

$$\frac{\partial}{\partial w_j} E = \frac{\partial}{\partial w_j} \sigma(wx+b) = \sigma'(wx+b)(1-\sigma(wx+b))$$

$$= g(1-g) \frac{\partial}{\partial w_j} (wx+b)$$

$$= g(1-g) \cdot \frac{\partial}{\partial w_j} (wx+b) = g(1-g) \frac{\partial}{\partial w_j} (wx_1 - \frac{w}{m} x_m + db)$$

$$= g(1-g) x_j$$

i.e. $w_j x_j$ is the only term not constant w.r.t w_j .

$$\therefore \frac{\partial E}{\partial w_j} = \frac{\partial}{\partial w_j} [-g \ln y - (1-y) \ln(1-y)]$$

$$= \left[-y \frac{1}{y} \frac{\partial}{\partial w_j} y + (1-y) \frac{\partial}{\partial w_j} \ln(1-y) \right]$$

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

$$= -y(1-y)x_j + y^2(1-y)x_j$$

$$= -y x_j + y^2 x_j$$

$$= -(y - y^2) x_j$$

$$\therefore \frac{\partial E}{\partial b} = -(y - y^2) \text{ (similar calculation)}$$

In summary gradient is $\nabla E = -(y - y^2)(x_0, -x_1)$

Gradient descent step :

$$w_i' = w_i - \alpha (-(y - y^2)x_i)$$

$$\Rightarrow w_i' = w_i + \alpha (y - y^2)x_i$$

$$b' = b + \alpha (y - y^2)$$

logistic reg algorithm

Pseudo code in next page

1) start with random weights:
 w_1, \dots, w_m, b

2) for every point (x_1, \dots, x_n) :

2.1 For $i=1$ to n

update $w_i' \leftarrow w_i - \alpha (g_i - y_i)$

$$g_i' \leftarrow g_i - \alpha (y_i - g_i)$$

repeat until error is small.

Perception vs gradient descent

gradient descent

change

$$w_i \leftarrow w_i + \alpha (g_i - y_i) x_i$$

Perception algo-

In perception
not every
point changes
the weight
only misclassified
ones

If x is classified

change w_i to $w_i + \alpha x_i$
(if +ve)

$$w_i - \alpha x_i$$

(if -ve)

If correctly classified

$$y - g = 0$$

If misclassified,

$$y - g = 1 \text{ if positive}$$

$$y - g = -1 \text{ if negative}$$

Are both same?

The main diff is that
in gradient descent y can take
any values.

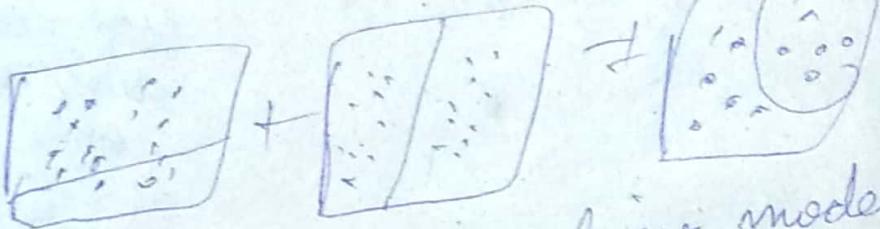
If a point is correctly classified perceptron algorithm stops. But an incorrect descent a misclassified point tells the line to come closer & a correctly classified point tells it to go farther away.

Von linear regions

Line set of points equally likely. Here is where neural networks

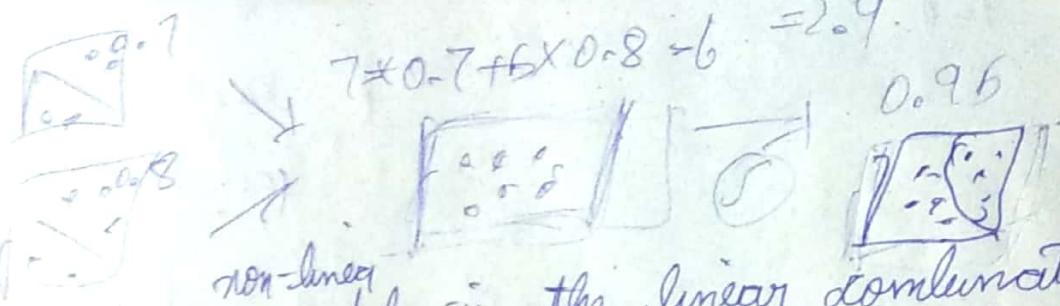
come into play.

Combining regions: we are gonna combine 2 perceptions into a 3rd complicated one.

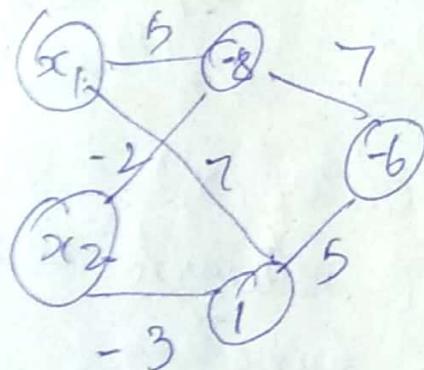


we are gonna combine 2 linear models to non linear model. we combine as follows: consider a point in blue region in 1st diag alone. Its prob lets say is 0.8 & prob of a point in 2nd diag in blue space is 0.7. When we combine both we get $0.5 > 1$. To squash it between 0 to 1, we can use sigmoid so $G(0.7, 0.5) = 0.81$

If we want we can add weights too to the 2 linear models -



The third model is the linear combination of the existing 2 models.



$$\frac{1}{1+e^{-x}} = 0.88$$

$$0.88 + 0.88 e^{-x} = 1$$

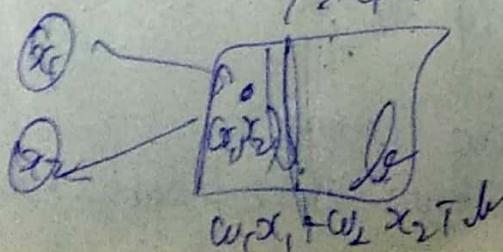
$$0.88 e^{-x} = \frac{1}{0.88}$$

$$-60.01(1/e) = -60.107$$

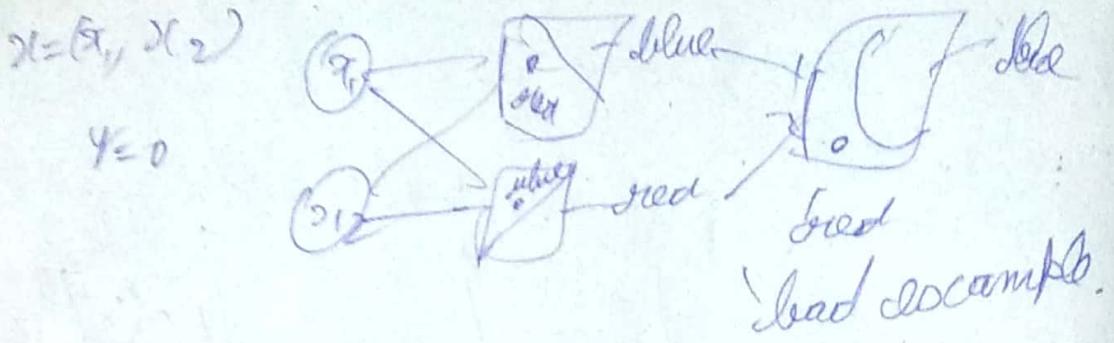
models.

Feedforward

Feedforward is the process neural networks use to turn input to output.
→ plots & outputs probability.



we can keep adding more layers & nodes to form complex networks. PNN's are where linear models combine to form non-linear models which further combine to form even more non-linear



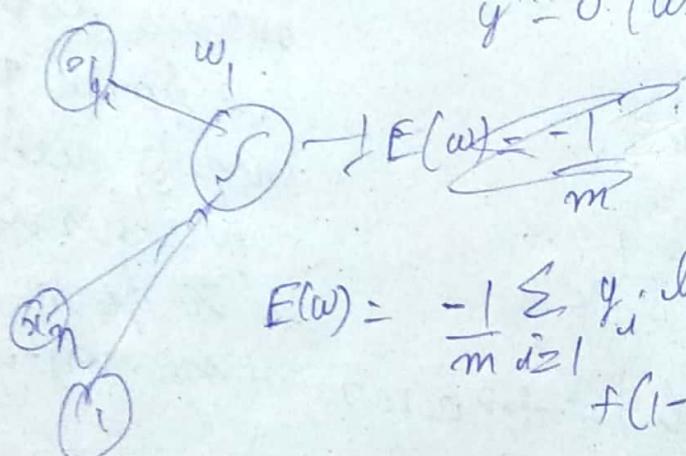
For a neural net with 2 layers

$$y^2 = \sigma \begin{pmatrix} w_{11}^{(2)} \\ w_{21}^{(2)} \\ w_{31}^{(2)} \end{pmatrix} \sigma \begin{pmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

$$y^2 = \sigma w_2 \sigma w_1 (x)$$

To train neural network we need
to define error function

$$E(w) = \sigma(wx + b)$$



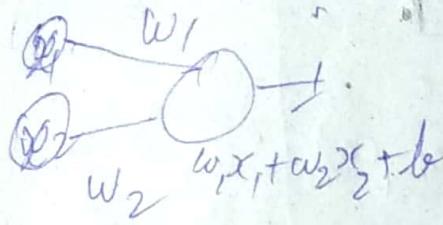
Backpropagation

In a nutshell backprop consists of:

- going in feedforward operation.
- comparing output of model with desired output.

In feedforward calculating error
in feedforward up backwards to spread error

- to each of the weights
- use this to update weights
- continue till model is good.



$$E(\omega) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(g_i) + (1-y_i) \ln(1-g_i)$$

For MCP to reduce error. Calculate $E(\omega)$, & walk in direction of $-\nabla E$ to minimize the error.

Backprop math

$$g = \sigma(wx + b)$$

$$E(\omega) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(g_i) + (1-y_i) \ln(1-g_i)$$

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_m}, \frac{\partial E}{\partial b} \right)$$

Here g
In MCP. $g = \sigma w^{(3)} \circ \sigma w^{(2)} \circ w^{(1)} x$
3 layers error func is same but

∇E is much more complex.

$$\nabla E = \left(\frac{\partial E}{\partial w^{(1)}}, \dots, \frac{\partial E}{\partial w^{(l)}} \right)$$

More formally for a MCP with l layers

$$g = \sigma w^{(l)} \circ \sigma w^{(l-1)} \circ \dots \circ \sigma w^{(1)} x$$

$$w^{(1)} = \begin{pmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{pmatrix}$$

$$w^{(2)} = \begin{pmatrix} w_{11}^{(2)} \\ w_{21}^{(2)} \\ w_{31}^{(2)} \end{pmatrix}$$

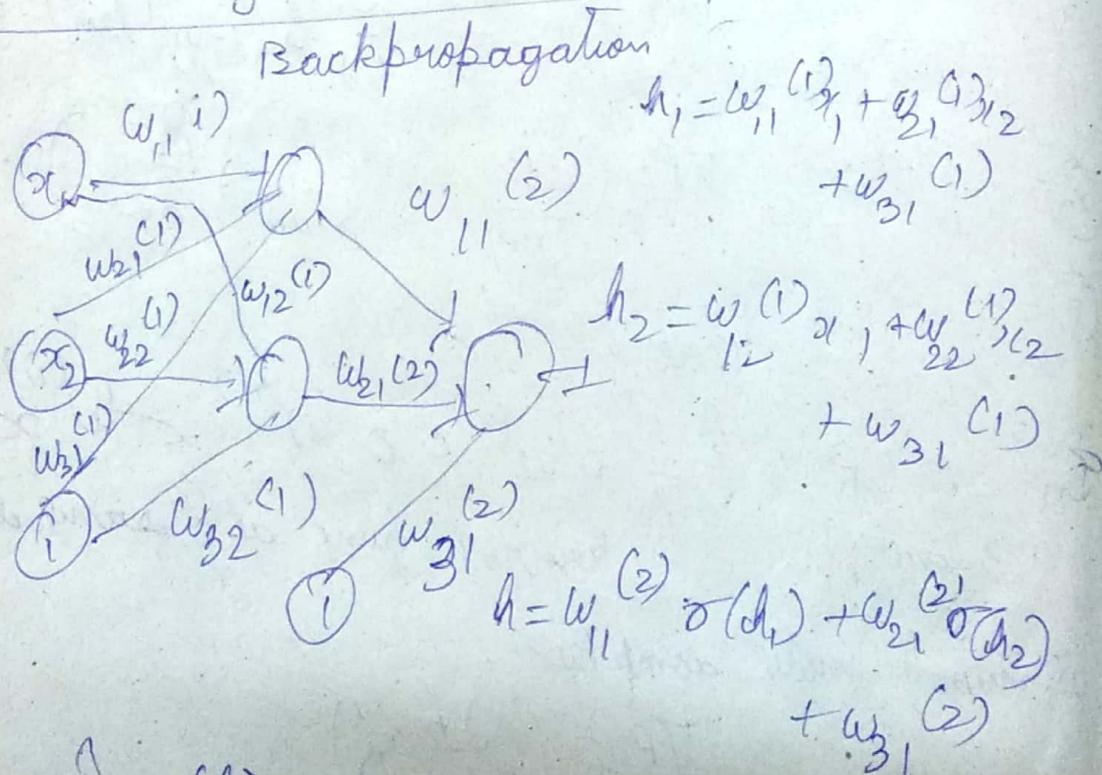
$$\nabla E = \begin{pmatrix} \frac{\partial E}{\partial w_{11}} & \frac{\partial E}{\partial w_{12}} & \frac{\partial E}{\partial w_{13}} \\ \frac{\partial E}{\partial w_{21}} & \frac{\partial E}{\partial w_{22}} & \frac{\partial E}{\partial w_{23}} \\ \frac{\partial E}{\partial w_{31}} & \frac{\partial E}{\partial w_{32}} & \frac{\partial E}{\partial w_{33}} \end{pmatrix}$$

Chain rule
 $\textcircled{1} \rightarrow \textcircled{2} \rightarrow \textcircled{3}$
 $B = g(A)$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial w_{ij}} \cdot \frac{\partial w_{ij}}{\partial A} = \frac{\partial B}{\partial C} = \frac{\partial B}{\partial A} \frac{\partial A}{\partial C}$$

$w_{ij}^{(k)}$ & $w_{ij}^{(k)} - \frac{\partial E}{\partial w_{ij}^{(k)}}$ when functions

compose derivatives are multiplied. Backprop is same as function composition & backprop is calculating derivatives & multiplication.



$$g = \sigma(h)$$

i.e. $g = \sigma(w^T \odot \sigma \odot w(x))$.

$$E(w) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(g_i) + (1-y_i) \ln(1-g_i)$$

Now we need to calculate derivative

of error function w.r.t. each of the weights.
using chain rule.

$$E(\omega) = E(\omega_{11}^{(1)}, \dots, \omega_{31}^{(2)})$$

$$\frac{\partial E}{\partial \omega_{11}^{(1)}} = \left(\frac{\partial E}{\partial \omega_{11}^{(1)}} - \frac{\partial E}{\partial \omega_{31}^{(2)}} \right)$$

$$\frac{\partial E}{\partial \omega_{11}^{(1)}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial h_1} \frac{\partial h_1}{\partial \omega_{11}^{(1)}} \quad \text{as}$$

feedforwarded is just composition of functions
so we take a small part of MLP

$$\begin{array}{l} \textcircled{1} \xrightarrow{\omega_{11}^{(1)}} h_1 = \omega_{11}^{(2)} \sigma(h_1) + \omega_{12}^{(2)} \sigma(h_2) \\ \textcircled{2} \xrightarrow{\omega_{21}^{(2)}} h_2 = \omega_{21}^{(2)} \sigma(h_1)(1 - \sigma(h_1)) \\ \textcircled{1} \xrightarrow{\omega_{31}^{(2)}} h_3 = \omega_{31}^{(2)} \sigma(h_1)(1 - \sigma(h_1)) \end{array}$$

MSE vs dog. loss -

$$\text{Sum of squared errors: } E = \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^m (y_i - g_i)^2$$

μ - set of all data points. y - output units.

First we perform inside sum over j . For each output unit we find diff between y_j & g_j sq & add

Second - for each data point find inner sum of squared diff for each output unit -

Then sum of the squared diff for each data point. \Rightarrow Bias: 1) errors are always +ve
2) larger errors more penalized

$$y_j^{(t)} = f(\sum_i w_{ij} x_i^{(t)})$$

$$E = \frac{1}{2} \sum_{\mu} \sum_j [y_j^{(t)} - f(\sum_i w_{ij} x_i^{(t)})]^2$$

We want the network to error to be as small as possible & weights acts as variables.

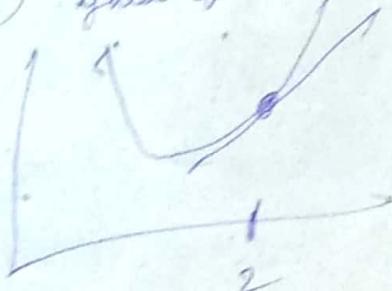
$$\hat{y}_j^{(t)} = f(\sum_i w_{ij} x_i^{(t)})$$

We want to step in direction that minimizes the error the most. This direction can be found by calculating gradient of squared error. Gradient is rate of change. To calculate rate of change we turn to derivative. Derivative of a function $f(x)$ gives $f'(x)$ that returns slope of $f(x)$ at x .

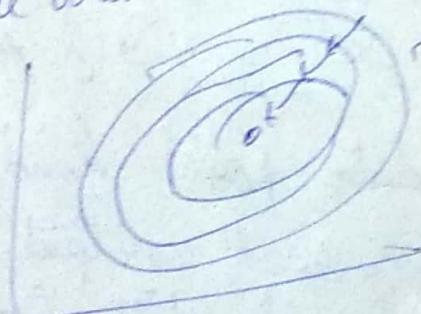
$$f(x) = x^2$$

$$f'(x) = 2x$$

$$f'(2) = 4$$



Gradient is first derivative extended to functions with more than one variable.



- Topographical map whose points on a contour don't have same depth & darker contours have

larger errors.

: weights will go wherever the error is less, so we will get stuck at local minimum, so we can add momentum to

overcome it.] In MSE why don't we take $\|g - \hat{g}\|_2^2$ instead of square - One advantage of squaring is that it penalizes outliers more.

Gradient descent of SSE: math

$$E = \frac{1}{2} \sum_{\mu} (y^\mu - \hat{g}^\mu)^2 = \frac{1}{2} \sum_{\mu} (y^\mu - f(\sum_i w_i x_i))$$

w-data - we can think of scanning the below matrix one row at a time, calculating diff. & summing them up. Now let us consider just one example -

$$E = \frac{1}{2} (y - f(\sum_i w_i x_i))^2$$

↑
↓ Δw = gradient direction.

$$w_j = w_j + \Delta w_j$$

$$\Delta w_j \propto -\frac{\partial E}{\partial w_j}$$

$$\Delta w_j = -n \frac{\partial E}{\partial w_j}$$

Learning rate.

$$\begin{aligned} \frac{\partial E}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} (y - \hat{g})^2 = \frac{\partial}{\partial w_j} \frac{1}{2} (y - g(w_j))^2 \\ &= (y - g) \frac{\partial}{\partial w_j} (y - g) \\ &= -(y - g) \frac{\partial g}{\partial w_j} \end{aligned}$$

(using chain rule)

But $g = \phi(h)$ where $h = \sum_i w_i x_i$

$$\frac{\partial E}{\partial w_j} = -(y - g) \frac{\partial g}{\partial h} \frac{\partial h}{\partial w_j}$$

Applying chain rule we get

$$= \frac{\partial}{\partial h} \left(\frac{\partial g}{\partial h} \right)$$

$$= (-y - \hat{y}) \cdot g'(w_d) \frac{\partial}{\partial w_d} \sum a_i^* x_i^* \quad \text{as } a_i^* = w_i^* x_i^*$$

$$\frac{\partial}{\partial w_d} \sum a_i^* x_i^* = x_i^* \quad \therefore \frac{\partial E}{\partial w_d} = -(y - \hat{y}) g'(w_d) x_i^*$$

$$w_d = a_d + \eta (y - \hat{y}) g'(w_d) x_d^* \\ \text{as } \Delta w = -\text{gradient}$$

General algo for updating with gradient descent
Get weights to zero.

For each record in training data:

- Make a forward pass through network $\hat{y} = f(\sum a_i x_i)$

- Calculate error for output and $\delta = (y - \hat{y}) h(x)$

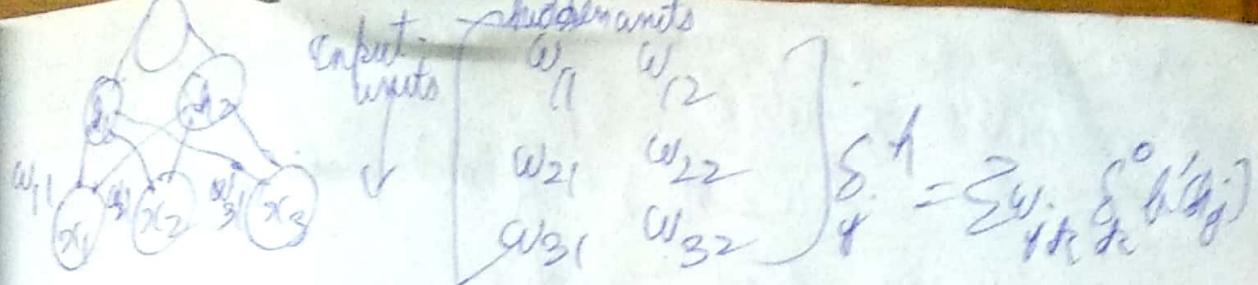
- Update weight step, $\Delta w_i = \eta a_i \delta x_i$.

- Update the weights $w_i = w_i + \eta \Delta w_i / m$ where m

m is no. of records that enables to reduce large variations in training data.

- Repeat for e epochs.

Initializing weights - we have to initialize them to be small so they are in linear region & not squashed at high and low ends. Initialization should also be random to avoid symmetry. A good value for scale is \sqrt{Nn} , n - number of input units.



$\Delta w_{dkj} = \eta \delta_j^h x_i$: w_{dkj} weights between input and hidden layer.

$$\Delta w_{dkj} = \eta \delta_{\text{output}} v_{in}$$

$Q_{0.1}$
 0.4 -0.2
 0.1 0.3

$$h = \sum w_i x_i = 0.1 \times 0.4 - 0.2 + 0.1 \times 0.1 + 0.3 = -0.02$$

$$a = f(h) = \text{Sigmoid}(-0.02) = 0.495$$

$$g = f(w \cdot a) = \text{Sigmoid}(0.1 \times 0.495) = 0.512$$

$$d'(w \cdot a) = f(w \cdot a)(1 - f(w \cdot a))$$

$$\text{error term } \delta^o = (y - g) d'(w \cdot a)$$

$$= (1 - 0.512) \times 0.512$$

$$\times (1 - 0.512) = 0.0122$$

To calculate error term for hidden unit we scale error term from output unit by w .

$$\text{for hidden error term, } \delta_j^h = \sum_k w_{jk} \delta^o h'(z_j)$$

$$\delta^h = w \delta^o f'(h) = 0.1 \times 0.122 \times 0.495 \times (1 - 0.495) = 0.003$$

Now since we have the errors we need to calculate gradient steps.

$$\Delta w = \eta \cdot g^0 \cdot h'(d) = 0.0302$$

For input to hidden units:

$$\Delta w_i = \eta \cdot g^0 \cdot h' \cdot x_i$$

$$= (0.5 \times 0.003 \times 0.1)$$

$$0.5 \times 0.003 \times 0.3$$

$$= (0.00015 \times 0.00045)$$

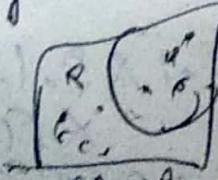
One advantage is that max derivative of sigmoid is 0.00025, so errors can be reduced by atleast 75%. & errors are solved down to 93.75%.

If we have a lot of layers, using sigmoid will reduce weight steps to tiny values near the input. This is the vanishing gradient problem.

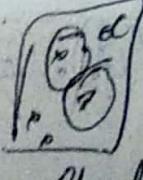
~~Training optimization~~
Errors - underfitting - trying resolution too simple & won't work. Overfitting - performs well on train set but not on test set
underfitting - bias, overfitting - variance
early stopping
overfitting Classifier is too specific.



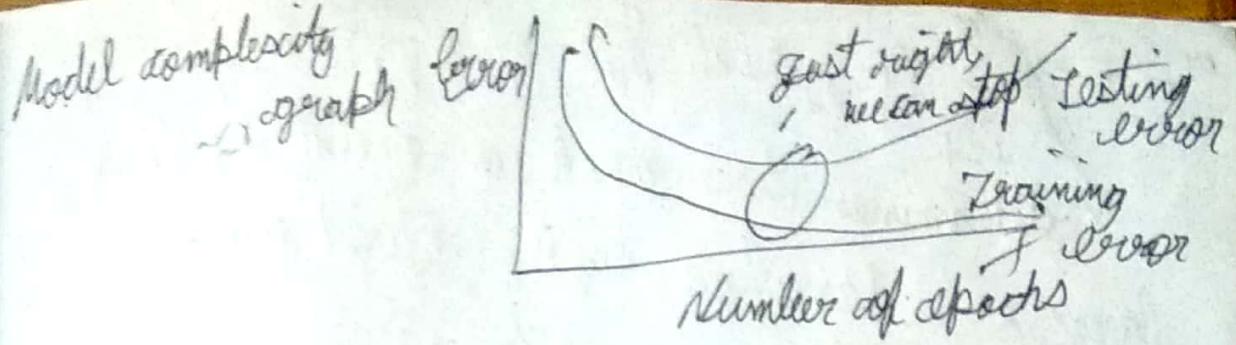
Epoch 1



Epoch 20



Epoch 100



$$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

which equation gives smaller error?

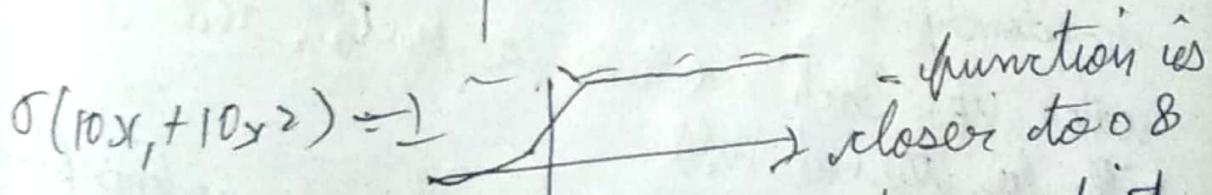
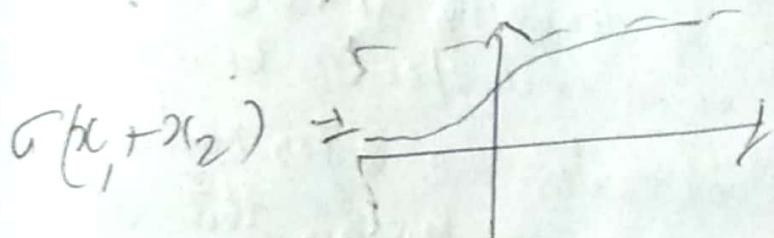
$$x_1 + x_2$$

$$10x_1 + 10x_2$$

Answer is $10x_1 + 10x_2$ as $x_1 + x_2 = \sigma(1+1) = 0.8$
 $\sigma(-1-1) = 0.12$

$10x_1 + 10x_2 = \sigma(10+10) = 0.99$ q (recall close to 1 great precision)

$(-1-1) = \sigma(-10-10) = \sigma(-20) = 0.000021.$



but is much steeper & hard to do gradient descent. Bad models give good accuracy but overfit - solution to this issue is regularization. large coefficients \rightarrow overfitting. Penalize large weights (w_1, \dots, w_n).

$$\text{do } E = -\frac{1}{m} \sum_{i=1}^m (t - y_i) \ln(1 - y_i) + y_i \ln(y_i) + \lambda (|w_1| + |w_2|)$$

$$E = -\frac{1}{m} \sum_{i=1}^m (1 - g_i) \ln(1 - \hat{g}_i) + g_i \ln(\hat{g}_i) + \lambda (\alpha_1^2 - \alpha_n^2)$$

λ determines how much we want to penalize large λ - penalize high. First equation is L1 reg & alone one is L2.

L1

$$\text{Sparsity} = (1, 0, 0, 1, 0)$$

Good for feature selection

Makes many weights to zero

L2

$$\text{Sparsity} (0, 0, 0, 0, 0)$$

Here it does not aim on making other weights 0 but rather keeping the weights homogeneously small.

Dropout

sometimes 1 part of the network

train

has more weights & ends

doesn't train

up dominating others.

So in each epoch we can

randomly turn off nodes so other nodes can catch up. Each node has a prob. that it will be turned off.

Randomrestart: Start from some random place & do gradient descent. This increases probability we will arrive at a global minimum or a good local minimum.

Vanishing grad: In sigmoid function, gradient vanishes at the ends (becomes 0). For a MLP

$$\frac{\partial E}{\partial w_1^{(1)}} = \frac{\partial F}{\partial g_1} \frac{\partial g_1}{\partial h} \frac{\partial h}{\partial a_1} \frac{\partial a_1}{\partial w_1} \rightarrow \text{do gradients}$$

Become really small. We make very tiny steps & never converge.

In alternative to avoid vanishing gradient is to use a different activation function e.g. Hyperbolic tangent.

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

Relu



$$\text{relu}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

derivative is 1 if non-zero

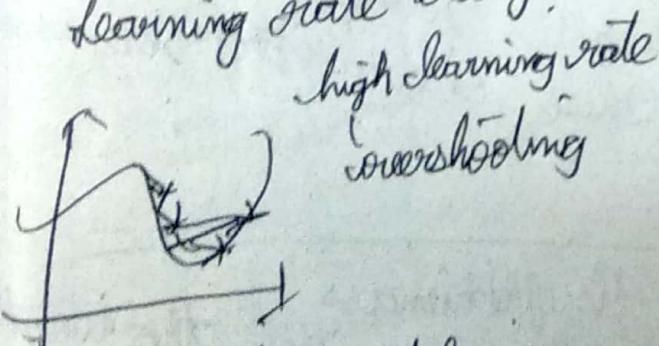
for -

Batch predictions :- In each epoch gradient descent we take entire

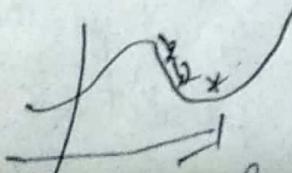
entire data and pass it through the entire network we find the predictions & then backpropagate. If we have huge data it would take lots of memory and huge memory for just one step. There are many such steps.

Stochastic gradient descent - Split data into batches & pass them through the network.

Learning rate decay:



may make model slow



low learning rate

A good model is where learning rate decreases as we approach convergence.

Rule:- 1) If steep : long steps

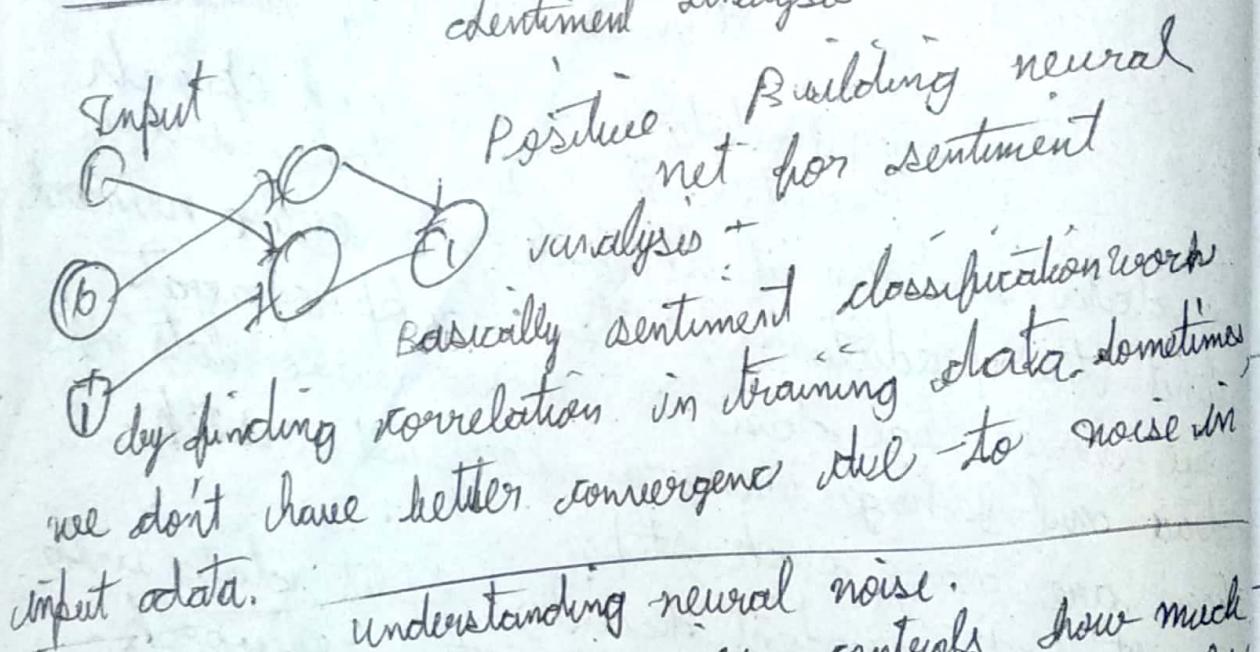
2) If plain : small steps.

Momentum = move with momentum to avoid getting stuck at local minimum. usually momentum can be taken as average of previous few steps.

$$\text{step}(n) = \text{step}(n) + \beta \text{step}(n-1) + \beta^2 \text{step}(n-2)$$

This way steps that happened long time ago may matter less than those that happened recently

Sentiment analysis



Basically sentiment classification work

by finding correlation in training data. sometimes we don't have better convergence due to noise in input data.

Understanding neural noise

Noise vs signal: Each weight controls how much each input contributes to the next layer. If lets say one of the words in our review has value 18 & it is nothing (' ') empty string or it is a period etc. These characters are noise.

Understanding inefficiencies

If a lot of inputs are 0 then the weights sum contributes nothing and computation leads to inefficiency.

Further noise reduction

Stoner noise & semantics

what if we can weigh positive words or negative words most common than other words & we can remove words that do not have much predictive words like "a", "the"

producing

CNN image is viewed as a 28×28 matrix, white pixels are encoded as 2^{255} , black as 0 etc. To convert our 25×25 MIP we have to convert image to vector. This is called flattening. e.g. 25×25 to $1 - 1000$

This is called softmax.
Since image classification is multiclass we can use categorical cross entropy as loss function.

Motel reactivation

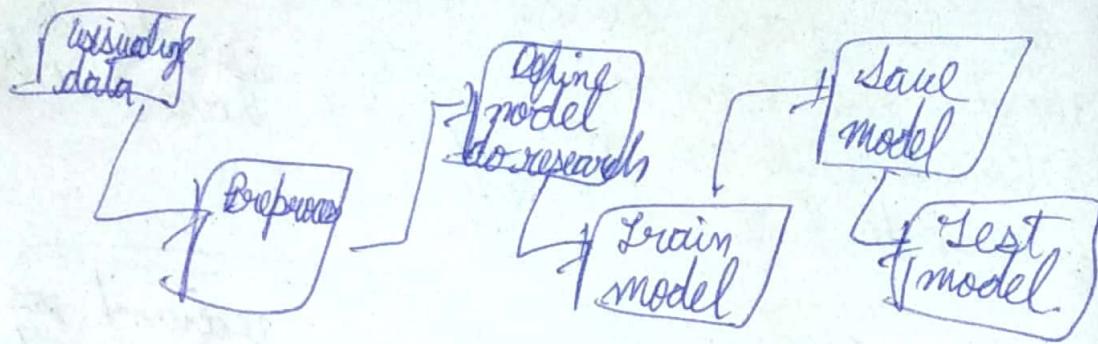
Training set Validation set

Train on training set & calculate train loss & validation loss. Validation set is not used for backpropagation. It is only used to evaluate if model is generalizing well. Since it is not used to update weights it helps us identify if model is overfitting.

Test set cause model selection

Overfitting - we also need a test set because model selection is based on train and validation set. we need a test set so that we can truly see if model is generalised & not biased towards validation set.

Image classification steps



MLP vs CNN Unlike MLP's CNN

know that pixels in close proximity are heavily correlated than those far apart.

Local connectivity

- MLP - lot of parameters, use only fully connected layers. * when we flatten the image to a vector for MLP we lose all 2-D information, knowledge of how pixels are located relative to each other.

CNN - uses sparsely connected layers.

- Accept matrices as input.

Instead of having more hidden nodes & will connect them connected to input layer seen in CNN there are lesser hidden nodes with each node only focusing on a smaller region of image. Thus we have fewer parameters. This will enable each hidden node share a common group of weights.

CNN's preserve spatial information. They process image as a whole or in blocks rather than input one by one as in MLP.

Spatial patterns - shape & color.

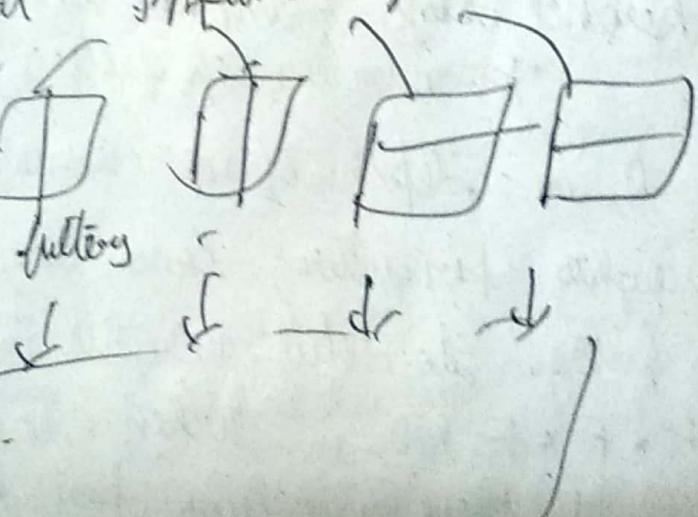
shape - intensity = (black or white).

edges can be detected by abrupt changes in intensity.

High pass filters - High frequency image zone where intensity rapidly changes when moving from one pixel to others. High pass filters enhance such parts of image (emphasizes edges).

edge detection filter $\begin{array}{|c|c|c|} \hline 0 & -1 & 0 \\ \hline -1 & 4 & -1 \\ \hline 0 & -1 & 0 \\ \hline \end{array}$ It is important that values sum to 0. Centre pixel is more important. Negative weights (-) increase the contrast in the image.

convolution layer
filters are initialized at random & so all the patterns detected. we define a loss function like categorical cross entropy for multi class. & update the filters to learn the patterns needed.



activation map
(filtered image)

stride & padding
stride - by how much the filter moves over the image.

Increasing depth

Resizing images is needed as CNN's need same size input. we resize such that spatial dimensions equal a power of 2^2 or divisible by a small power of 2^2 .

Conv layers help increasing depth of the image & max pooling decreases $\times 8$ Y-coordinates. Increasing depth helps extract useful features.

Number of parameters in a conv layer: It depends on filters, out channels kernel size & input shape -

K - number of filters, F - height & width
 D_{in} - depth of previous layer. $F \times F \times D_{in}$
 weights per filter, conv layer is composed of K filters. So total weights in conv layer is $K \times F \times F \times D_{in}$. When taking bias into account, there is 1 bias per filter which means total no. of weights = $K \times F \times F \times D_{in} + K$.

Shape of conv layer: $\frac{W_{in} - F + 2P}{S} + 1$

Image augmentation:

scale invariant - we don't want pred. to depend on size of object.

rotation invariance - not depend on angle of object
translation invariance - not depend on position of object.

Imagenet - 10 million hand labelled images drawn from 1,000 image categories
Alexnet pioneered ReLU & dropout to avoid overfitting.

VGG-16 819 - 16 819 layers - 3×3

conv, 2×2 pooling 83 FC layers.

Alexnet - 11×11 windows.

Resnet - similar to VGG but latest version has 152 layers. More layers - vanishing gradient problem. skip layers were added which enabled gradients to take a shorter route & avoid vanishing gradient problem.

Transfer learning

VGG-net - we can take an existing model & apply it to new dataset. If data set is small & has distinct shapes & features similar to Imagenet we can leverage most layers of VGG & just use final trainable layers.

Opinion to use transfer learning

		Similar	Different
size of data	large	Fine-tun.	
		fine tune + pretrain	start fresh comnet
small	small	rand wts comnet	start fresh comnet

weight initialization

When we use transfer learning we use least pre-trained weights for image classification.

But for non-pretrained we need to figure out way to initialize weights.

All weights' initialization would have ~~no~~ negative effect as it is impossible to take the weights or nodes and cause of error as because the gradients are all the same it is difficult to figure out how much weight should be updated.

If weights are all $1's$ same thing happens. But $1's$ are more harmful as they are a very large value & it is difficult for backprop to update & bring the weights down.

All this can be avoided by proper weight initialization by initializing weights to distinct values. This can be done by picking a random value from uniform distribution.

General rule: But initializing weights to just any value like 0.81 is not efficient. Rather since weights multiply with inputs, they should have some relation with weights i.e. inversely proportional. More ~~is~~ number of units mean that the weights ~~are~~ should be smaller.

A general rule for setting weights in network is to set them to be close to zero without being too small - i.e. weights should be in range $[-g, g]$ where $g = \gamma Nn$ (n - no. of inputs)

Autoencoders

CNN through its previous layers discards spatial information & isolates high level info about the content of the image. This can be viewed as data compression (to feature map).

Now autoencoder has 2 parts: encoder & decoder. Autoencoder reduces the dimensionality of any input. This compressed format can be used for various purposes that are faster than original image. Autoencoders help in denoising image.

Learnable downsampling

Conv layers preserve spatial information. A series of conv & max pooling layers downsample the spatial dimensions. This constitutes the encoder.

This can be reversed by reversing down sampling. This can be used as decoder that gives the original image. We can use unpooling to get original vectors. This can be done by linear interpolation that is nearest neighbours.

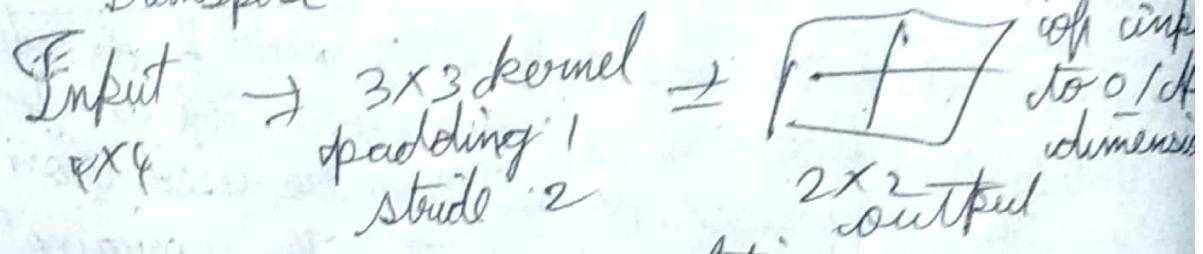
1	2
3	4

1	1	2	2
1	1	2	2
3	3	A	9
3	0	4	6

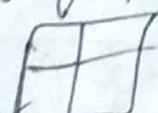
Linear interpolation is a crude idea as we are just copying numbers. A more realistic method have variety and learnable parameters.

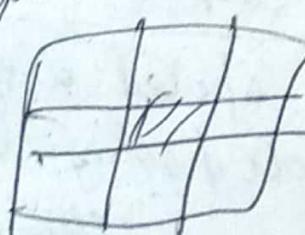
A better way to upscale the image is Trans convolutions. This method does not use pre-defined convolutions but rather has learnable parameter convolutions but rather has learnable parameters. They are also called deconvolution layers. But they do not undo convolutions but rather upsample the input using learnable filter weights.

Transpose convolution



Now for transpose convolution

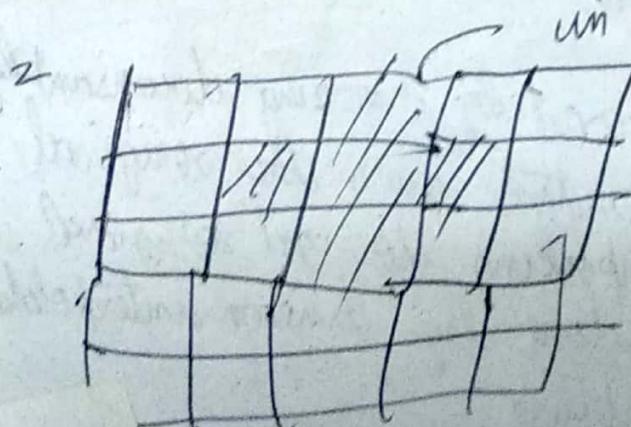
 - we take one pixel from this 3x3 apply a 3x3 filter to it which in turn gives a 3x3 output.



Now for next pixel we take center & since stride is

in case of overlap values would be summed together.

- result is a 5×5



upsampling & denoising

Sometimes blocky checkerboard patterns appear in decoders. They maybe due to overlap. We can apply denoising to remove such patterns. Avg Nearest neighbours could be used to interpolate & avoid checkerboard pattern.

style transfer: A trained CNN has already learnt to represent content in later layers. But what about style? In case of a painting style would be brush strokes, etc.

VG919

First the content image is passed through VG919 & it goes through feedforward till it reaches the deepest convolutional layer. This layer gives the content representation.

When the style image is passed through VG919 different layers extract different styles. The content & style image can be combined to give the output image.

A content loss is defined to compare the content extracted from VG919 with original content image.

content image - C_c .

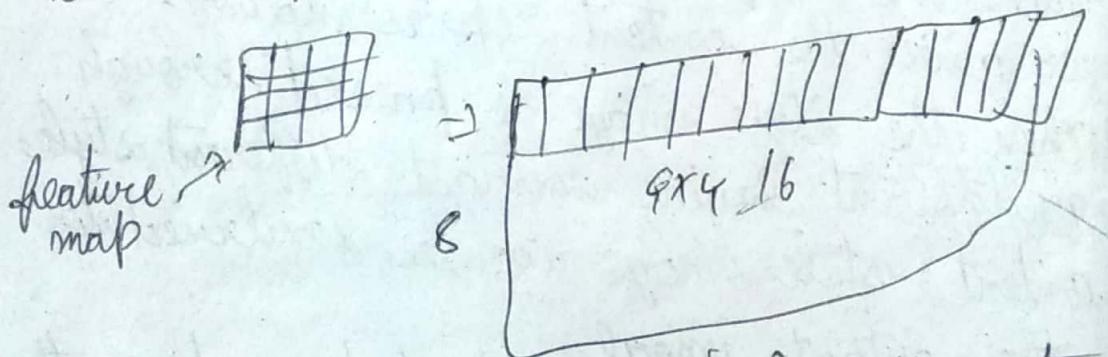
Target image - T_c
content

$$Content = \frac{1}{2} \sum (T_c - C_c)^2$$

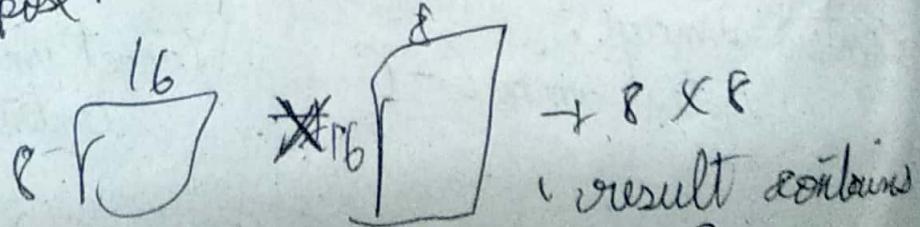
Here we try to optimize but not do training
but changing target image to match source image
(content).

We want to do the same thing for styles.
i.e. compare the target image style representation
style image. This is done by looking at correlations
between features in individual layers. This includes
calculating how similar features in single layer
are. This gives a multi scale style represent.
correlations at each layer are given by gram matrix.

4×4 \rightarrow $4 \times 4 \times 8$ - it has 8 feature map
input convolutional layer that we want to find
correlations between.
First step involves vectorizing the layer.



Next step involves multiplying the outcome by
its transpose.



is localized into which means result won't vary
if images are shuffled in space.

so compute style loss we need to find MSE between the gram matrices from style image S from target image -

$$L_{\text{style}} = \alpha \sum_i w_i \cdot (T_{d,i} - S_{d,i})^2$$

↓ weights. only that changes

Total loss = $L_{\text{content}} + L_{\text{style}}$ - minimize them.

we need to add constant weights $\alpha \beta \rho$ as content and style loss are calculated differently and we need to weigh them differently.

$$\text{Total loss} = \alpha L_{\text{content}} + \beta L_{\text{style}}$$

α more as it is more
style is more important.

$\frac{\alpha}{\beta}$ - lesser α/β ratio more stylized the image.

$$L_{\text{content}} = \frac{1}{2} \sum_i (F_{ij}^d - p_{ij}^d)^2$$

activation of filter at position j in layer d .
derivatives of loss w.r.t activations in layer d

$$\frac{\partial L_{\text{content}}}{\partial F_{ij}^d} = \begin{cases} (F_i^d - p_i^d) & \text{if } F_{ij}^d > 0 \\ 0 & \text{if } F_{ij}^d \leq 0 \end{cases}$$

Then gradient w.r.t. image \vec{x} can be computed using back prop & \vec{x} is changed till it's zero.

same response as a certain layer in the CNN.
the original image ϕ .

style : ϕ^l - original image \rightarrow image
that is generated. A G_{ij}^l -style representations.
layer l. $E_d = \frac{1}{4N^2M^2} \sum_{ij} (G_{ij}^l - A_{ij}^l)^2$
 N, M - height & width of feature
maps.

Total style loss is $L_{\text{style}}(\phi^l, \vec{s}^l) = \sum_d w_d E_d$
 $d=0$, weights of
contribution of each layer to total loss.

$$\frac{\partial E_d}{\partial F_{ij}^l} = \begin{cases} \frac{1}{N^2M^2} (F^l)^T (G_{ij}^l - A_{ij}^l), & \text{if } F_{ij}^l > 0 \\ 0, & F_{ij}^l \leq 0. \end{cases}$$

Skin cancer detection :-

Precision is fraction of relevant instances
among those retrieved & recall is the fraction
of total number of relevant instances that
were retrieved. For example if a image of
12 dogs is were identified as dogs but
ground truth is only 5 dogs, then precision = 5/12,
recall = 5/12.

Sensitivity : Of all sick people how many
did we diagnose as sick -

Specificity : Of all healthy how many
we diagnosed as healthy.

sick Diagnosed
 healthy True Positive
 sick False
 healthy Positive

$TP = \text{correctly diagnosed as sick}$ Sensitivity = $\frac{TP}{TP+FN}$

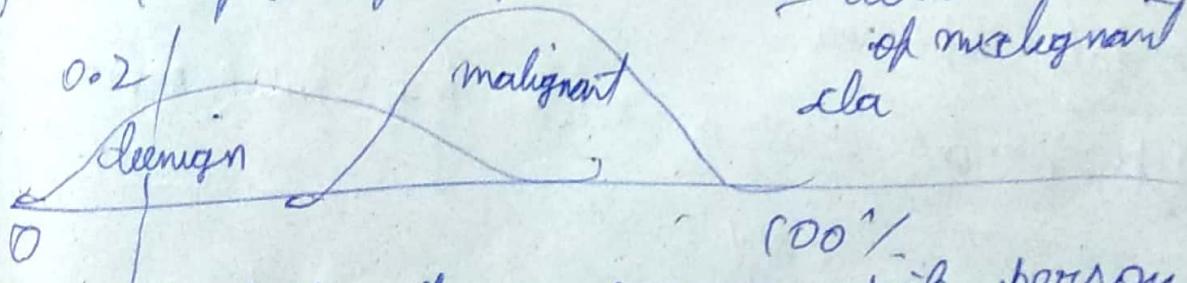
$TN = \text{healthy people diagnosed as healthy}$ Specificity = $\frac{TN}{TN+FP}$

$FP = \text{healthy people diagnosed as sick}$

$FN = \text{Sick people incorrectly classified as healthy}$. Recall = $\frac{TP}{TP+FN}$

Threshold for detecting melanoma $\$0.2, \0.8 or $\$0.5$. Precision = $\frac{TP}{TP+FP}$

Below is histogram of predictions model gives for set of images of lesions.



No threshold of 0.2 ensures sick person gets treatment.

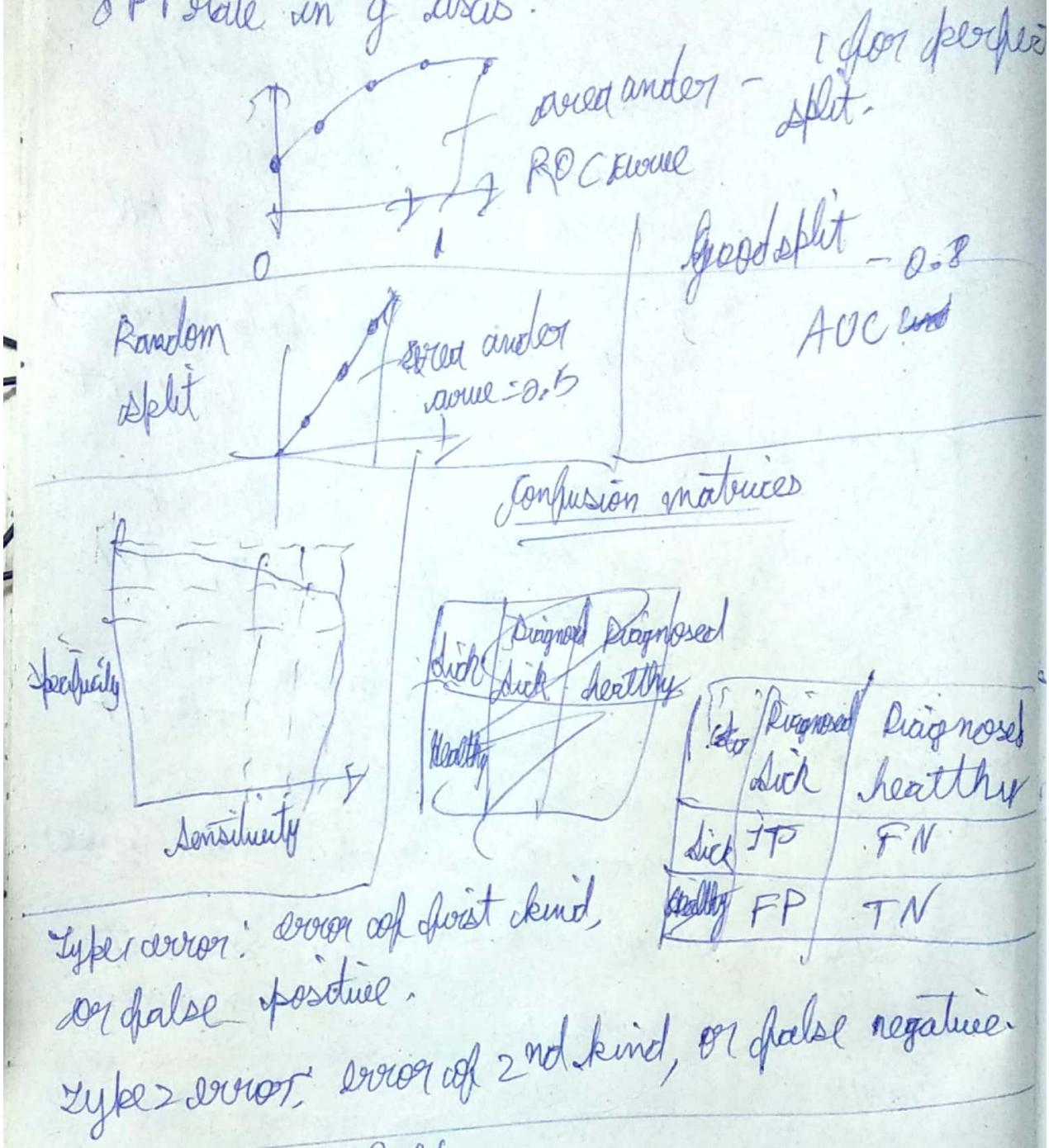
ROC curve True positive rate = $\frac{\text{True P}}{\text{All positives}}$
 Perfect split - 1

Good split - 0.8

Random split - 0.5 or 0.2

$F1_{\text{rate}} = \frac{FP}{\text{All negatives}}$

Plot with points d.e. true positive rate on x axis
& FPR rate on y axis



RNN

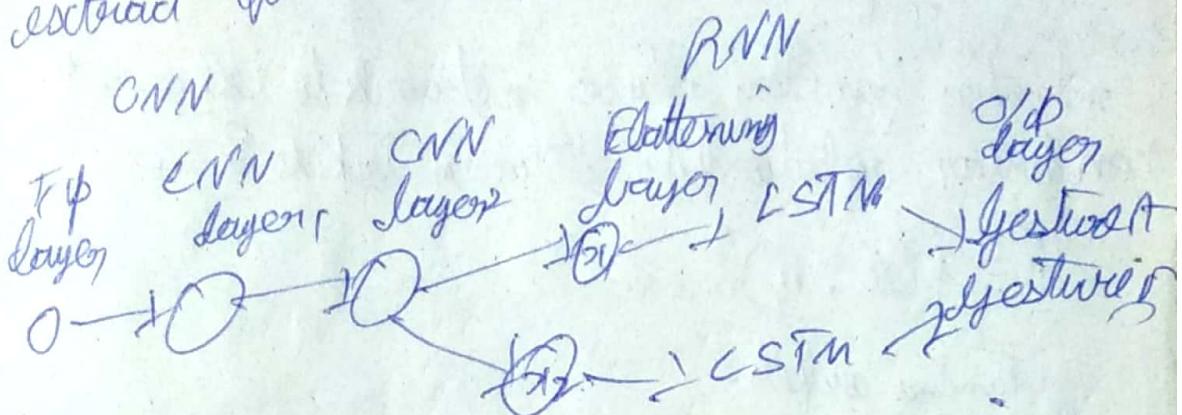
RNN capture temporal dependencies (dependencies over time). RNN - are called recurrent as we perform same task for each element in input sequence.

Early neural networks were limited in that they were not able to model temporal dependencies. So recurrence was added.

FFNN - quick reminders

~~O -> O~~ + ~~O/P~~
 1/O P hidden
~~O -> O~~
 extract features & RNN for memory.

CNN's and RNN's
 can be combined
 CNN's can be used to



Neural networks are non linear function approximators.
 FFNN are static mapping as there is no memory &
 output just depends on inputs & weights, which
 means given same inputs & weights same output

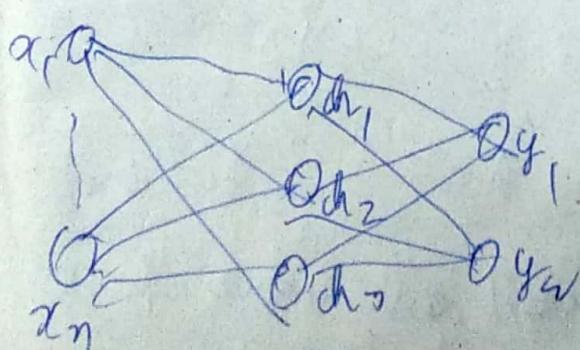
occurs.

feedforward
 backpropagation

This process is
 repeated as many
 times as needed.

Notation w_k - weight matrix k , w_{ij}^k = i^{th} element of
 weight matrix k .

Feedforward



- step 1 - finding α
 - step 2 - finding β
- w_{ii} - connects x_i to h_i ,

$$[h'_1 \ h'_2 \ h'_3] = [x_1 \ \dots \ x_m] \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & \dots & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

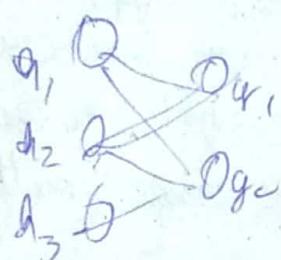
$\bar{h}' = \bar{x} \cdot \bar{w}$, $h' = x_i \cdot w_i$, $x_n \cdot w_n$
 The present h value from exploding we need to
 use an activation function. $h = \phi(h')$.

tan h, f' of f 's rule.

Activation function allows network to represent
 non linear relationships between inputs & outputs.

$$h \equiv \phi(x^T \cdot w)$$

Finding outputs:



$$[o_1, o_2] = [h_1, h_2, h_3] \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

2 columns are same,

20/40 \rightarrow 2 rows as 3 hidden units. $\bar{o} = \bar{h} \bar{w}$

unit $\phi(x^T \cdot w) = a$
 $x^T \cdot w$

$$h_1 = \phi\left(\sum_i x_i w_{i1}\right)$$

$$h_2 = \phi\left(\sum_i x_i w_{i2}\right)$$

$$x^T \cdot w$$

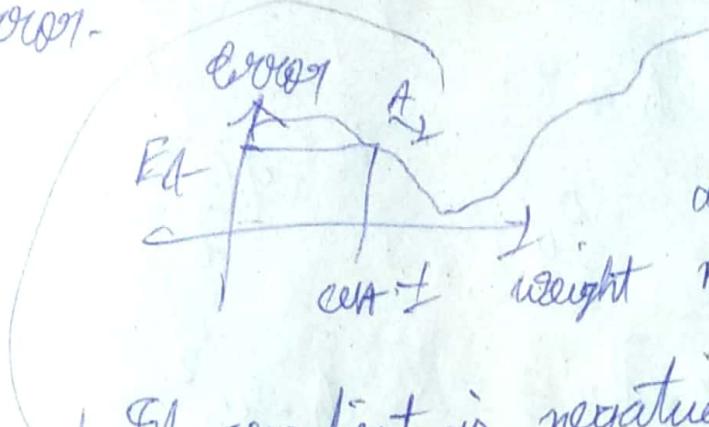
$$h_3 = \phi\left(\sum_i x_i w_{i3}\right)$$

softmax is used when we want all values of softmax to
 lie between 0 & 1 & their sum to be 1.

Backpropagation - SGD with chain rule

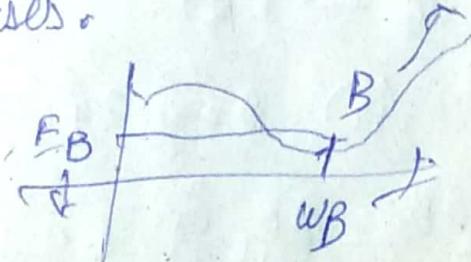
rule: find set of weights minimizing network error.

error-



We need to increase value of w_A to decrease error.

↓ If gradient is negative we need to change weight in that direction so that it increases.



To decrease error we need to decrease weight w_B . The gradient at B

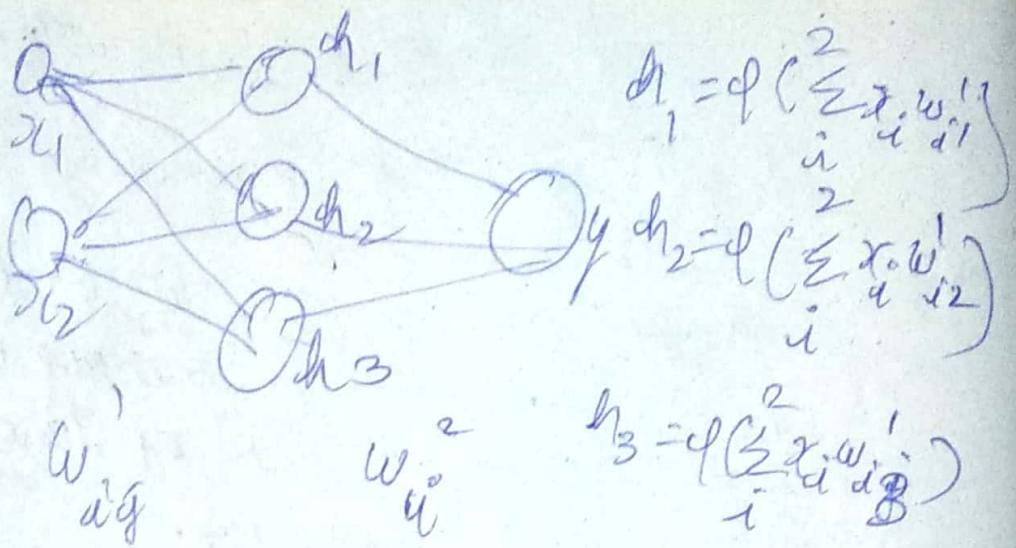
is increasing (positive). If we take a step in negative direction of gradient, then weight would be correctly decreased.

$$w_{\text{new}} = w_{\text{prev}} + \alpha \cdot \frac{\partial E}{\partial w}$$

$$w_{\text{new}} = w_{\text{prev}} + \nabla w (-E)$$

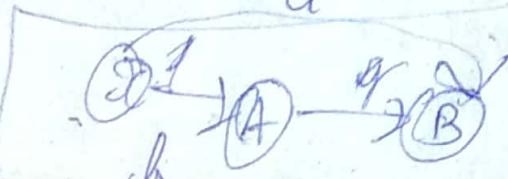
gradient - vector of partial derivatives of error w.r.t each weight.

$$\frac{\partial E}{\partial w_{ij}} = -\frac{\partial E}{\partial w_{ij}}$$



$$h = [h_1, h_2, h_3] \quad y = \sum_i h_i \cdot w_i^2$$

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



$$A = f(x), \quad B = g \circ f(x)$$

composing functions is

same as multiplying
partial derivatives,

output to hidden

$$\frac{\partial B}{\partial x} = \frac{\partial B}{\partial A} \frac{\partial A}{\partial x}$$

$$\frac{\partial y}{\partial w_i^2}$$

$$y = \sum_i h_i \cdot w_i^2$$

$$\frac{\partial y}{\partial w_i^2} = \frac{\partial}{\partial w_i^2} \left(\sum_i h_i w_i^2 \right) = h_i$$

$$\frac{\partial y}{\partial w_{ij}^1}$$

$$\frac{\partial y}{\partial w_{ij}^1} = \frac{\partial y}{\partial h_i} \frac{\partial h_i}{\partial w_{ij}^1}$$

$$\frac{\partial h_i}{\partial w_{ij}^1} = \frac{\partial \phi}{\partial (\sum_i x_i w_{ij}^1)} = \frac{\partial \phi}{\partial \sum_i x_i w_{ij}^1}$$

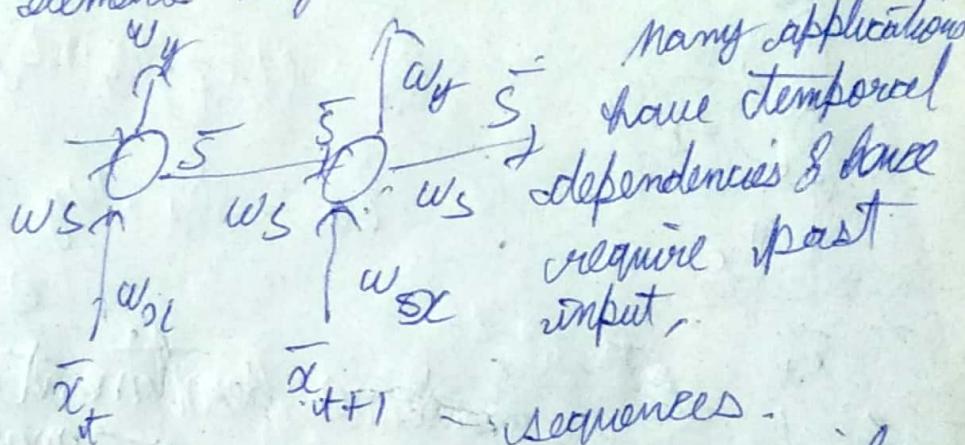
$$\frac{\partial \sum_i x_i w_{ij}^1}{\partial w_{ij}^1}$$

$$\frac{\partial \Psi}{\partial w_{ij}} = \omega_{ij}^2 \varphi_j^2 \cdot x_i$$

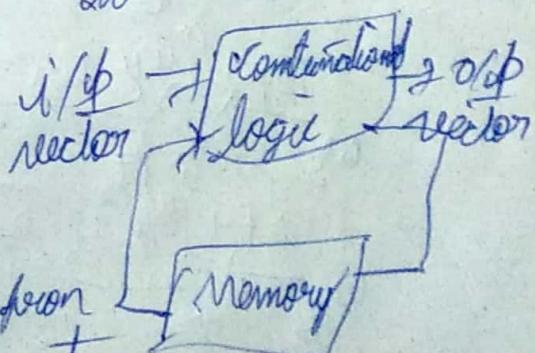
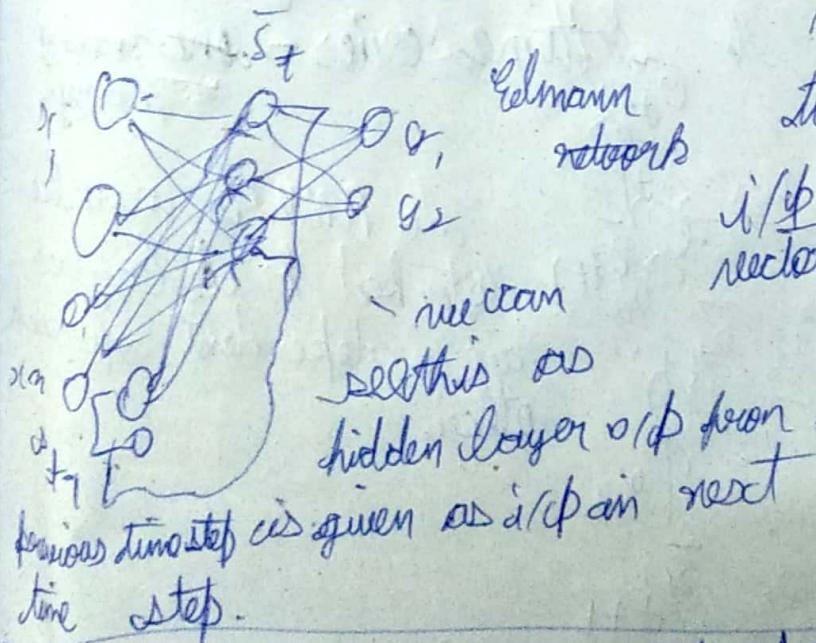
$\frac{\partial \text{loss}}{\partial w_{ij}}$ mini batch gradient descent - updating weights
every N steps. $\delta = \frac{1}{N} \sum_{ijk}^N \epsilon_{ijk}$

RNN - Recurrent neural networks.
They are called RNN as we perform same task
for all input elements. They include memory elements
called states.

many - many
many - one
one - many

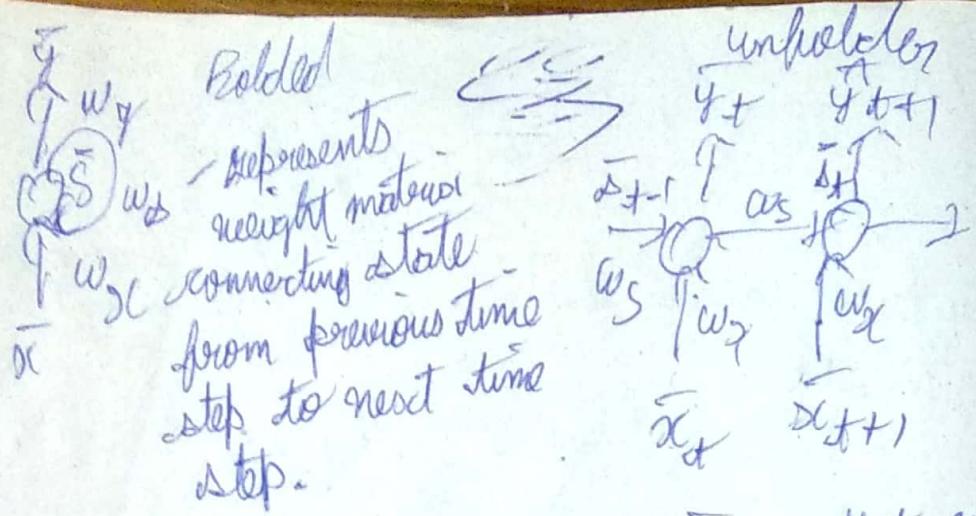


This is similar
to state machine

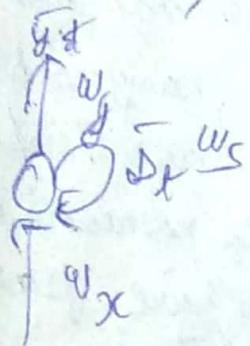


In FFNN inputs are considered to be independent.

But in ANN $\bar{q}_t = f(\bar{x}_t, \bar{w}_{t-1})$



\bar{x}_t - input vector at time t , \bar{y}_t - output vector at time t , \bar{d}_t - hidden state vector at time t

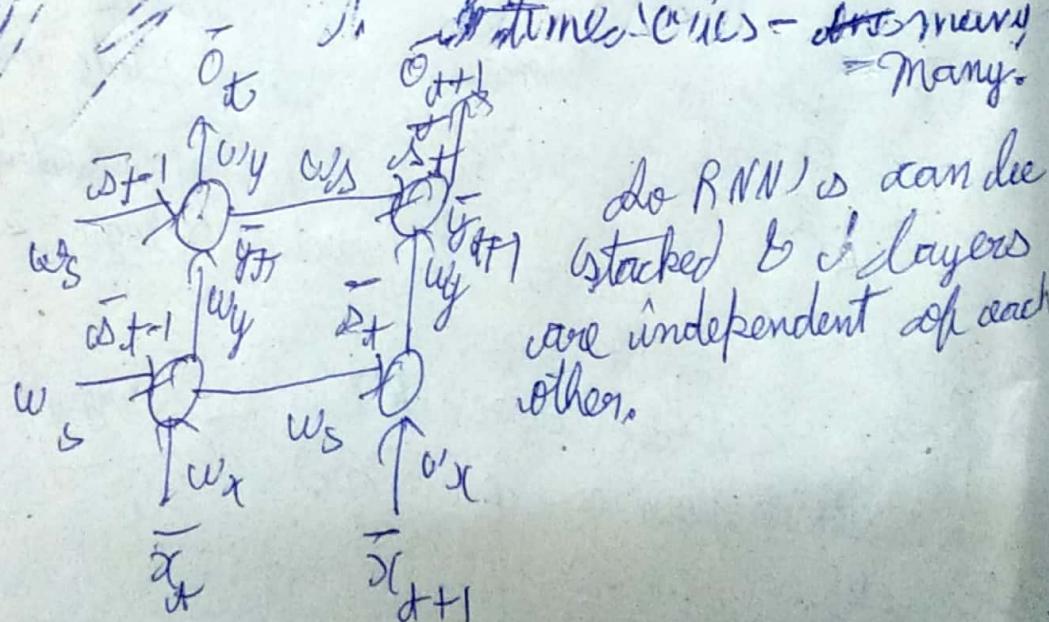


$$d_t = \sigma(\bar{x}_t \cdot w_x + \bar{d}_{t-1} \cdot w_d)$$

$$y_t = \bar{d}_t \cdot w_y \text{ or } \sigma(\bar{d}_t \cdot w_y)$$

diagram of an RNN cell showing inputs x_t , hidden state d_t , and output y_t . It includes weight matrices w_x , w_d , and w_y , and bias terms b_d and b_y .

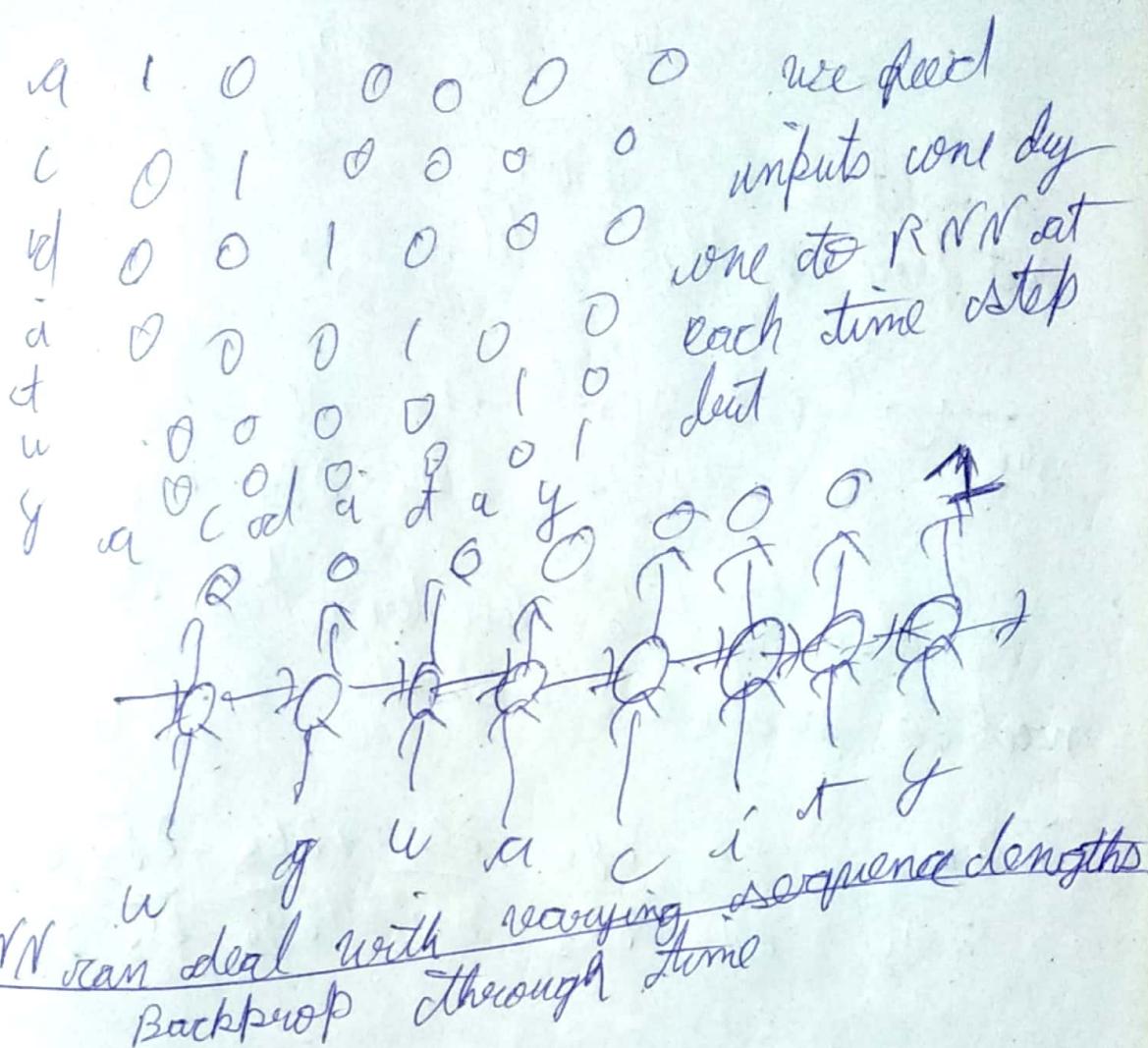
Many - GRU
 Many - LSTM



~~RNN~~ FFNN $\bar{y}_t = F(\bar{x}_t, w)$

RNN, old at direct $\bar{y}_t = F(\bar{x}_t, \bar{d}_{t-1}, \bar{d}_{t-2})$

RNN-example - sequence detection - word detector
Here we want to detect word adjacency.



Timestep 3

$$E_3 = (d_3 - g_3)^2$$

d - desired output w_y w_s $\frac{\partial E_3}{\partial w_y}$
 g - calculated output s_3 $\frac{\partial E_3}{\partial w_s}$ s_3
 w_{yx} $\frac{\partial E_3}{\partial w_{yx}}$

We also need $\frac{\partial E_3}{\partial s_3}$

To consider previous timestep - δs_2 , δt_3

$$E_1 = (d_1 - g_1)^2 \quad \text{Unfolding and BPTT} \quad E_3 = (d_3 - g_3)^2$$

$w_y - \frac{\partial E_2 - g_2}{\partial w_y}$ $w_s - \frac{\partial E_2 - g_2}{\partial w_s}$ $w_{yx} - \frac{\partial E_3 - g_3}{\partial w_{yx}}$

s_1 s_2 s_3

3 weight matrices to modify w_x, w_y, w_s

$$\text{At time } t=3, \frac{\partial E_3}{\partial w_y} = \frac{\partial E}{\partial y_3} \cdot \frac{\partial y_3}{\partial w_y}$$

Now let us find updates needed from $(d_3 - y_3)^2$ rest from unfolded model.

$$\begin{array}{c} s_1 \quad s_2 \quad s_3 \\ \downarrow \quad \downarrow \quad \downarrow \\ w_x \rightarrow O \rightarrow O \rightarrow O \rightarrow d_3 = (d_3 - y_3)^2 \\ w_y \quad w_z \quad w_s \end{array} \quad \frac{\partial E_3}{\partial w_s} = \frac{\partial E_3}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial w_s}$$

We also need to consider s_2 & s_1 . So we need to accumulate gradients.

$$\begin{aligned} \frac{\partial E_3}{\partial w_s} &= \frac{\partial E}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial w_s} + \frac{\partial E_3}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial w_s} \cdot \frac{\partial s_3}{\partial s_2} \cdot \frac{\partial s_2}{\partial y_3} \\ &+ \frac{\partial E_3}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial s_2} \cdot \frac{\partial s_2}{\partial s_1} \cdot \frac{\partial s_1}{\partial w_s} \end{aligned}$$

Accumulative grad was calculated considering all time steps (previous)

$$\text{A more generalized form is } \frac{\partial E_N}{\partial w_s} = \sum_{i=1}^N \frac{\partial E_i}{\partial y_N} \cdot \frac{\partial y_N}{\partial s_i} \cdot \frac{\partial s_i}{\partial w_s}$$

Now we need to adjust w_{s1}

$$E_3 = (d_3 - y_3)^2 \text{ at time } t=3$$

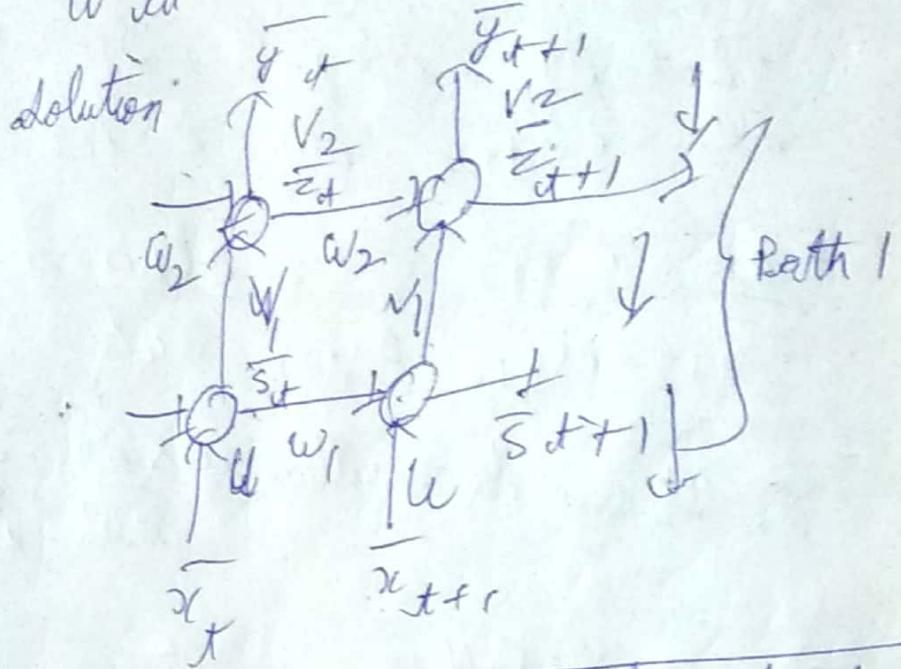
$$\begin{array}{c} s_1 \quad s_2 \quad s_3 \\ \downarrow \quad \downarrow \quad \downarrow \\ w_x \rightarrow O \rightarrow O \rightarrow O \rightarrow d_3 = (d_3 - y_3)^2 \\ x_2 \quad x_3 \end{array} \quad \frac{\partial E_3}{\partial w_x} = \frac{\partial E_3}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial w_x} + \frac{\partial E_3}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial w_x} \cdot \frac{\partial s_3}{\partial s_2} \cdot \frac{\partial s_2}{\partial y_3} \cdot \frac{\partial y_3}{\partial s_2} \cdot \frac{\partial s_2}{\partial w_x}$$

$$\frac{\partial E_N}{\partial w_x} = \sum_{i=1}^N \frac{\partial E_N}{\partial \bar{y}_N} \frac{\partial \bar{y}_N}{\partial \bar{s}_i} \frac{\partial \bar{s}_i}{\partial w_x}$$

v_2
↓
 w_i

w_1
↓
 u
↓
 x

what is update rule off matrix
w at time $t+1$ over 2 time steps

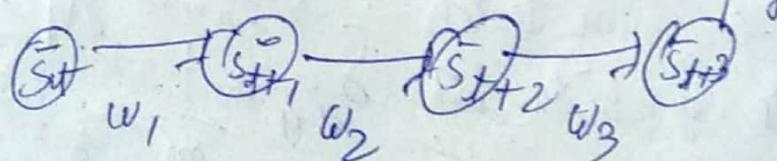


More math: eg. generalization w/ backprop

thru time - generalized form $\frac{\partial E_N}{\partial w_s} = \sum_{i=1}^N \frac{\partial E_N}{\partial \bar{y}_N} \frac{\partial \bar{y}_N}{\partial \bar{s}_i} \frac{\partial \bar{s}_i}{\partial w_s}$

$$\frac{\partial E_N}{\partial w_x} = \sum_{i=1}^N \frac{\partial E_N}{\partial \bar{y}_N} \frac{\partial \bar{y}_N}{\partial \bar{s}_i} \frac{\partial \bar{s}_i}{\partial w_x}$$

To generalize case let us consider



$$\frac{\partial y}{\partial w_3} = \frac{\partial y}{\partial s_{t+3}} \frac{\partial s_{t+3}}{\partial w_3}$$

$$\frac{\partial y}{\partial w_2} = \frac{\partial y}{\partial s_{t+3}} \frac{\partial s_{t+3}}{\partial s_{t+2}} \frac{\partial s_{t+2}}{\partial w_2}$$

$$\frac{\partial s_2}{\partial y_1} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial w_3}$$

$$\frac{\partial y}{\partial w_1} = \frac{\partial y}{\partial s_{t+3}} \frac{\partial s_{t+3}}{\partial s_{t+2}} \frac{\partial s_{t+2}}{\partial s_{t+1}} \frac{\partial s_{t+1}}{\partial w_1}$$

In BFTT we accumulate gradients. 2.

$$\frac{\partial g}{\partial w} = \frac{\partial y}{\partial w} + \frac{\partial y}{\partial w_1} + \frac{\partial y}{\partial w_2} + \frac{\partial y}{\partial w_3} . \text{ weight}$$

matrices are actually identical, $w_1 = w_2 = w_3$

$$\therefore \frac{\partial y}{\partial w} = \frac{\partial y}{\partial \bar{s}_{t+3}} \frac{\partial \bar{s}_{t+3}}{\partial w} + \frac{\partial y}{\partial \bar{s}_{t+2}} \frac{\partial \bar{s}_{t+2}}{\partial w}$$

$$+ \frac{\partial y}{\partial \bar{s}_{t+3}} \frac{\partial \bar{s}_{t+3}}{\partial \bar{s}_{t+2}} \frac{\partial \bar{s}_{t+2}}{\partial \bar{s}_{t+1}} \frac{\partial \bar{s}_{t+1}}{\partial w}$$

$$\therefore \text{It can be written as } \frac{\partial y}{\partial w} = \sum_{i=t+1}^{t+3} \frac{\partial y}{\partial \bar{s}_i} \frac{\partial \bar{s}_i}{\partial \bar{s}_{t+3}}$$

$$\therefore \frac{\partial y}{\partial w} = \sum_{i=t+1}^{t+N} \frac{\partial y}{\partial \bar{s}_{t+N}} \frac{\partial \bar{s}_{t+N}}{\partial \bar{s}_i} \frac{\partial \bar{s}_i}{\partial w}$$

Case example for $i = t+1$, we

$$\text{derive } \frac{\partial y}{\partial w} = \frac{\partial y}{\partial \bar{s}_{t+3}} \frac{\partial \bar{s}_{t+3}}{\partial \bar{s}_{t+1}} \frac{\partial \bar{s}_{t+1}}{\partial w}$$

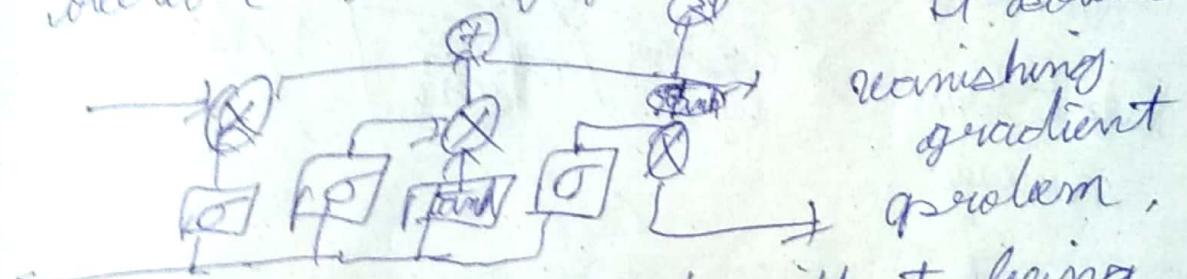
We can use mini batch
we calculate gradients that $\delta_{ij} = \frac{1}{m} \sum_k \delta_{ijk}$
update only after a certain number of steps.

For longer number of timesteps when we
forget update gradients it will lead to vanishing
gradient issue where gradients become very

small. Another issue is exploding gradients where gradients become too large. We can perform gradient clipping where at each timestep we check if gradient is $>$ threshold if yes we normalize the gradient.

CSTM were proposed to

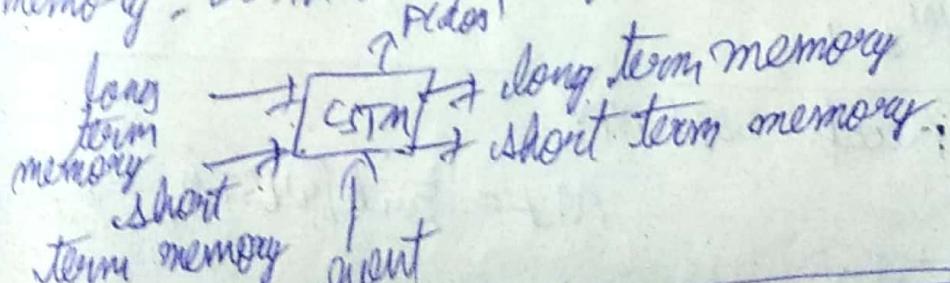
overcome vanishing gradient problem.



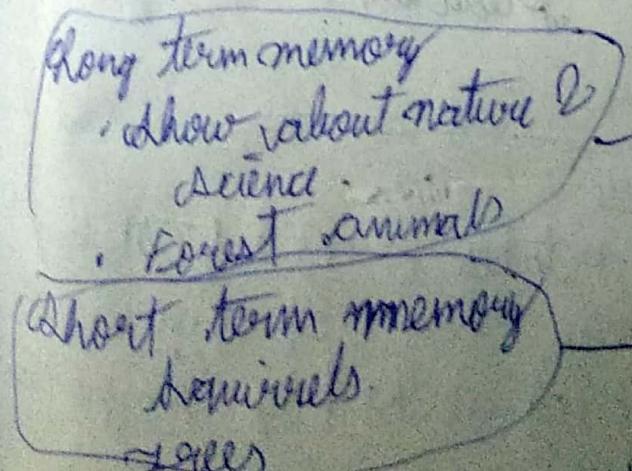
It solves
vanishing
gradient
problem.

Inputs can be stored without being forgotten. The CSTM cell is fully differentiable.

LSTM
ANN suffers from vanishing gradients problem. This is because ANN's only have short term memory. LSTM have long term memory.



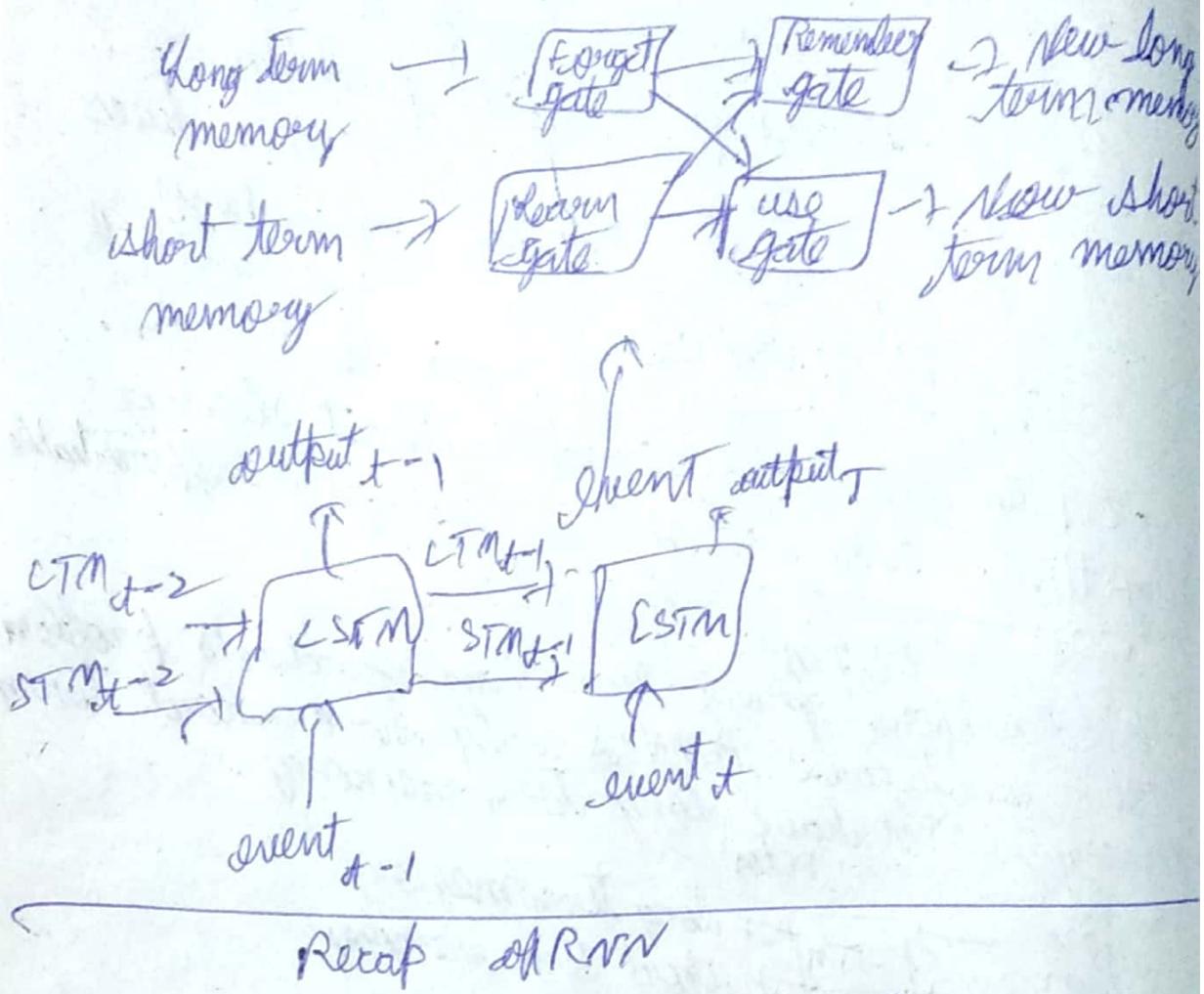
Basics of CSTM



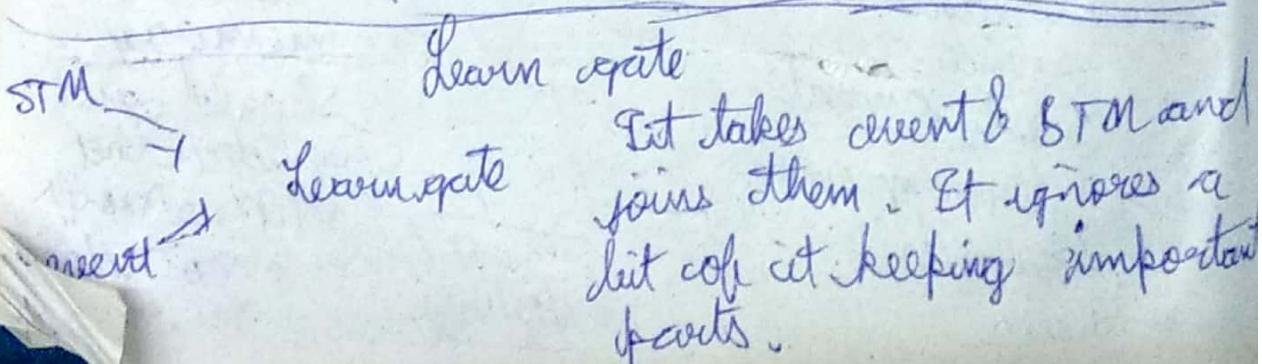
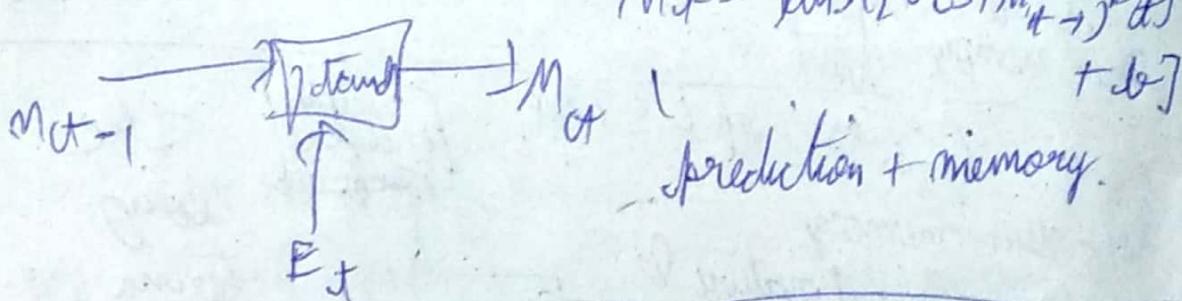
The long term memory should give an idea that O/P is wolf & not dog

Long term memory can be used to update short memory.

LSTM has 4 gates forget gate, learn gate, remember gate & use gate -

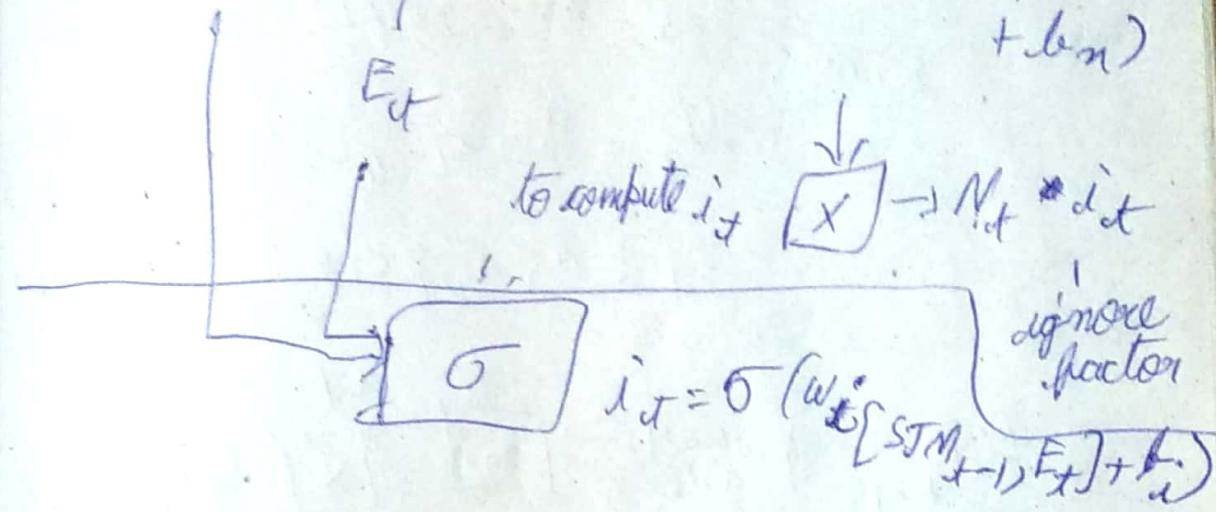


Recap of RNN



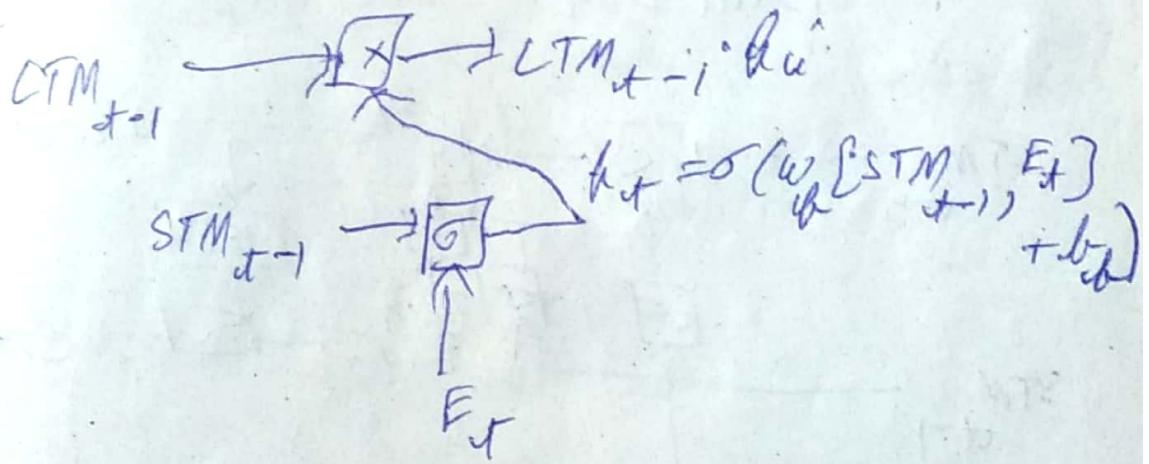
It takes event & STM and joins them. It ignores a bit of it keeping important parts.

$$STM_{t-1} \rightarrow \boxed{f \text{ tanh}} \rightarrow N_t = \tanh(w_n [STM_{t-1}, E_t] + b_n)$$

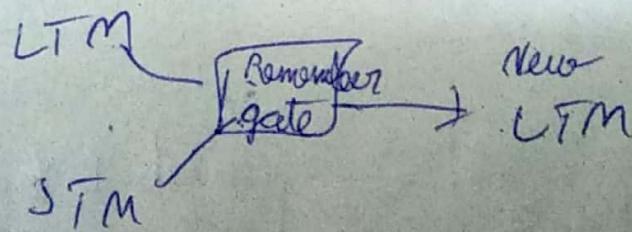


Forget gate

Takes long term memory as input & decides what to keep & forget.

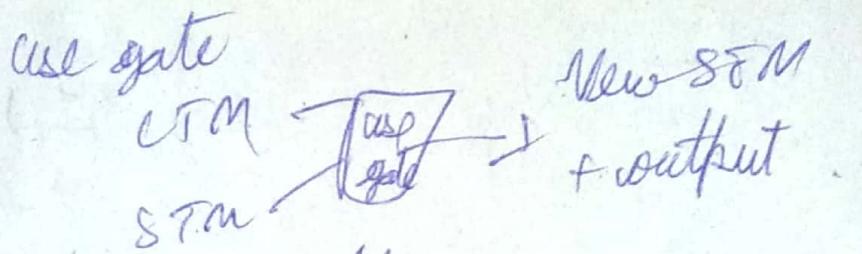


Remember gate combines LTM from forget gate & STM from learn gate & computes new long term memory



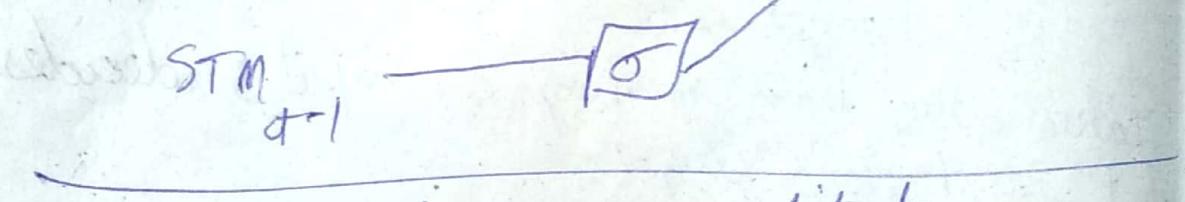
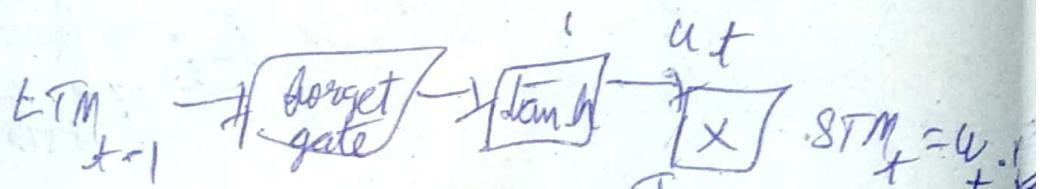
$$LTM_{t-1} \rightarrow \boxed{\text{Forget gate}} \rightarrow f_t \cdot LTM_t = LTM_{t-1} \cdot h_{t-1} + N_t$$

$$STM_{t-1} \rightarrow \boxed{\text{Learn gate}} \rightarrow E_t$$

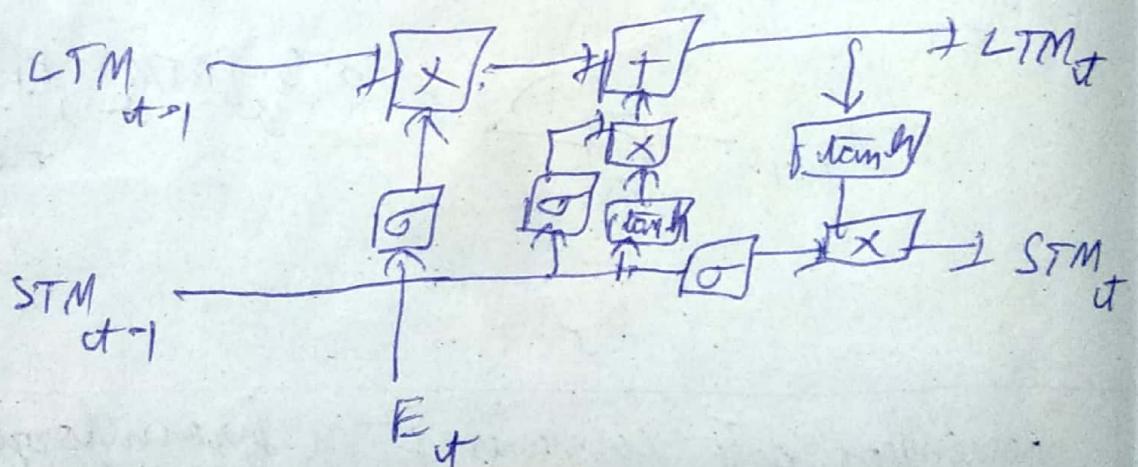


$$h_t = \tanh w_u LTM_t + b_u$$

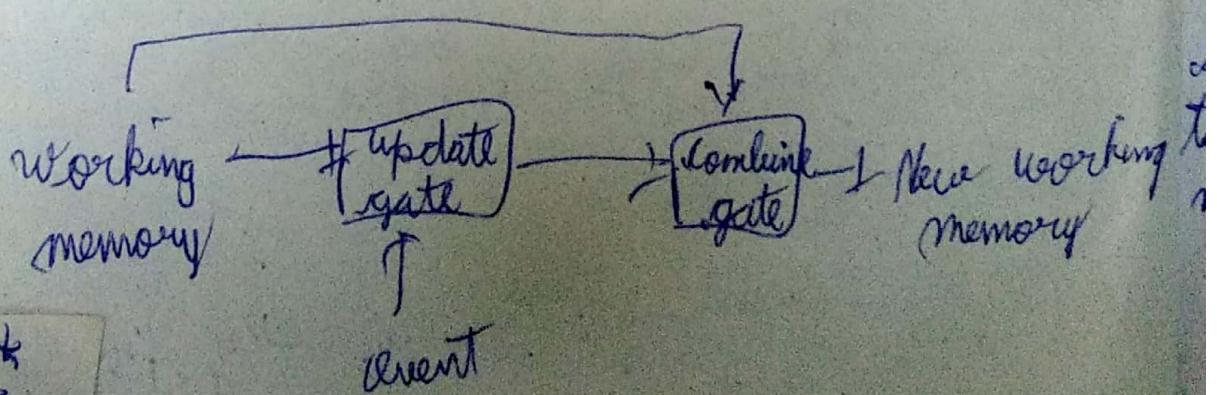
$$v_t = \sigma(w_v \{STM_{t-1}, E_t\} + b_v)$$



Final LSTM architecture

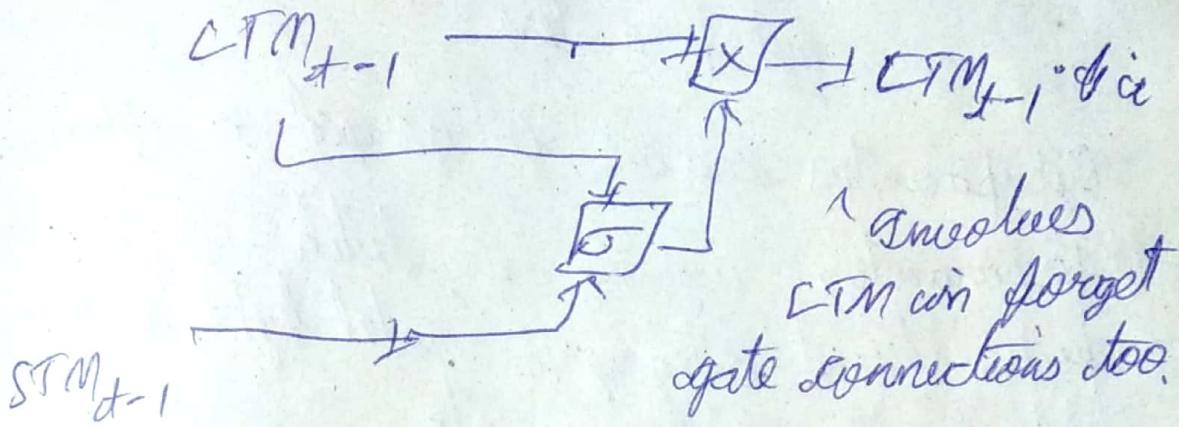


Q/RW



35

Backprop connections



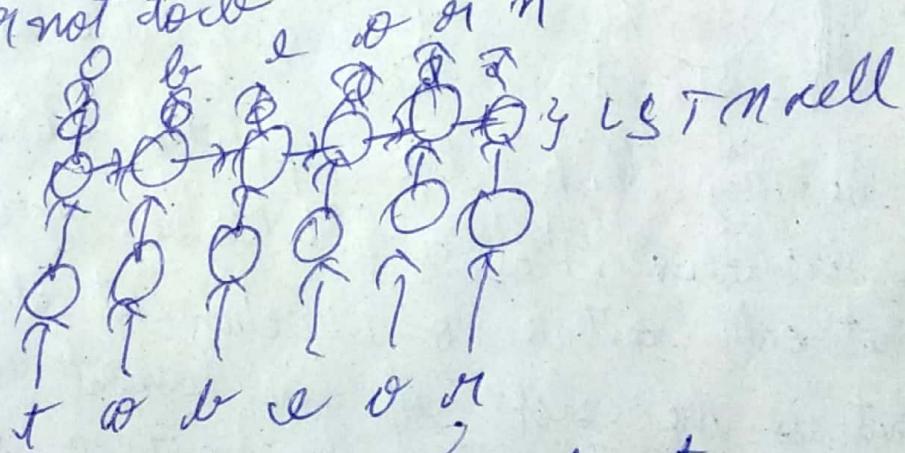
^ involves
LSTM can forget
gate connections too.

~~Character-wise RNN~~

- Predict next character in sequence

J P

To be or not to be



onehot encoded vectors

Sequence batching: In RNN we train on sequence data. Instead of whole sequence if we split & take smaller sequences, matrix operations would be more efficient eg: $[1 \ 2 \ 3 \ 4 \ 5 \ 6]$

We can split it into 2 sequences:

$[1 \ 2 \ 3]$

$[4 \ 5 \ 6]$

den

The batch size = no. of sequences.

Hyperparameters

Hyperparameters can be categorized as optimizers hyperparameters - learning rate, batch size, number of epochs. Second category is model hyperparams. includes number of layers etc.

Learning rate

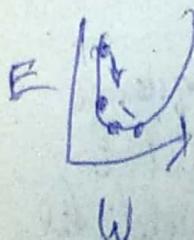
Good starting point - 0.01.



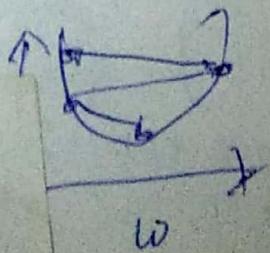
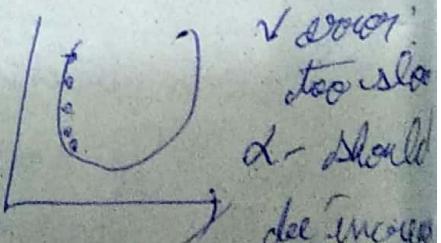
Learning rate is the multiplier used to push weight in right direction.

If learning rate is chosen to

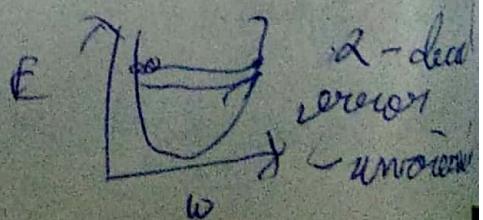
be very high, the updated value will cross more than ideal weight value. Gradient does not only contribute to direction but also to value. That is the slope of line tangent to curve at that point. The higher the point is, the more steep the slope is & larger the value of grad.



Validation on: Reversing learning rate: good



✓ error;
decrease
d - good



Learning rate decay - Reduce error rate during training for convergence.

Minibatch size

It has effect on computing resources, training speed.

online (stochastic)

$\boxed{f(x_1) \quad x_2 \quad x_3}$ - one example
is fed &

with that forward pass, error calc & backprop is done.

Batch - entire training set, minibatch - 1-(online), same size as input (batch). For minibatch often sizes are ; 3, 4, 8 — 256

↓
small minibatch sizes have more resources but cost noises that help not settling on local minima

Word embeddings

converting words to numerical vectors.
Done to reduce dimensionality.

word2vec - learns to map words to embeddings that contain semantic meaning,
eg + word embeddings can learn relationships between present & past tense.

Walking → \vec{x}^0 walked
Swimming → \vec{x}^0 swam

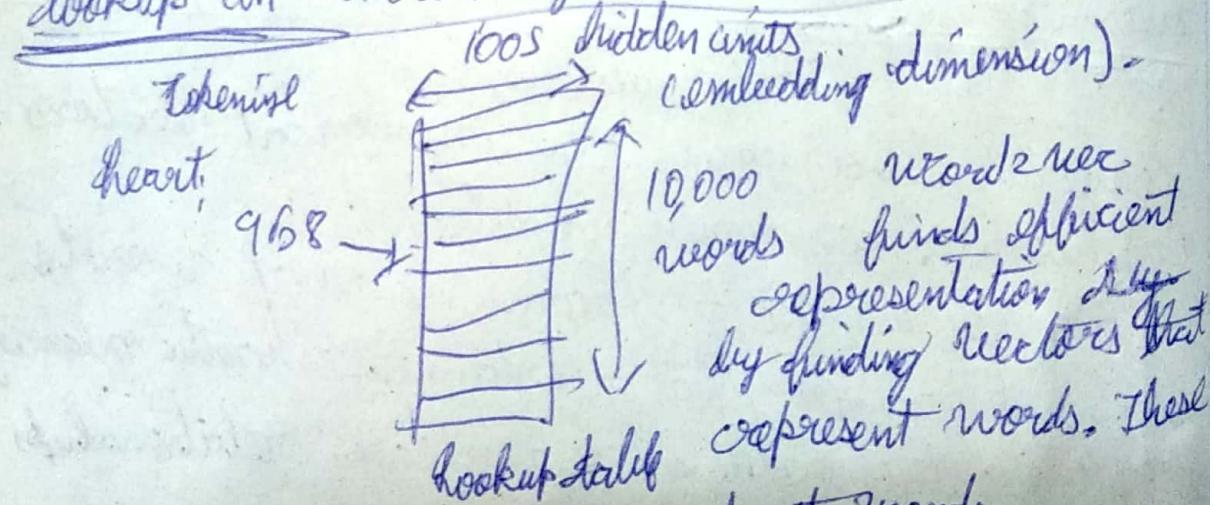
or

golden

embedding weight matrix; Imagine if our words is long (1000's of words), one-hot encoding them will result in a large vector where only 1 value is 1 and others are 0. This is computationally inefficient as dot product operations occur on values that are 0.

Inputs - embedding \rightarrow embedding weight layer \rightarrow embedding weight matrix.

e.g. $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ Values of hidden layer can be directly looked at from weight matrix as each row multiplication of weight matrix with (hot vector gives row where value in one hot is 1). Each word can be encoded as an integer which is an index that can be used to lookup in embedding weight matrix.

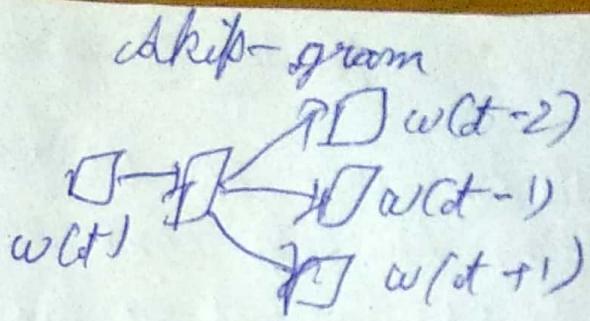
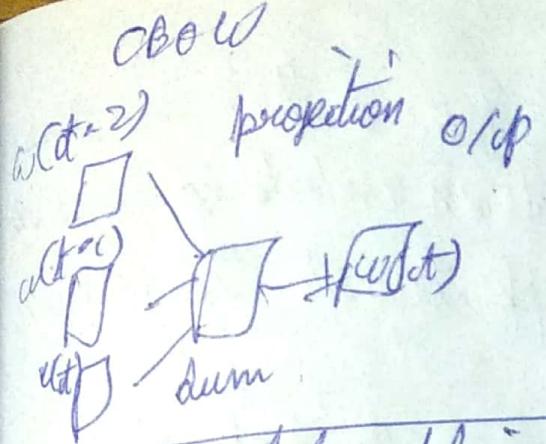


Word2vec is implemented in 2 ways:

- (BOW(cont. bag of words))
- (skip gram)

given context predict missing word.

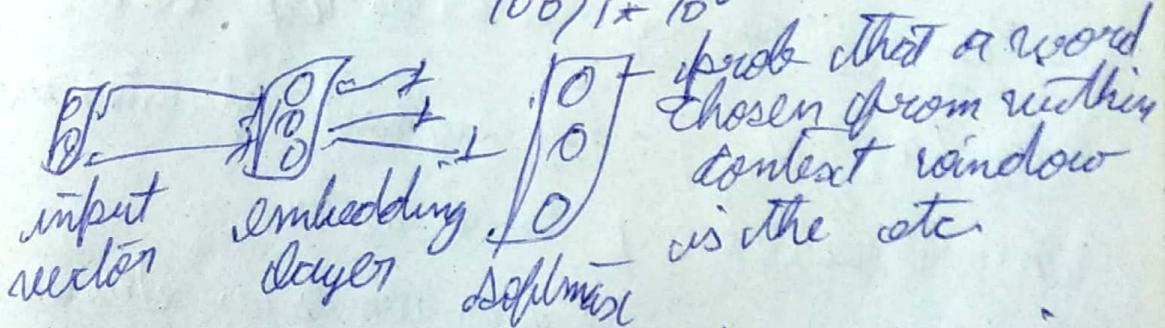
give word of interest, predict context



subsampling: $P(w_i) = 1 - \sqrt{\frac{ct}{d(w_i)}}$

We need to remove words that do not provide much context. We can discard such words & this process is called subsampling. For each word w_i we can discard with prob calculated as above.

For example: At a dict with 1 million words (10^6)
 Learn - appears 700 times
 $t = 1 \times 10^{-4}$. prob. that we will discard word "Learn"
 $1 - \sqrt{\frac{(1 \times 10^{-4})}{700/1 \times 10^6}} = 0.622$
 $- 62\%$



closest words can be found using cosine similarity

$$\text{Cos}(\theta) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

for \vec{a}, \vec{b}

validation words can be encoded as \vec{a} .
 Using embedding table, & then similarity with each word vector \vec{b} is calculated.

on update

Negative Sampling

We have only one or few true outputs in softmax layer. So only some weights in embedding layer can be updated. Instead we can approximate the loss in softmax layer.

To map inputs to embedding weights we can use another embedding layer to map outputs to embedding layer. We need a loss func. that considers only the target & few noisy target words.

$$-\log \sigma(w_o^T v_{w_t}) - \sum_i^N \log \sigma(-w_i^T v_{w_t})$$

w_o^T - embedding vector for o/p target word.

v_{w_t} - embedding vector for i/p word.

$\log \sigma(w_o^T v_{w_t})$ - log sigmoid of inner product of o/p word vector & i/p word vector.

$\sum_i^N w_i \sim p_n(w)$ sum over words w_i drawn from noise distribution $w_i \sim p_n(w)$.

The first term in loss pushes prob that our network will predict correct word to 1. The second term $\log \sigma(-w_i^T v_{w_t})$ will regulate or push prob of noise words to 0.

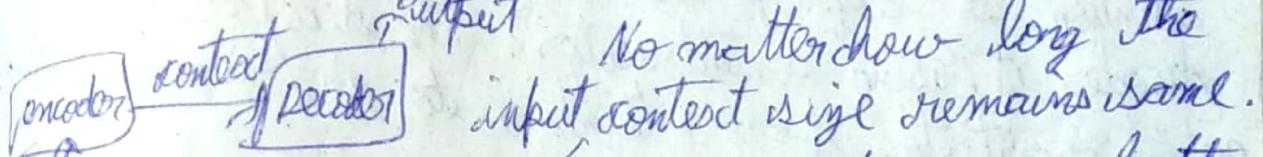
Sentiment analysis with RNN

Attention

Humans don't process an image at once but focus on selective parts, combining the fixations to form an internal representation.

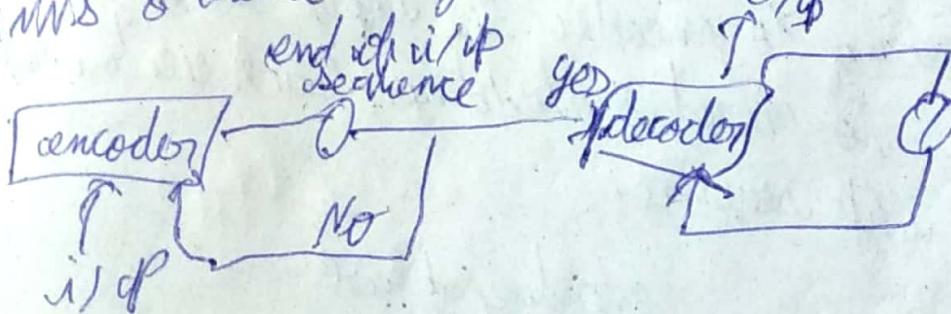
Attention enables sequence - seq models to look at important parts of input & output instead of looking at whole input once & generating o/p

Encoders & decoders

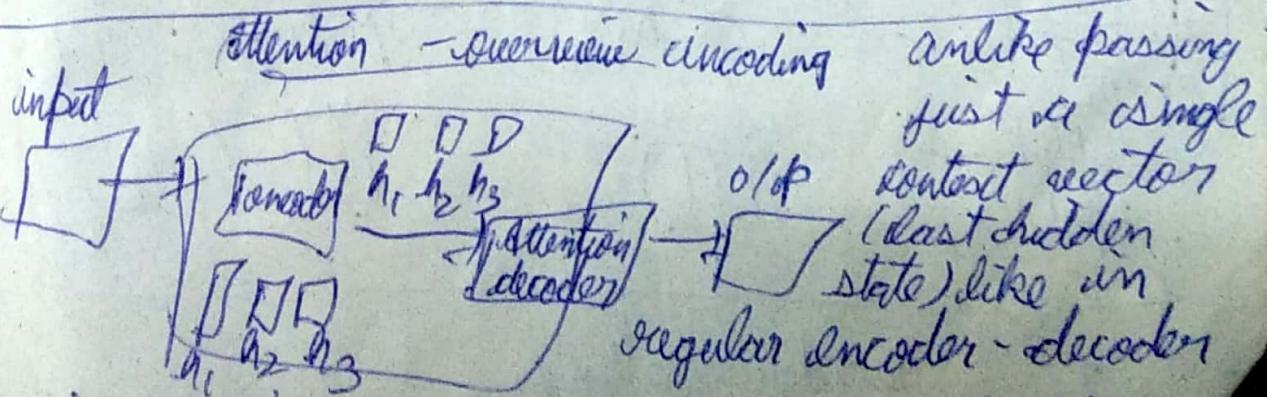


No matter how long the input, context size remains same.

encoder & decoder are both RNNs & hence they have loops:



A major limitation is that the context vector produced by encoder no matter how shorter along the i/p sequence is. Choosing a reasonable size for this vector is a challenge.



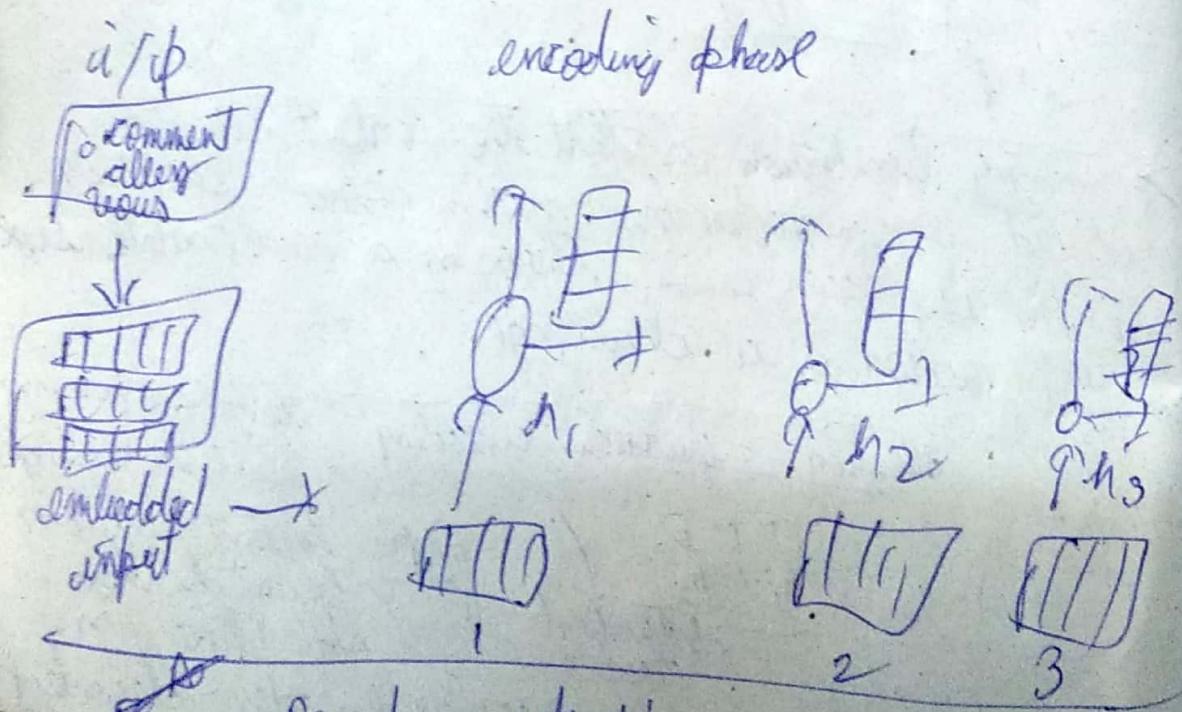
In attention mechanism we pass all hidden

states. This is advantageous as longer sequences have longer context vectors & hence enable capture of information. The main advantage is that attention mechanism enables each hidden state to focus on a word. For example first h_1 represents first word & so on. Although each hidden state represents other info preceding it.

Attention decoder - At every time step it pays attention to appropriate part of i/o sequence & decodes it (using context vector). The decoder knows which part to pay attention to in training process.

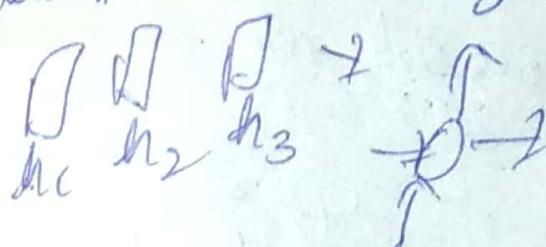
Attention encoder - in depth

Before feeding to encoder the input sequence is translated to word embeddings. Each vector represents words in i/o sequence.



In regard without attention only the

Last context vector is fed to decoder in addition to embedding of `<end>` token.



It uses a scoring function

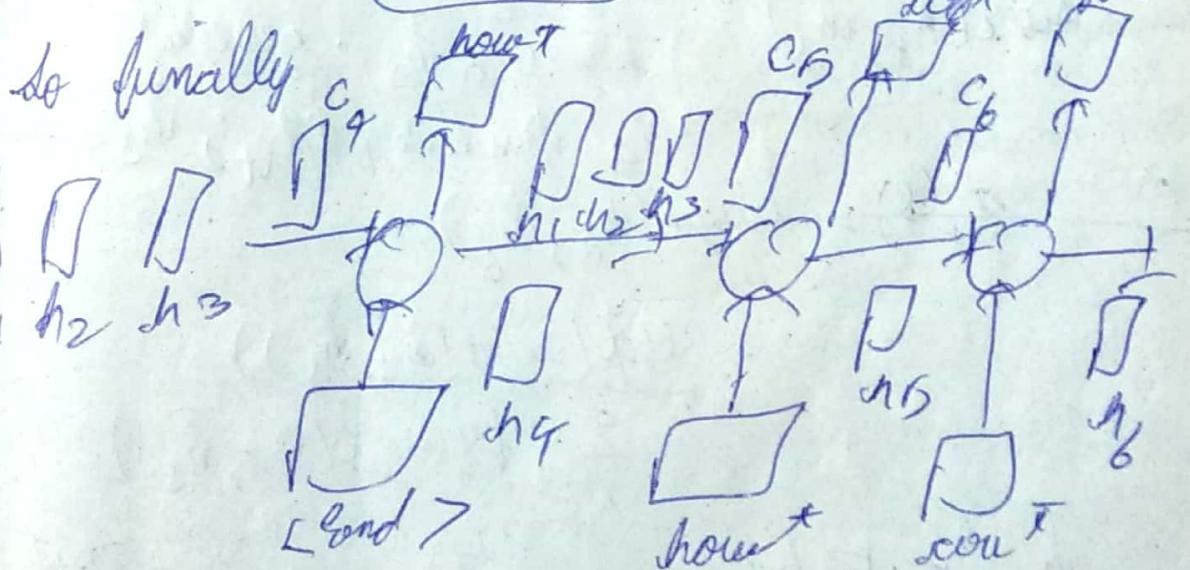
to score the hidden states. The α passes it through softmax

whose output determines

how much these context vectors will be expressed in attention vectors before decoder produces o/p.

Hidden state
* softmax scores

$$[] + [] + [] = [] \text{ context vector for time step } q.$$



2 types of attention - addition & multiplication

Bahdanau

Luong

$$d_{ij} = v_a^T \tanh(w_a s_{i-1} + u_a h_j) \quad h_j - \text{hidden state of encoder at previous step}$$

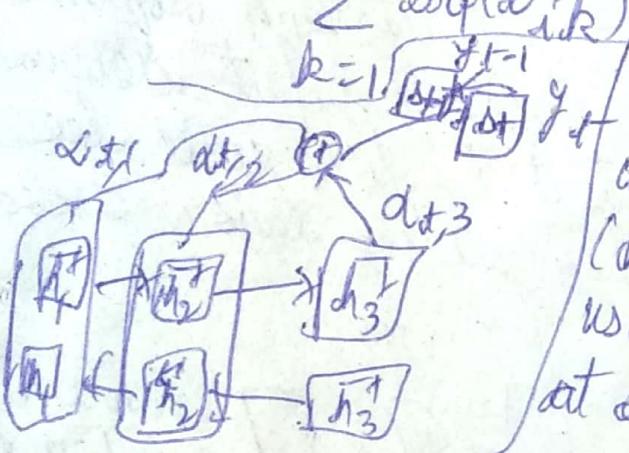
s_{i-1} - hidden state of decoder at previous step

V_a, W_a, U_a - weight matrices learnt during training process.

then they are passed into softmax,

$$d_{ij} = \frac{1}{\sum_{k=1}^{T_2} \exp(d_{ik})}$$

$$c_i = \sum_{j=1}^{T_2} d_{ij} \cdot h_j$$



Multiplicative attention uses only RNN at top (hidden state). This leads us to use stack of RNN at encoder & decoder.
product of hidden states of encoder & decoder

Scaling attention

$$\text{score}(h_t, \bar{h}_s) = \begin{cases} h_t^T \bar{h}_s & \text{dot} \\ h_t^T W_a \bar{h}_s & \text{general} \end{cases}$$

$V_a^T \tanh(W_a [h_t; \bar{h}_s])$ const

$$a_t(s) = \text{align}(h_t, \bar{h}_s)$$

$$= \frac{\exp(\text{score}(h_t, \bar{h}_s))}{\sum_s \exp(\text{score}(h_t, \bar{h}_s))}$$

$$\bar{h}_s = \tanh(W_c [c_t; h_t])$$

- final output

Attention scoring function: It takes hidden state of decoder at decoding step t & hidden state of encoders at each encoding step

$$\text{score}(h_t, \bar{h}_s) = \text{score}([h_t; h_{t-1}; h_{t-2}]) = \boxed{\quad}$$

context vector

It calculates context vector for all encoder hidden states at once but rather for example taking one:

$$\begin{matrix} a \\ h_1 \end{matrix} \cdot \begin{matrix} b \\ h_2 \end{matrix} = \boxed{\quad}$$

the multiplication attention scoring method is

$$\text{score}(h_t, h_i) = h_t^T h_i$$

Now use this method

assume both are ~~same~~ we use ~~same~~ embedding = $\boxed{\quad}$ - Attention

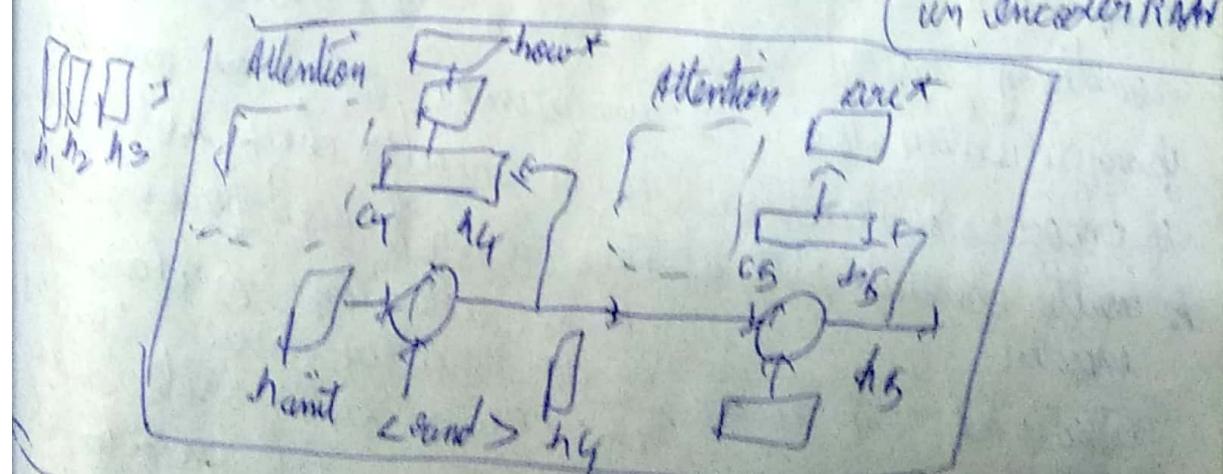
score for each encoder state with n.t current decoder hidden state

But for machine translation encoder & decoder will not have same embedding space. So an alternative method is $\text{score}(h_t, f_j) = h_t^T w_j h_t$

Here a weight matrix is introduced. It is a linear transformation that enables input & output to have diff embeddings.

$$\begin{matrix} \boxed{\quad} \\ \text{1x}m \end{matrix} \times \begin{matrix} \boxed{\quad} \\ m \times n \end{matrix} \times \begin{matrix} \boxed{\quad} \\ m \times t \end{matrix} \xrightarrow{\text{1- no of elements in sequence}} \text{no of hidden units in encoder RNN}$$

Multiplicative attention decoding phase



Attention scoring function - additive (concat)
Here we concat the decoder hidden state & encoder state & proceed

$$\text{score}(\boxed{\quad}, \boxed{\quad}) = \boxed{\quad} \xrightarrow{\text{score}}$$

Formally concat = $\alpha^T \tanh(\alpha_a [h_t; \tilde{h}_S])$

hidden state at time t collection of encoder hidden states.

Attention - computer vision applications - caption generation.

Other attention methods

Attention is all you need - uses a single component called transformer & improve instead of overhead in encoder-decoder architecture. A transformer does not look at the input seq one by one like encoder but rather takes all of them at once & produces output sequence. Transformer does not use RNN.

Transformer - Transformer also

Encoder Decoder uses encoder-decoder architecture but uses feedforward networks instead of RNN & it uses something called self attention. This allows parallelisation too. Transformer contains a stack of encoders & decoders, each encoder has sublayers & multi headed self attention layer & feed forward layer.

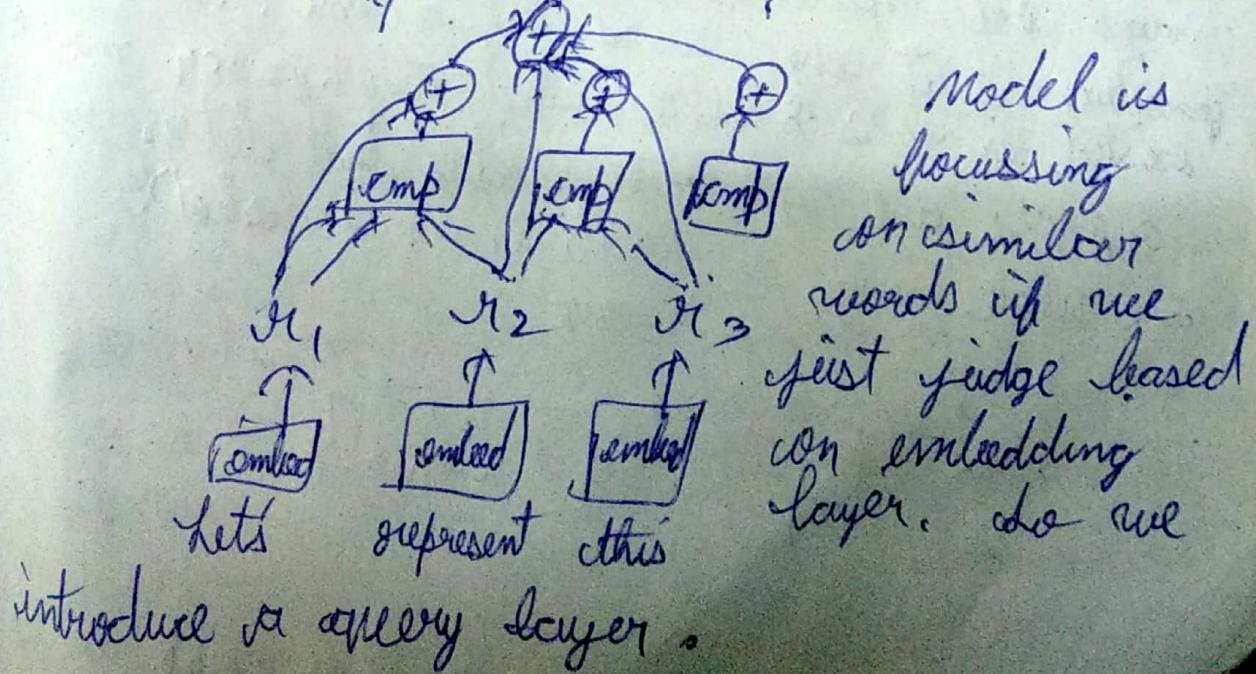
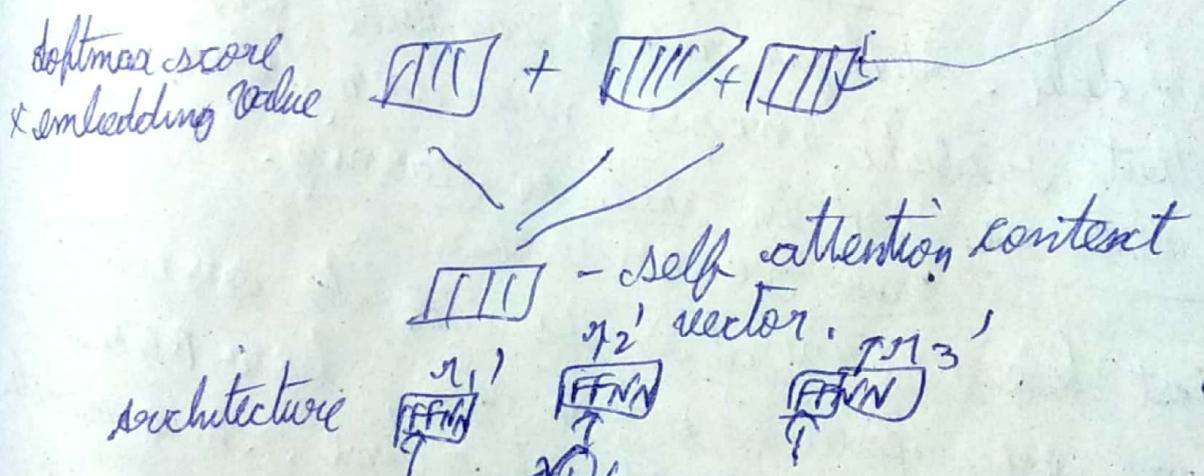
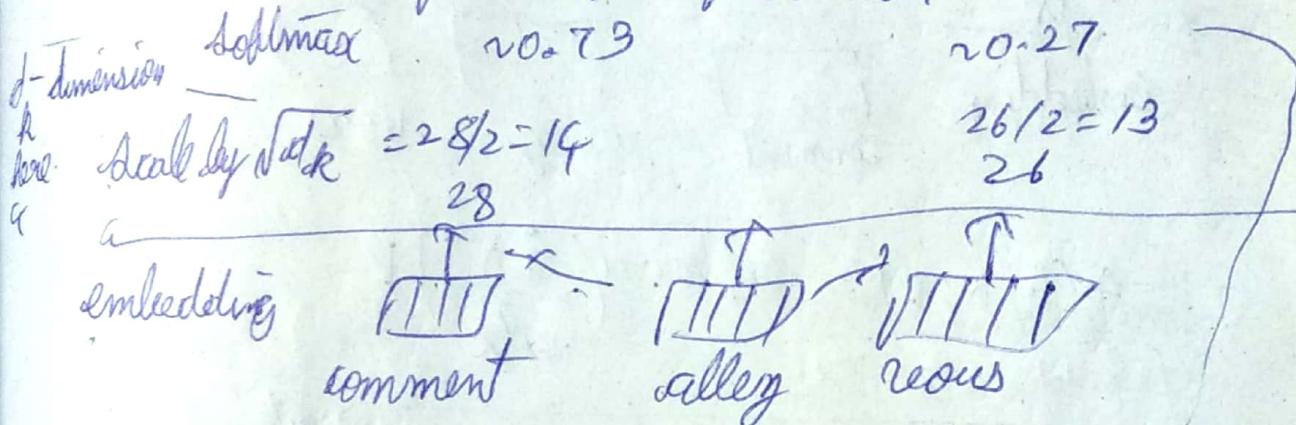
encoder here attention mechanism is on encoder side & the feed forward self attention is on decoder side like previous attention mechanisms. At encoder it pays attention to relevant parts when reading the input. Decoder too has attention components.

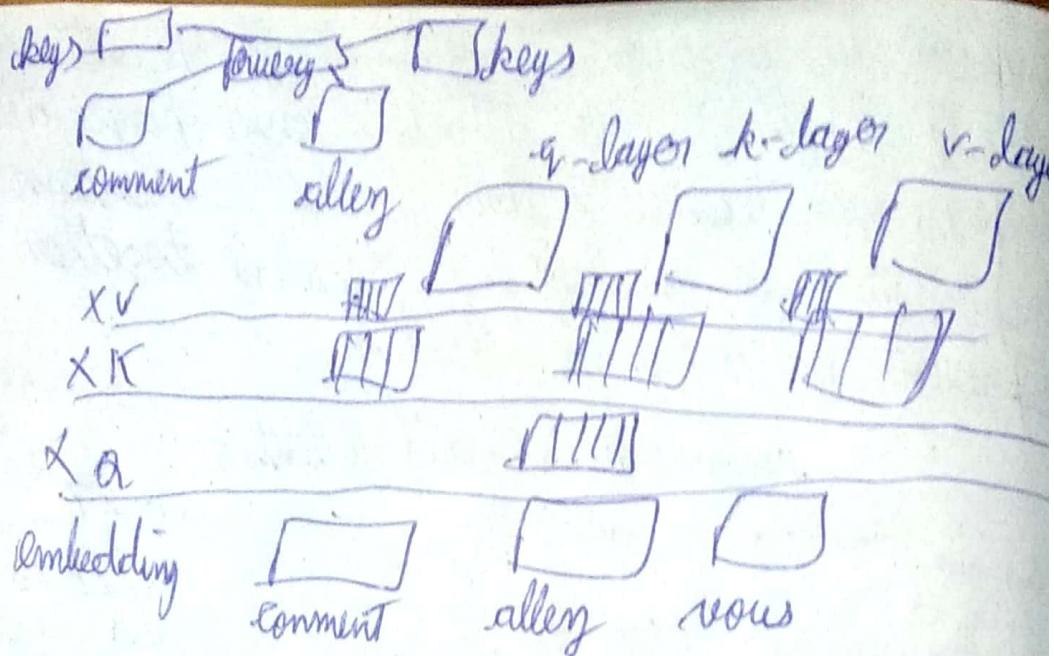
Decoder
feed forward
encoder-decoder attention
self attention

It has 2 attention components.

encoder-decoder attention helps to look at relevant part of inputs & self attention only pays attention to previous decoder outputs. The 3 attention components (in encoder & in decoder) together represent multiplicative attention.

Transformers & self attention



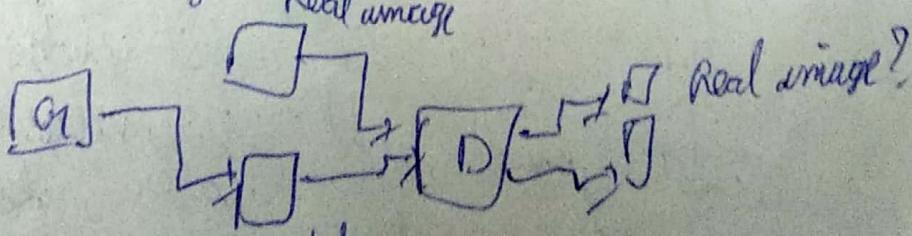


GAN can generate real world data. A GAN model can take input text & output realistic images never seen before. GAN can also be used in imitation learning.

How GANs work

Just like RNNs it is possible to generate images one pixel at a time like how RNNs produce text one word at a time. These models existed in 90's & were called autoregressor models.

GAN's are generative models that let us render the whole image in parallel in one shot.



Generation network takes random noise as i/p.

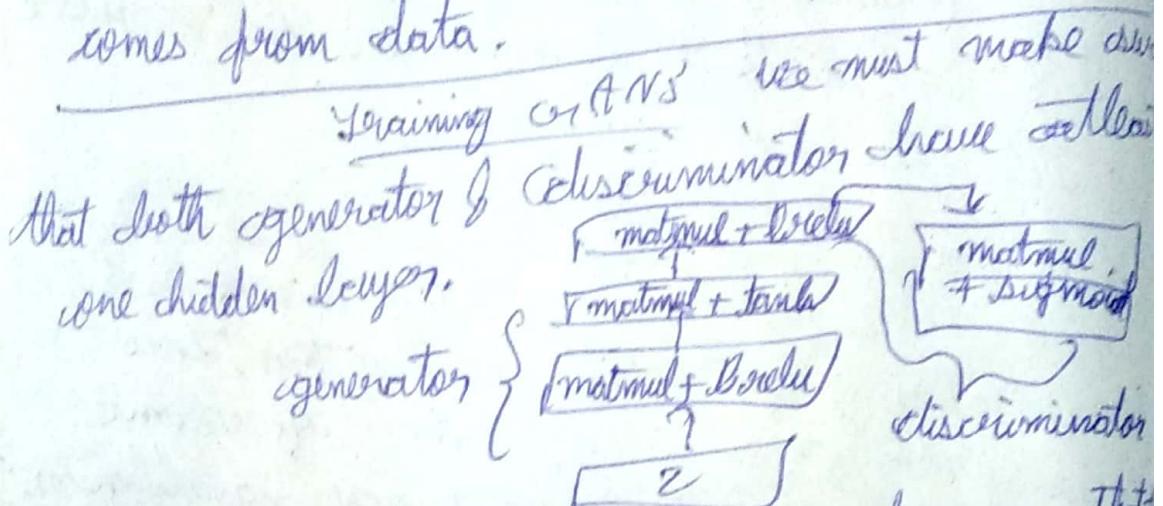
\mathcal{D} runs it through a differentiable function that transforms noise to realistic structure.

Discriminator guides the generator. During training process, discriminator is shown real images & fake images. It outputs prob that the input is a real image. It outputs a prob close to 1 to real images & near 0 to fake images. \mathcal{D} has to output prob of 0.5 everywhere for convergence.

\ Games & Equilibrium : The adversarial in GAN stems from the fact that generators & discriminators compete with each other. This comes from game theory. Aim of game theory is for agents to reach equilibrium. example : rock, paper, scissors. If both players choose their moves uniformly at random then even if they change their strategy it wouldn't affect the outcome & they would tie (equilibrium).

Generators & discriminators have their own loss functions. So for GAN there is a value function. Generator wants to minimize the value function & discriminator wants to maximize the value function. we attain equilibrium when value func is at minimum for generator & at max for discriminator. This is a saddle point.

If we analyse game of GAN's the local member discriminator occurs when it incorrectly thinks that input is real rather than fake. This probability is given by ratio of data density at input & model density induced by generator at input. This ratio determines how much prob mass comes from data.



Silky relu is used as they make sure that gradient flows through entire architecture. For GAN's 2 optimization algo. must run simultaneously. d-loss, g-loss.

Discriminator uses BCE loss as its just classification. As sigmoid output are unstable at time we multiply labels by output values closer to 1 but ≤ 1 (0.9) so that output goes to 0 if misclassified accordingly (eg: 0.9).

Generator also uses BCE loss but labels flipped.

To scale GAN's to work for

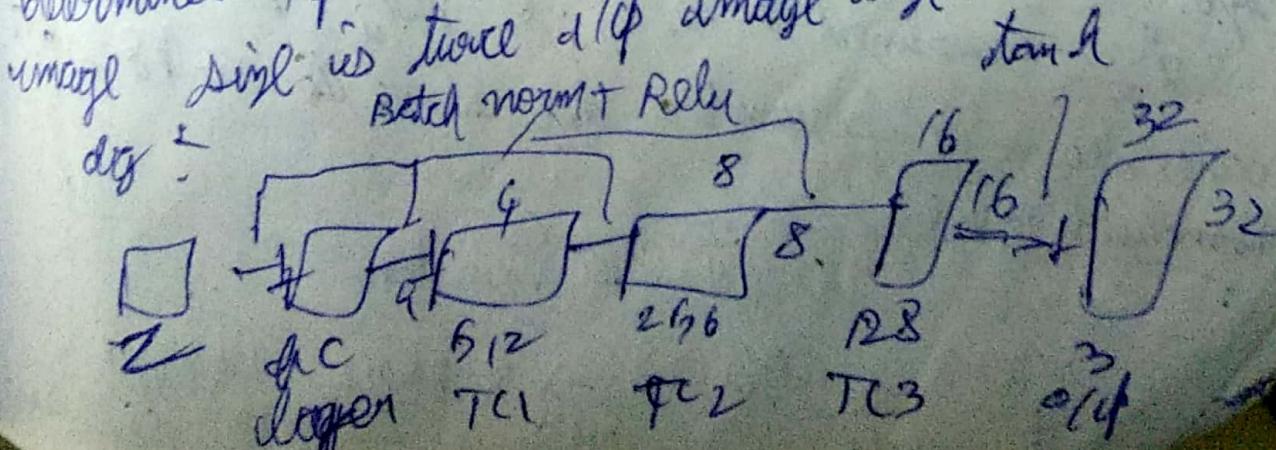
large images we are convolutional layers.
 discum { matmul, + sigmoid
 generator } conv + relu
 conv + tanh
 conv + softmax
 reshape
 z

DCGAN - discriminator :-

$d/dp \rightarrow \text{conv}_1 \rightarrow \text{conv}_2 + \text{conv}_3 + \text{flatten}$

Every hidden layer fed $\sigma(\text{sig} \leftarrow \text{fc layer})$
 has leaky relu & batch norm. Batch normalization scales the
 output of the layer to have mean 0 &
 variance 1.

DCGAN generator - As usual input is
 a latent vector z which is fed to GAN to generate
 realistic images. We need to upsample the inputs
 to generate the realistic image. Instead of linear
 interpolation we can use transposed convolutions.
 For transposed convolution, for every 1 pixel in
 o/p we move 2 pixels in d/dp image. stride -
 determines o/p image size. For stride of 2 o/p
 image size is twice d/dp image size.



Batch normalization

In training we normalize each layer's o/p by its mean & std dev. This is batch statistics. Batch norm normalizes o/p of previous layer by sub the batch mean dividing by batch std dev.

Internal covariate shift: It refers to change in distribution of inputs to next layers. Training is more efficient when dist of inputs to next layer is similar.

Mean $\hat{}$ average value of o/p of any layer by passing them non-linear act func -

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

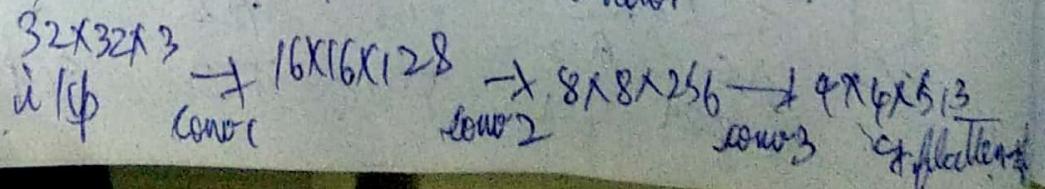
$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\text{Normalized } \hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

ϵ is a small constant that is added to avoid divide by 0. It also slightly increase variance for each batch.

\hat{x}_i - normalized value, $y_i \leftarrow \gamma \hat{x}_i + \beta$. γ & β are learnable parameters that serve to scale & shift normalized value.

DCGAN - discriminator



Fix 2 pic & cycle gan -
cycle gan is used for image to image translation

Image cop date (Training data)

GAN → New images
cycle gan can also be used random noise for image colorization, super resolution. It is basically also used to transform (map) image in one domain to another domain.

Designing loss functions.

Here we need to compare true & generated image.
eg - Image colorization, we can use Euclidean distance comparing pixel by pixel, $L(g, I) = \frac{1}{2h, w} \sum_{i,j} (g_{ij} - I_{ij})^2$
In practice this is too simple.

Off Nectab

Generator tries to generate fake images to fool Disc P

$$\min_{\text{gen}} \max_{\text{disc}} E_{z, x} [\underbrace{\log(D(\text{gen}(z)))}_{\text{fake}} + \underbrace{\log(1 - D(x))}_{\text{real}}]$$

If discriminator outputs 0.5 it means generated & real image are indistinguishable. It means generator has learnt or mapping to map a latent vector \mathbf{z} to a realistic image.

Latent vector is also called noise

space, a compressed feature level representation of image. For image to image translation the dataset is paired image data.

$$D \quad \square \quad G: x_d \rightarrow y_d$$

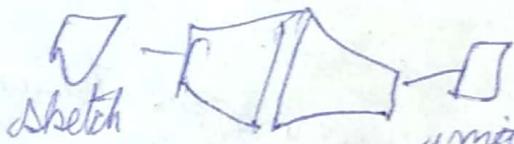
$x_d \quad y_d$ Gx should be indistinguishable

from real image y .

Pixel 2 pixel : Rather than like an regular GAN, pix2pix architecture involves some modifications in generator. For translating an input sketch to an image.

The learned

feature representa-



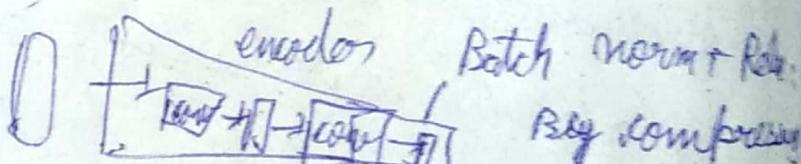
Sketch

feature
representation

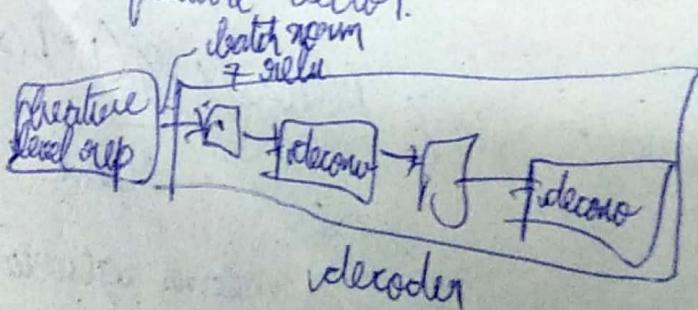
image will help us in gener-

ating image from the sketch

Encoder +



input the encoder captures the content of the image in small feature vector.



Discriminator

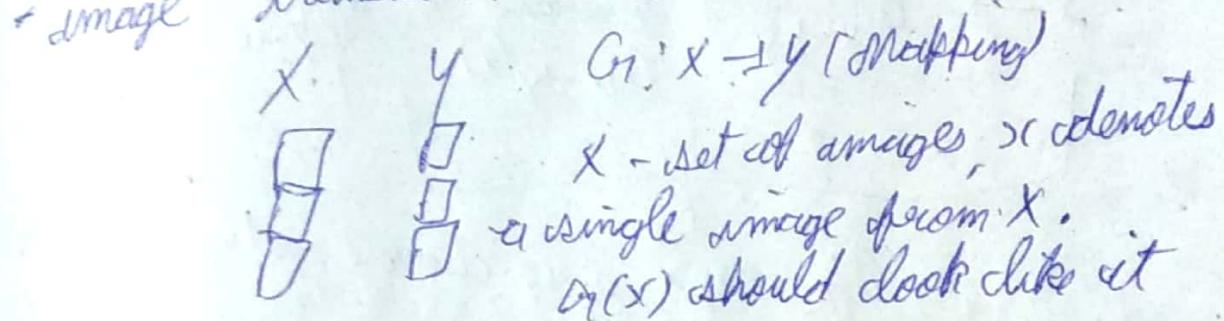
pix2pix

Instead of taking a single

image as input & outputting a prob for real or fake the discriminator takes a pair of images & outputs 0.1 if real fake pair

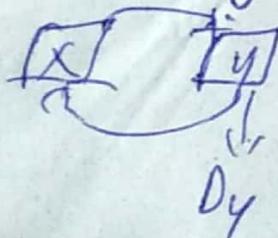
$$\min_{G} \max_{D} \mathbb{E}_{\substack{x \sim p_{\text{data}}(x) \\ z \sim p_z(z)}} [\log D(x, G(z))] + \mathbb{E}_{\substack{x \sim p_{\text{data}}(x) \\ z \sim p_z(z)}} [\log (1 - D(x, G(z)))]$$

cycle gan's & pix2pix works only with paired images & does mapping $G_1: X \rightarrow Y$. But paired image collection is time consuming. So instead of pairs of images we have sets of input & output. Cycle consistent networks help in image translation in this case.



comes from set Y . Unpaired data such as mode collapse. For example neural net might map multiple horses to same zebra image.

Adversarial loss can be used as a constraint we need to add another constraint. Inverse mapping $G_{1Y} \circ G_1: G_1Y \rightarrow X$. So this ensures cycle consistency



$$G_{1X \rightarrow Y}: G_1X \rightarrow Y$$

$$G_{1Y \rightarrow X}: Y \rightarrow X$$

$$G_{1Y \rightarrow X}(G_{1X \rightarrow Y}(x)) \approx x$$

Consider style transfer. It can be considered as mapping from X domain to Y domain. Here cycle consistency is needed. For this we can use 2 discriminators, $D_Y \circ D_X$. D_Y takes

in real image and generated image $G_{x \rightarrow y}(x)$ as if y
 \rightarrow_x takes in $G_{y \rightarrow x}$ x takes in generated image.
 $G_{y \rightarrow x}(y)$ & x (real image) in x domain as input

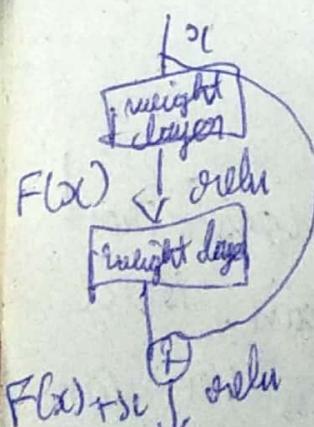
$$x \rightarrow G_{x \rightarrow y} \rightarrow \text{generated} \rightarrow G_{y \rightarrow x} \rightarrow x$$

$G_{x \rightarrow y}(x)$

If $x \neq y$ are same then forward
cycle consistency holds good. Diff between $x \neq y$
is called cycle consistency loss.

cycle consistency = adversarial + forward and
loss losses backward cycle
losses.

Generator for cyclegan. Here we use residual blocks which only connect x & y .
 \rightarrow_x one layer with input of an earlier layer.



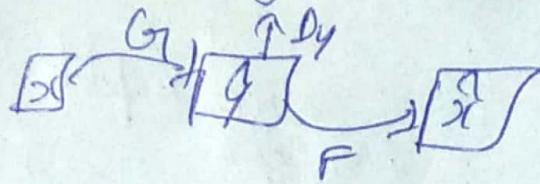
residual function: $F(x) = M(x) - x$, diff betw mapping applied to x & original x .

$F(x)$ - typically 2 some layers + norm layer + relu inbetween.

$$\therefore M(x) = F(x) + x, y = F(x) + x$$

Idea is that it is easier to optimize dual function than optimizing $M(x)$.

D & G_t losses - In addition to disc loss we need to calculate cycleconsistency loss. This loss ensures reproducibility



$G - X \rightarrow Y$ $F - Y \rightarrow X$.

Here as opposed to cross entropy loss, least squares can be used. GANs treat disc. as classifier with sigmoid cross entropy loss function. This can lead to vanishing gradient problem. To overcome this we can use least squares loss func.

discriminator loss - mean squared errors between 0/1 of discriminator given an image & target value 0 or 1, depending on whether it should classify image as fake or real. generator losses - If generator we might generate fake images that look like they belong to set of X images but are based on real images in Y & vice versa. Real class can be computed by looking at 0/1 of disc as its applied to fake images. Generator aim is to classify fake images as real images.

cyclegan:

discriminator training:

- 1) compute discriminator loss D_X on real images
- 2) generate fake images that look like domain X based
- 3) compute fake loss for D_X

- 4) compute total loss & perform backprop & D_X opt.

- 5) Repeat 1-4 switching domains

- 6) Add up generator & reconstruction loss & perform backprop

Generator training

- 1) Generate fake images in X based on images in Y .
- 2) Compute generator loss based on how D_X responds to fake X .
- 3) Generate reconstructed Y' images based on fake X images in step 1.
- 4) Compute cycle consistency loss by comparing reconstruction with real Y images.

- 5) Repeat 1-4 switching domains
- 6) Add up generator & reconstruction loss & perform backprop