

\Rightarrow DL is part of a broader family of ML methods based on artificial neural network (ANN) with representation learning

\Rightarrow Representation learning: ML we need to give feature while AI detect feature automatically from given duty

Ex Dog or cat ML need features like size, colour (need feature eng)
 \Rightarrow AI detect extract feature automatically (no need to feature engineering)

\Rightarrow DL algo uses multiple layers to progressively extract higher-level features from raw inputs Ex in image lower layers may identify edges while higher level may identify the concepts relevant to human such as digits or letters or face

ML & DL is used for find relation b/w input and output mostly statistical technique for this while DL follows logical

Normal network

Object detection
Image segmentation
Speech recognition, chat box

Some complex problems

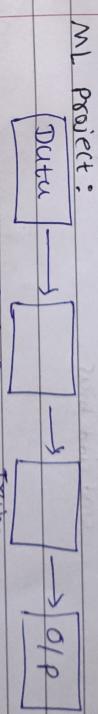
Like Images object detection

Neural network problems

ML model

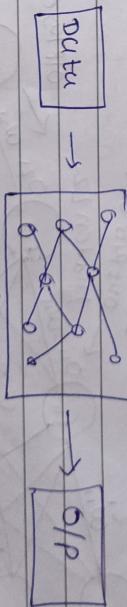
High performance

How do data science techniques scale with amount of data?



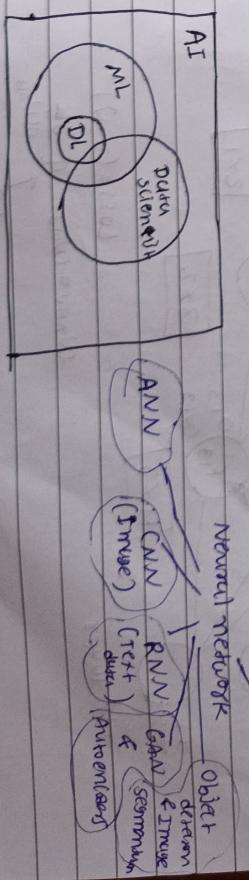
Feature selection
(Feature extraction)

DL project



This combine

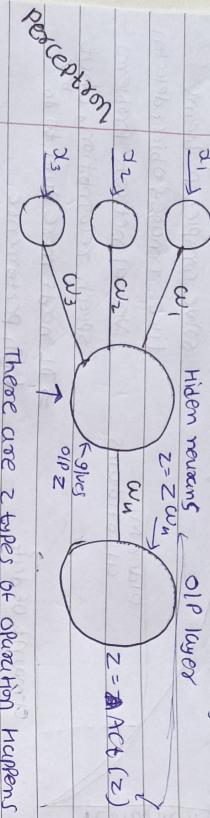
Feature extraction &
Classification or Regression



* **Neural network work**

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

$$z = \text{Act}(y)$$



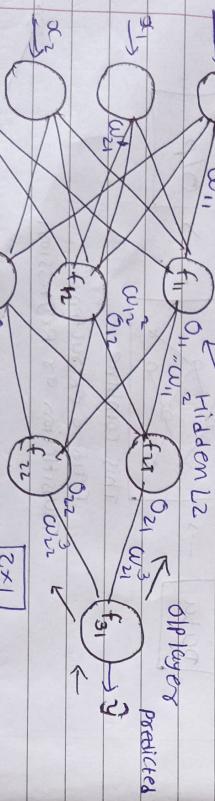
Step 1: $\Theta \rightarrow$ Step 2: $\Theta \rightarrow \text{Act}(\sum_i w_i x_i)$

$$\frac{\partial L}{\partial w_{i3}} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{i3}}$$

w_{i3} →

* **Multilayer Neural Network**

weights
hidden L1 outputs or activation function



→ w_{21} is impacting o_{31} so $\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{21}}$

$$\rightarrow w_{21}^2 = w_{21}^2 - \eta \frac{\partial L}{\partial w_{21}^2}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{11}} + \frac{\partial L}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}} + \frac{\partial L}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}}$$

one const
+ sec con
1st den

$$\text{I/P layer}$$

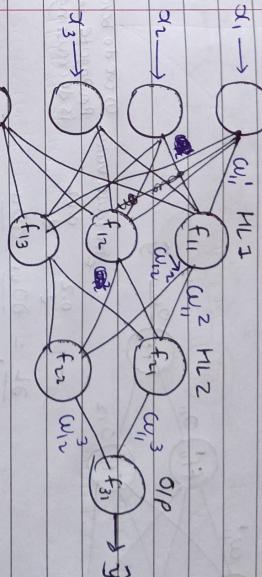
$$\text{LOSS} = (y - \hat{y})^2$$

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial \text{LOSS}}{\partial w_{\text{old}}}$$

LR

* **Chain Rule**

$$\text{grad}_n = \frac{\partial L}{\partial w_{nk}}$$



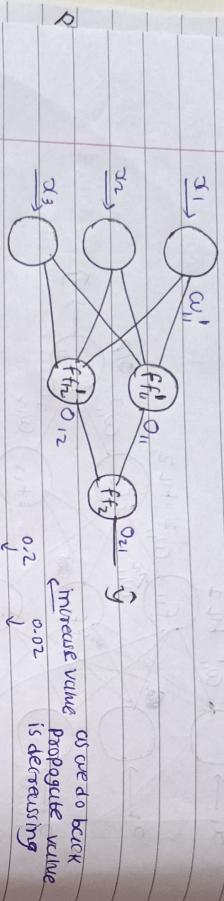


PAGE NO.: _____

PAGE NO.: _____

* Vanishing gradient Problem

$$0 < \frac{\partial L}{\partial z} \leq 0.25$$



$$\Rightarrow \text{For sigmoid function } \sigma(z) \text{ range bet'n 0 to 1}$$

$$\text{and } \sigma'(z) \text{ from } 0 \text{ to } 0.25$$

So as the hidden layer increase $\frac{\partial L}{\partial w_{11}}$ decreases

and w_{11} will decrease

so sigmoid function gives vanishing gradient problem

Some Problem Ex tanh

tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

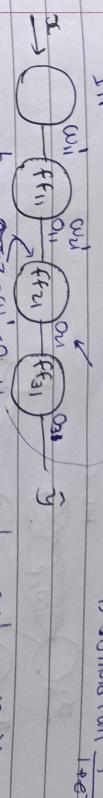
but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

but tanh has derivative 1 - tanh^2 which is less than 1

* Exploding gradient Problem

$$0 < \sigma'(z) < 0.25$$



$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}}$$

let activation fun
is sigmoid fun $\frac{1}{1+e^{-z}}$



$$z_1 = w_1 A_0 + b_1 \rightarrow A_2 = \hat{y} = w_2 (w_1 A_0 + b_1) + b_2$$

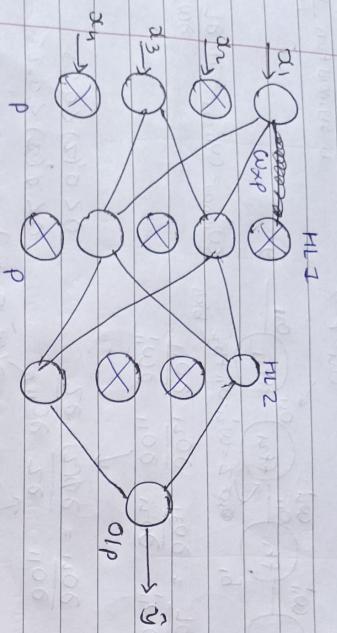
$$A_1 = g(z_1) = z_1$$

$$A_2 = g(w_2 A_0 + w_2 b_1 + b_2)$$

$$y = w_3 (w_2 A_0 + w_2 b_1 + b_2) + b_3$$

without or linear activation function only use on linear regression
and logistic regression model can't work on non-linear regression
without non-linear fun we cannot do non-linear data classification

* Drop out & regularization (Reduce overfitting)



Drop out ratio $p \Rightarrow$ to improve overfitting

$0 \leq p < 1$ we use drop out method

In this method we deactivate some neurons & also activate some neurons.

For test time we activate all features and all neurons and multiply \vec{w} with every weight

\Rightarrow For P value we can use hyperparameter optimization (ex: cross validation)

\Rightarrow For every epoch we remove p fraction of neuron in each

HL randomly.

\Rightarrow using this method we can avoid overfitting by not focusing on single thing & focusing on every things.

\Rightarrow this method is like random forest here we uses diff

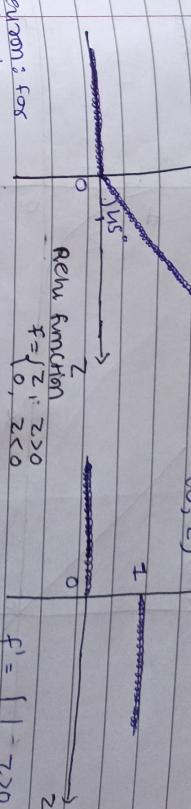
neural networks instead of diff decision trees

\Rightarrow Total possible neural network 2^n (without using p)
approx 100-diff neural network possible for 100 epochs

without or linear activation function only use on linear regression
and logistic regression model can't work on non-linear regression
without non-linear fun we cannot do non-linear data classification
problems that are not linearly separable

Rectified Linear Unit (ReLU) Activation Function

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



dead neuron: for

any given input

$o_p = 0$

Region for dying neuron: $z_i < 0$

If high LR

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

is either 0 or 1

\Rightarrow when one of derivative is zero then $w_{ij} = 0$ when

This cause dead neuron or dead activation function

\Rightarrow to overcome this problem we use leaky ReLU. In this we will add some value (very small) so stop for $f < 0$

have some value > 0 because it $z < 0$ that means $w_{ij} < 0$ or $b_{ij} < 0$ and since our neuron is normalized so $z < 0$. One small and for any input $z < 0$ \Rightarrow Solution

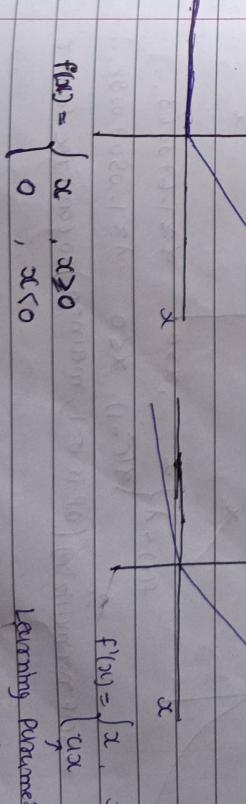
Leaky ReLU $f = \begin{cases} z, & z > 0 \\ 0.01, & z \leq 0 \end{cases}$

use variance of relu $f'(x)$

$$f'(x) = \begin{cases} 1, & x > 0 \\ 0.01, & x \leq 0 \end{cases}$$

Learning Parameters

when $c = 0.01$ It's become



Leaky ReLU

If $c = 0$ ReLU
If c is learnable parameter
& become prelu

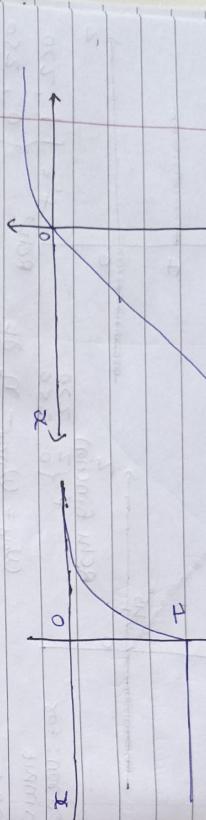


Exponential Linear units (ELU) function (activation)

sigmoid activation function

PAGE NO.:

$$f(x) = \begin{cases} x, & \text{if } x \geq 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \end{cases}$$



$$f(x) = \begin{cases} x, & \text{if } x \geq 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \end{cases}$$

- ⇒ No dead ReLU issues, continuous & differentiable
- ⇒ mean of op is close to 0, zero centered
- ⇒ computationally expensive because $\alpha(e^x - 1)$ take time

\Rightarrow Not zero centered data will we reduce the efficiency of weight update.

\Rightarrow computationally expensive ($\because \frac{\partial}{\partial x} \alpha(e^x - 1)$)

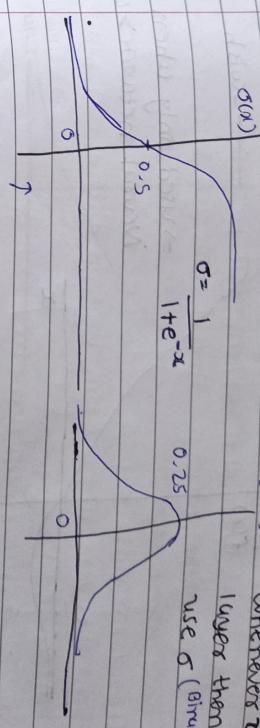
\Rightarrow $\sigma'(x)$ is range to 0 to 0.25 so vanishing gradient problem occurs.

$\frac{\partial L}{\partial w_i} = \frac{1}{2} \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_i}$ $0 < \alpha, b_1, b_2 < 1 \Rightarrow$ so both gradient are +ve or $\rightarrow 0$ and for all weight we need to do it together and converge slowly

\Rightarrow $\frac{\partial L}{\partial w_i} = \frac{1}{2} \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_i}$ \Rightarrow op from activation not normalized (from 0 to 1)

\Rightarrow hence this problem occurs (from zero centered)

* tanh activation function

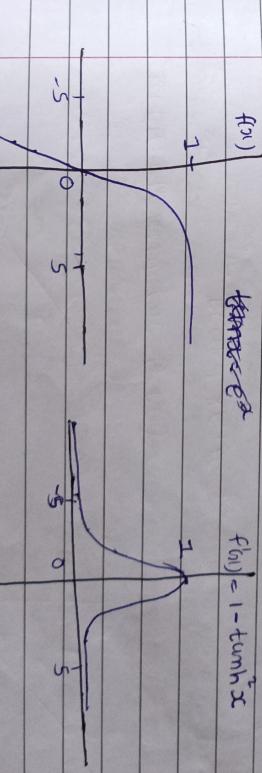


$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

PAGE NO.:

* SELU (scaled Exponential Linear unit)

$$f(x) = \lambda \begin{cases} x, & x \geq 0 \\ \alpha e^x, & x \leq 0 \end{cases}$$



$$f(x) = \lambda \begin{cases} x, & x \geq 0 \\ \alpha e^x, & x \leq 0 \end{cases}$$

$$f(x) = \lambda \begin{cases} x, & x \geq 0 \\ \alpha e^x, & x \leq 0 \end{cases}$$

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

- ⇒ almost like sigmoid but $f(x)$ ranges to 0 to 1
- ⇒ delta is zero centered which is better than sigmoid
- ⇒ sum after vanishing gradient occurs.

QUESTION

ANSWER

Swish (Asset-Gated) Activation Function

PAGE NO.: _____

Softplus Activation Function

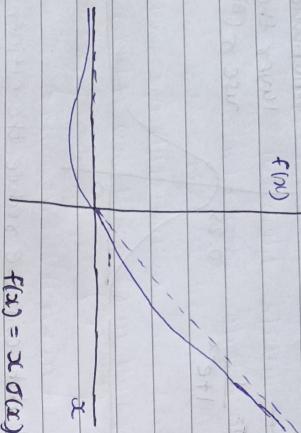
PAGE NO.: _____

$f(x) = \text{ReLU}$

$f(x) = \ln(1 + \exp x)$

Softplus

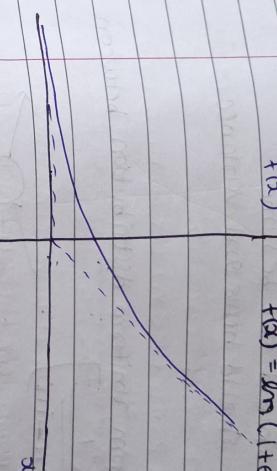
Rectifier



$$f(x) = x \sigma(x)$$

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

⇒ used only when
Neural network > no layers



★ Softmax

$$S(x_i) = \frac{e^{x_i}}{\sum_{k=1}^K e^{x_k}} = \frac{e^{x_i}}{\sum_{k=1}^K e^{x_k}}, \quad i=1, 2, 3, \dots, K$$

⇒ Softmax function is used when we have more than 2 op layers (For 2 op layers σ)

OP layers

$$S(x_1) = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}}$$

$$S(x_2) = \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}}, \quad S(x_3) = \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}}$$

$$S(x_4) = \frac{e^{x_4}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}}$$

$$\Rightarrow \sum_{i=1}^K S(x_i) = 1$$

For C + OOP, the best choice is ReLU
more complex and difficult to implement
also requires additional memory space



Various weight initialization Techniques

key points i) weights should be small

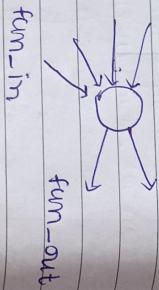
ii) weights should not be same

iii) weights should have good variance

i) uniform distribution

weights are taken from uniform distribution between

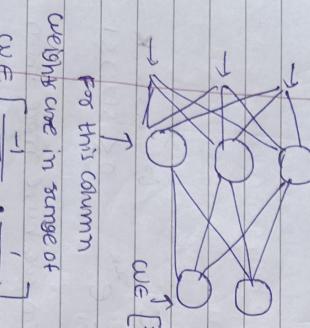
$$\text{range } \left[\frac{-1}{\sqrt{\text{fun_in}}}, \frac{1}{\sqrt{\text{fun_in}}} \right]$$



ii) Xavier / Goopt Uniform

$$\sigma = \sqrt{\frac{2}{(\text{fun_in} + \text{fun_out})}}$$

normal distribution



for this column

weights are in range of

$$w \in \left[\frac{-1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]$$

3) He init work nicely with ReLU

i) He uniform

$$w_{ij} \in U \left[-\frac{6}{\sqrt{\text{fun_in}}}, \frac{6}{\sqrt{\text{fun_in}}} \right]$$

ii) He Normal

$$\sigma = \sqrt{\frac{2}{\text{fun_in}}}$$

\Rightarrow work nicely with $\sigma = (\text{ReLU})$

2) Xavier / Goopt work nicely with σ & tanh

i) Xavier / Goopt Normal

$$w_{ij} \in N(0, \sigma)$$

ii) normal distribution

Stochastic Gradient Descent (SGD)

$$\text{One step} = \Delta w_j - \eta \frac{\partial L}{\partial w_j}$$

at time

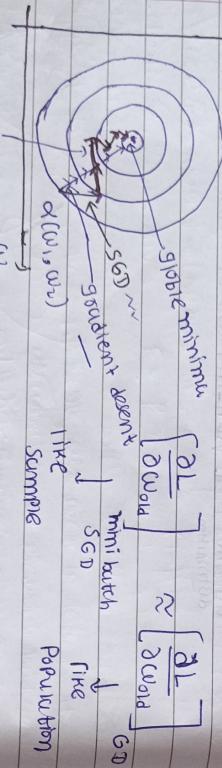
Total no of data points

\Rightarrow If we consider only one data point and finding \hat{y} and update the weight is called SGD.

\Rightarrow If we consider K data points ($K < n$) at a time and finding $L = \sum_{i=1}^K (y_i - \hat{y})^2$ and update the weights is called mini SGD (no of updates per epoch = batch).

\Rightarrow If we take all n data at a time and find $L = \sum_{i=1}^n (y_i - \hat{y})^2$ and update weights then its called gradient descent (GD).

(w) Top view

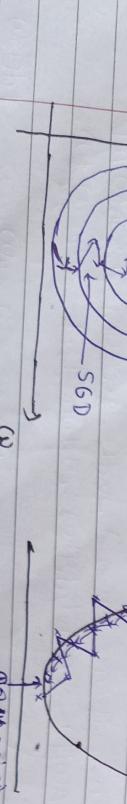
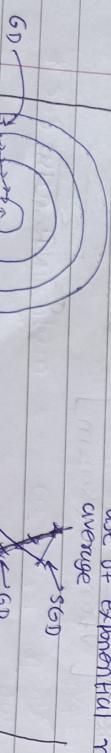


In 1 epoch - n iterations
mini SGD



SGD with momentum

(w) Top view \Rightarrow Here we remove noise with use of exponential moving average



SGD with momentum

Next \rightarrow

be understood as the number of times the algo scans entire data whereas iteration is number of times a certain batch passed via algorithm. Ex. a dataset having 20 samples, batch size = 5, epoch = 5 there will be 5 batches (row) each batch will pass by the algo

\Rightarrow For Batch size = 1000 and total data 10000

let we have 20 epochs

so in 1 epoch \rightarrow iteration will be $\frac{10000}{1000} = 10$ iteration

to focus on a particular part of population

so total iteration = $10 \times 20 = 200$

happens

\Rightarrow In case of GD there are 1 iteration in one epoch $\sum_{i=1}^n (y_i - \hat{y})^2$

\Rightarrow For SGD iteration = no of SGD

iterations = no of records

(y - \hat{y}) \rightarrow GD

in 1 epoch - n iterations mini SGD

batched

SGD

time



PAGE NO.: 11

\Rightarrow We have some delta point at different - different time intervals

time	t_1	t_2	t_3	-
delta	u_{t_1}	u_{t_2}	u_{t_3}	-
after denoising	v_1	v_2	v_3	-

higher perimeter

$$v_t = \beta v_{t-1} + (1-\beta) s_t$$

Smooth then
value at point t

\Rightarrow Choice of β will determine
how many delta point that we

avg the value of v_t as

shown below

$$v_1 = c_1$$

$$v_2 = \beta v_1 + (1-\beta) u_2$$

$$v_3 = \beta v_2 + (1-\beta) u_3$$

$$= \beta(\beta v_1 + (1-\beta) u_2) + (1-\beta) u_3$$

$$v_n = (\beta)^n [v_1 + \beta^{n-1} u_2 + \dots + u_n]$$

$$\beta = 0.9, \text{ ave over 10 delta point}$$

$$\frac{\beta}{1-\beta} = 10$$

$$\Rightarrow$$
 Higher β value means previous point
take more importance

$$v_1, v_2, \dots, v_n$$

$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$v_n = (1-\beta) v_{n-1} + \beta v_{n-2} + \dots + \beta v_1 + \beta u_n$$

$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_t}$$

$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$w_t = w_{t-1} - \eta v_{t-1}$$

$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$w_t = w_{t-1} - \eta v_{t-1}$$

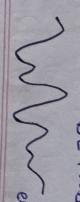
$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$w_t = w_{t-1} - \eta v_{t-1}$$

$$\Rightarrow$$
 Here adding the fraction of the previous update to the current update will make process a bit faster.

$$w_t = w_{t-1} - \eta v_{t-1}$$

Non-Convex Func increases occurs for DL Prob $\beta \approx 0.95$ for smoothening of curve for ML problems
Ex: Linear logistic regression



$$b_t = b_{t-1} - \eta v_{dt}$$

where $v_{dt} = \beta v_{db_{t-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$

$\Rightarrow \beta$ is showing that how fast velocity equal to $(\frac{1}{-\eta})$ plus to velocity

$$\Rightarrow$$
 Instead of ~~computing~~ $\frac{\partial L}{\partial w_t}$ we replace it as v_{wt}

$$\Rightarrow$$
 ~~minibatch~~ ~~for every iterations~~ $\frac{\partial L}{\partial w_t}$

~~and over for every batch we compute~~

\Rightarrow On iteration t inside epoch compute derivative of initially v_{dw} and v_{db} and updates weight and bias $w = w - \eta v_{dw}$ $b = b - \eta v_{db}$

\Rightarrow In case of minibatch gradient descent when we update the model parameter after iterating through all the

delta points in the given batch, thus the direction of the update will have some variance which leads to oscillations.

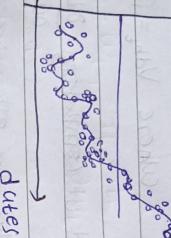
\Rightarrow Due to this oscillation it is hard to converge and it slow down the process of attaining it. To combat this we use momentum.

\Rightarrow Momentum helps us in not taking the direction that does not lead us to convergence.



- ⇒ EWA is technique to find hidden trend in time series
- base data
- Recent point has more weight than older one
- ⇒ at any fix point, with time weight is decreases.

Temperature
dates



Nesterov Accelerated Gradient (NAG)

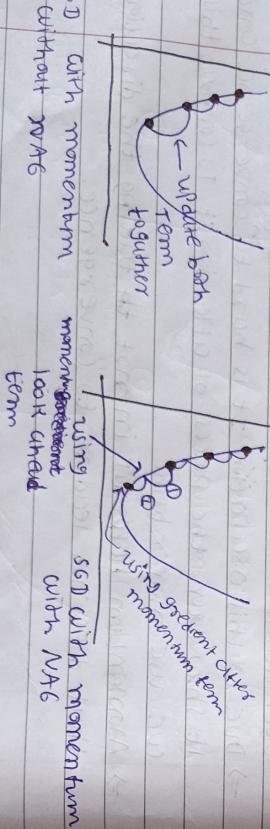
$w_t = w_{t-1} - \eta \nabla w_t$ depends on two terms.

①

history of velocity + gradient at that point

⇒ For NAG update one by one instead of both together.

- First ⇒ There is momentum term ② is applied and then we apply gradient term at that point.



SGD with momentum

using SGD with momentum with NAG without NAG

look ahead

team

- Disadvantage of NAG
- ⇒ sometimes it stuck with local minimum problem.

Adaptive gradient descent (Adagrad optimizers)

Adadelta and RMSprop

(Root mean square prop)

PAGE NO.:

$$w_t = w_{t-1} - \eta_t \frac{\partial L}{\partial w_{t-1}}$$

main idea behind Adadelta
and RMSprop is same to deal with vanishing LR by taking the weighted avg of gradient step

$$\eta_t = \frac{\eta}{\sqrt{S_{w_{t-1}} + \epsilon}}$$

η = initial LR
 $S_w = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$

$\eta_t = \frac{\eta}{\sqrt{S_{w_{t-1}} + \epsilon}}$ ← Exponentially weighted avg (EWMA)
 $\beta = 0.9, 0.95$

$\Rightarrow \eta_t'$ is decreases as time goes
 ϵ = very small & very number to avoid $\eta_t = 0$ effect

\Rightarrow after some time η_t very large and η_t' decrease and become very small so $w_t \approx w_{t-1}$ this is disadvantage of Adagrad optimizers.

\Rightarrow Adagrad is work best for sparse data (more without normalization) and data with different scale (without normalization)

$$S_{ab} = \beta S_{ab} + (1-\beta) \left(\frac{\partial L}{\partial b} \right)^2$$

\Rightarrow good for sparse (zero) data
 \Rightarrow not working well with smooth data so we use Adagrad for this type of data

$$S_{ab} = \beta S_{ab} + (1-\beta) \left(\frac{\partial L}{\partial b} \right)^2$$

then we update weight and bias

$$C_{ab} = C_{ab-1} - \eta_t' \frac{\partial L}{\partial b_{t-1}}$$

$$w_t = w_{t-1} - \frac{\eta_t}{\sqrt{S_{ab} + (1-\beta) \left(\frac{\partial L}{\partial b} \right)^2 + \epsilon}} \frac{\partial L}{\partial w_{t-1}}$$

$$\eta_t' = \frac{\eta}{\sqrt{S_{ab} + \epsilon}}$$

For Batch GD
 η_t' calculate for whole batch

\Rightarrow For sparse data η_t' mostly zero and we get $\eta_{avg} = m$ There are no zero (1) so b update each time so movement is jumping

\Rightarrow To avoid this we use Variable m

\Rightarrow Adagrad is not working on non-convex optimization but

RMS works \nwarrow
RMS has no disadvantage.

Adam optimizer (Adaptive momentum Estimation)

Initialising with zero initial values
 \Rightarrow $V_{dw} = 0$, $V_{db} = 0$, $S_{dw} = 0$, $S_{db} = 0$
 on iteration t ,
 compute $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial b}$ using current minibatch

$$\left. \begin{array}{l} V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \frac{\partial L}{\partial w} \\ V_{db} = \beta_1 V_{db} + (1-\beta_1) \frac{\partial L}{\partial b} \end{array} \right\} \text{momentum functions}$$

$$\left. \begin{array}{l} S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \left(\frac{\partial L}{\partial w} \right)^2 \\ S_{db} = \beta_2 S_{db} + (1-\beta_2) \left(\frac{\partial L}{\partial b} \right)^2 \end{array} \right\} \text{RMS Prop}$$

$w_p = w_{t-1} - \frac{\eta \cdot V_{dw}}{\sqrt{S_{dw} + \epsilon}}$ combine
 momentum \rightarrow
 (for convergence)
 smoothing function Adam weight

$$b_t = b_{t-1} - \frac{\eta \cdot V_{db}}{\sqrt{S_{db} + \epsilon}} + \text{RMSProp}$$

$$\text{Bias correction } V_{dw}^{\text{correct}} = \frac{V_{dw}}{1-\beta_1^{(t)}} \text{ sum for } V_{dw}$$

on momentum

$$\text{Bias correction } S_{dw}^{\text{correct}} = \frac{S_{dw}}{1-\beta_2^{(t)}} \text{ sum for } S_{dw}$$

1.0000000000000000

Minibatch

Stochastic gradient descent

• All gradients with size 2000

Different types of Loss functions

\Rightarrow Loss function: For only one data point $(y - \hat{y})^2$ is loss
 \Rightarrow cost function: For batch size t , $\sum_{i=1}^t (y_i - \hat{y}_i)^2$ is cost function
 \Rightarrow Error function: Loss & cost function called as error fun

- i) squared error ~~connection~~ Loss

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

$$\text{Cost function} = J = \frac{1}{t} \sum_{i=1}^t (y_i - \hat{y}_i)^2 \leftarrow \text{MSE}$$

mean squared error Loss

Advantages: In form of quadratic eqn

$$c x^2 + b x + c$$

↑

- i) For quadratic eqn has only one global minimum

→ MSE Loss penalize the model for making larger errors by squaring them.

Disadvantage:

- i) It is not robust to outliers (penalize more)

- ii) Absolute error Loss $L = |y - \hat{y}|$ ← modulus operator is non-differentiable

$$J = \frac{1}{t} \sum_{i=1}^t |y_i - \hat{y}_i| \leftarrow \text{MAE}$$

- iii) Cross entropy loss Logistic Regression

$$\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \quad \text{for binary cross entropy}$$

Compare to MSE

(for this need to call softmax)

disadvantage: i) may have local minima, not differentiable at 0

- iv) Huber Loss (combination of MSE & MAE)

$$\text{Loss} = \begin{cases} \frac{1}{2} (y - \hat{y})^2, & |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & \text{otherwise} \end{cases}$$

No of categories
categorical cross entropy

- v) Multiclass Cross Entropy Loss

$$L(x_i, y_i) = -\sum_{j=1}^c y_{ij} \log(\hat{y}_{ij})$$

y_{ij} is one hot encoded target vector

0.1 multiclass $y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{it}]$ total number of features

$y_{ij} = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ element is in class } j \\ 0, & \text{otherwise} \end{cases}$

\Rightarrow In huber loss If $|y - \hat{y}| \geq \delta$ means diff is greater than some value hyper parameter then it is outliers and we avoid unnecessary penalization and we use MAE instead of MSE

IIP

OIP

$$\begin{array}{ll} IIP & OIP \\ f_1 \quad f_2 \quad f_3 & y_1 \quad y_2 \quad y_3 \\ 2 \quad 3 \quad 4 & [1 \quad y_1 \quad 0 \quad y_2 \quad y_3] \leftarrow 3 \text{ bits} \\ 5 \quad 6 & [0 \quad y_1 \quad y_2 \quad 1 \quad y_3] \\ 7 \quad 8 \quad 9 & [0 \quad 0 \quad 1] \end{array}$$

 $\hat{y}_{ij} \Rightarrow$ Soft max function

$$f(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \leftarrow \text{probability value}$$

Cost function
For multiclass = $-\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$
Cross entropy

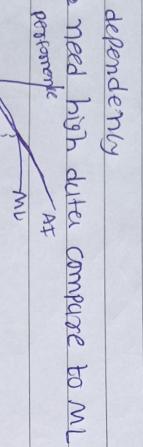
Maximize y_i \rightarrow Minimize $-\log(y_i)$

Deep Learning Vs Machine Learning

PAGE NO.:

PAGE NO.:

- 1) Data dependency
- 2) hardware dependency
- 3) training time
- 4) Feature selection
- 5) Interpretability



- 1) Data dependency
AI ~~need~~ need high data compare to ML
- 2) AI need GPU (- high computations)
ML does not need GPU It can be runs on CPU
- 3) DL is very slow but prediction time is very high

- 1) DL used Representation Learning and extract features automatically while ML doesn't

- 2) AI model is not interpretable
Ex for cat-dog model AI does not ans why It has
Select cat or dog
but ML can explain why using weights of ~~that~~ features

2015 → Google → Tensorflow → difficult to use

Keras → use as
User → Keras → Tensorflow
industry use Tensorflow
most

2016 → Face book → pytorch

AI Researcher uses most
and others one-line
FB made complex
framework Pytorch

drop down Application that convert neural network into
GUI base

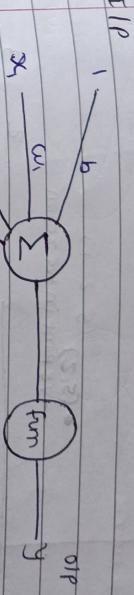
Ex Auto AI (Google)

Apple → Core ML

Microsoft → Custom Vision AI

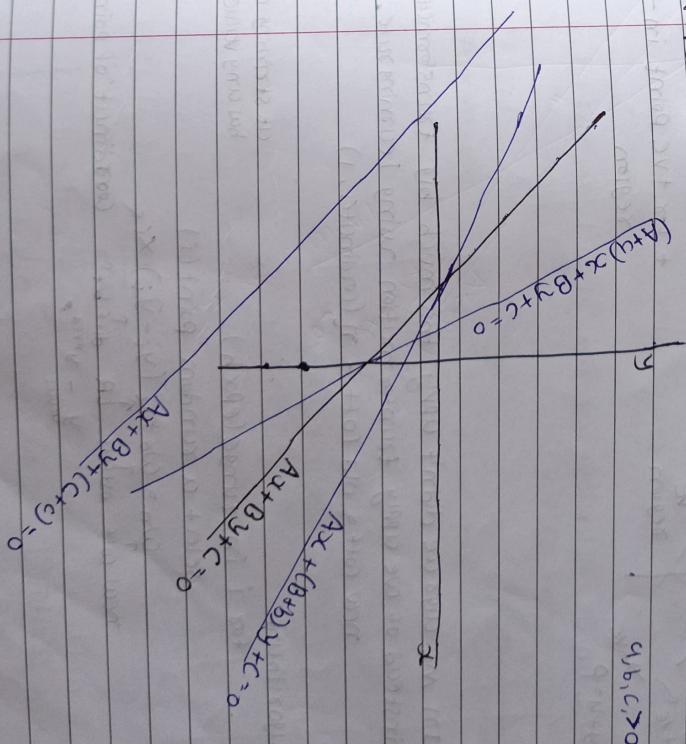
Perception

PAGE NO.:

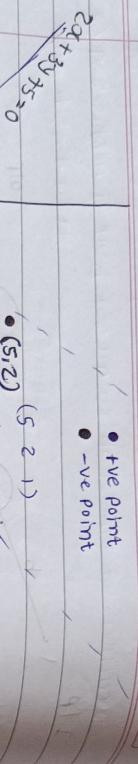


★ Perception trick

$$a, b, c > 0$$



- +ve point
- -ve point

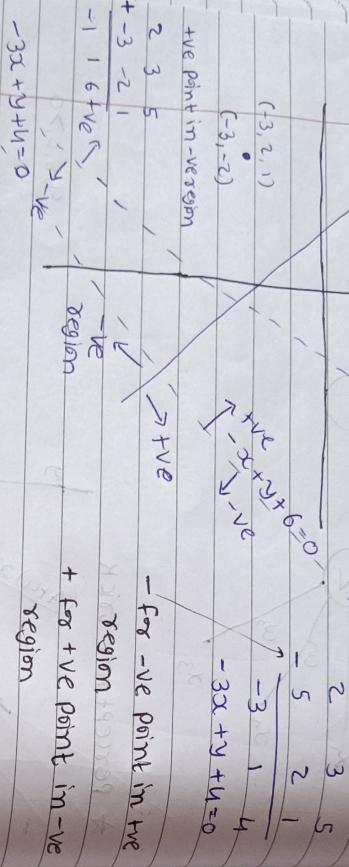


• (5, 2) (5, 2, 1)
-ve point in +ve region

\Rightarrow If this trick choose only those points which are correctly classified then line does not move and gives us wrong result
 \Rightarrow Here there are random fun used so every time we get diff line and we can't decide which is give more precise OP.

\Rightarrow So we use loss function to know how precise our solution is
+ve loss function for Perceptron

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) = L(w_0, w_1, b)$$



+ for -ve point in +ve region
+ for +ve point in -ve region

$$L(y_i, f(x_i)) = \max(0, -y_i f(x_i))$$

Actual OP

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

where $f(x_i) = w_0 + w_1 x_{i1} + w_2 x_{i2} + b$

(Top w column)
(actual feature)

Algorithm:

for i in range(epoch):

 select a random point(i)
 new corr = old corr - n(coordinates, 1)

Algorithm:

for i in range(epoch):

 ut starting w.

 Select a random point(i)

$w_n = w_0 + \eta (y_i - f_i) x_i$
new corr old corr \uparrow diff betw. Corrdinate of point
 $y_{corr} = y_{pred}$

$$w_r = w_r + \eta \frac{\partial L}{\partial w_r}$$

$$b = b + \eta \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f(x_i)} \frac{\partial f(x_i)}{\partial w_i}$$

$$\frac{\partial L}{\partial f(x_i)} = \begin{cases} 0, & y_i f(x_i) \geq 0 \\ -y_i, & y_i f(x_i) < 0 \end{cases} \Rightarrow \frac{\partial f(x_i)}{\partial w_i} = \frac{\partial L}{\partial w_i}$$

\Rightarrow This perception trick is not work in all cases this is rugged method

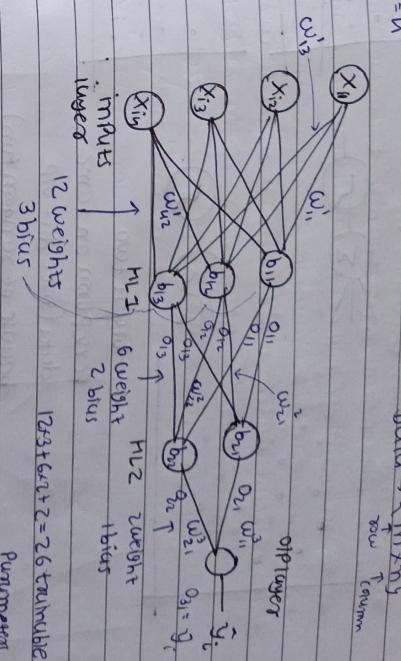
$$\frac{\partial L}{\partial w_i} = \begin{cases} 0, & y_i f(x_i) \geq 0 \\ -y_i x_i, & y_i f(x_i) < 0 \end{cases} \quad \frac{\partial L}{\partial b} = \begin{cases} 0, & y_i f(x_i) \geq 0 \\ -y_i, & y_i f(x_i) < 0 \end{cases}$$

Multilayer Perceptron (MLP) Notation

$n = n$

Data \rightarrow $l \times m \times n$

row
column



\Rightarrow bias b_{ij} , output o_{ij} , input x_{ijk} feature j node i in k^{th} layer

Total no. of features of k^{th} layer

\Rightarrow weight

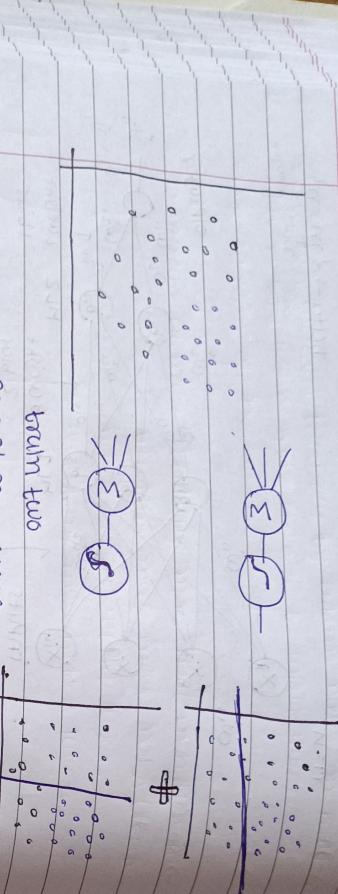
w_{ijk} \leftarrow input for k^{th} layer

Output of i^{th} node in k^{th} layer

node in p^{th} layer

Multilayer Perceptron intuition

PAGE NO.:
CALL NO.:
SUBJECT:
DATE:



brain two
Perception on same
data (gives diff op)

because of random fun

\Rightarrow use linear combination of 2 perception to get final probability

Ex:

GPA = 3
 $i \frac{q}{3} \rightarrow \hat{y}$

Credit = 5
 $\frac{5}{3} \rightarrow \hat{y}$

\Downarrow

$i \frac{q}{n} \rightarrow \hat{y}$

CG = 6
 $\frac{6}{6} \rightarrow \hat{y}$

Credit = 6
 $\frac{6}{6} \rightarrow \hat{y}$

\Downarrow

$i \frac{q}{n}$

CG = 2
 $\frac{2}{2} \rightarrow \hat{y}$

\Downarrow

\Downarrow

$i \frac{q}{n} \rightarrow \hat{y}$

CG = 3
 $\frac{3}{3} \rightarrow \hat{y}$

Analysts are same

IA

$\frac{3}{3} \rightarrow \hat{y}$

\Downarrow

\Downarrow

$i \frac{q}{n} \rightarrow \hat{y}$

Not exact value
(we are assuming)

Final result

as shown above

all point we find

$P(Y_1) = 0.7$

$Z = P(Y_1) + P(Y_2) = 1.5$

$\sigma(Z) = \frac{1}{1+e^{-Z}} = 0.42$

$b = 10$

$Z = 0.2(1) + 0.3(1) + 0.4(1)$

$+ 10$

$\sigma(Z) = \frac{1}{1+e^{-Z}}$

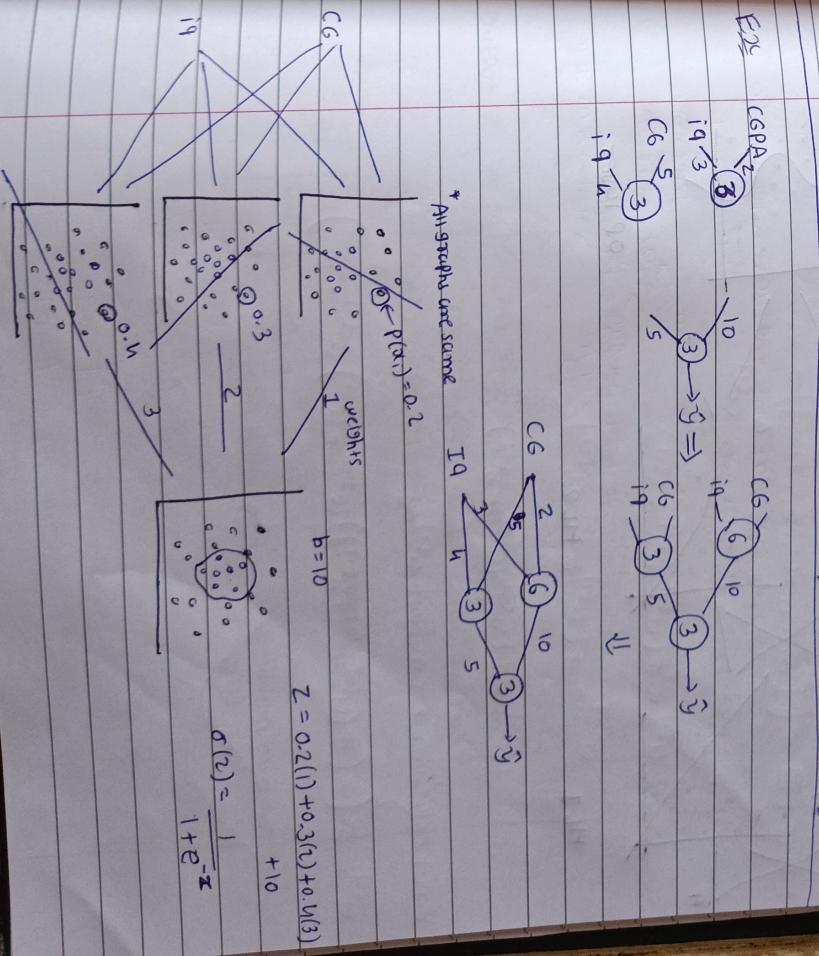
\rightarrow If perception 1's op is dominating then weighted

addition is used and we can also use bias.

$$Z = (0.7)(10) + (0.8)(5) + 3$$

$$\sigma(Z) = \frac{1}{1+e^{-Z}}$$

$$P_1 \quad P_2 \quad P_3 \quad \rightarrow y$$



Forward Propagation

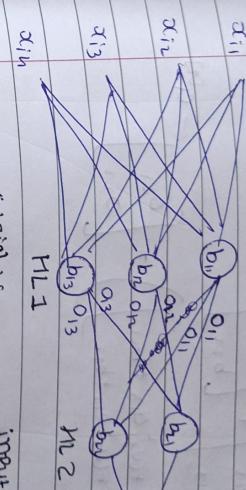
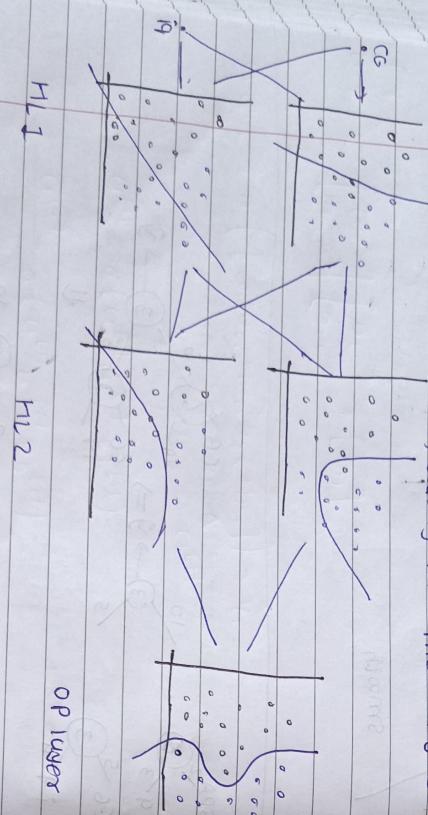
⇒ we can increase accuracy by adding nodes on hidden layers OR adding more feature to input layers
 ⇒ OR increase op nodes

To increase accuracy i) adding more node in hidden layer

ii) adding more input nodes (feature)

iii) increasing output layers

iv) adding more hidden layers



$$\text{For Layer 1: } \begin{aligned} & \text{weights: } \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}^T \quad \text{input for HL1: } \begin{bmatrix} x_{in} \end{bmatrix} \\ & \text{bias of HL1: } b_{in} \quad \text{ML1} \\ & \text{ML2} \end{aligned}$$

$$x_{in} \rightarrow 3 \times 1 \quad 3 \times 1$$

$$w_{in} \rightarrow 3 \times 3$$

$$b_{in} \rightarrow 3 \times 1$$

②

$$\text{output: } \sigma \left(\begin{array}{l} w_{11}x_{in} + w_{12}x_{in} + w_{13}x_{in} + w_{21}x_{in} + b_{in} \\ w_{21}x_{in} + w_{22}x_{in} + w_{23}x_{in} + w_{31}x_{in} + b_{in} \\ w_{31}x_{in} + w_{32}x_{in} + w_{33}x_{in} + w_{41}x_{in} + b_{in} \end{array} \right)$$

$$\text{output: } \sigma \left(\begin{array}{l} w_{11}x_{in} + w_{12}x_{in} + w_{13}x_{in} + w_{21}x_{in} + b_{in} \\ w_{21}x_{in} + w_{22}x_{in} + w_{23}x_{in} + w_{31}x_{in} + b_{in} \\ w_{31}x_{in} + w_{32}x_{in} + w_{33}x_{in} + w_{41}x_{in} + b_{in} \end{array} \right)$$

$$\boxed{\begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix}} \leftarrow \text{activation of layer } \text{ML1}$$

input layer $\begin{bmatrix} x_{in} \end{bmatrix}$ is activation of zeroth layer $\begin{bmatrix} x_{in} \end{bmatrix}$

INPUT FOR
HL2

bias of HL2

PAGE NO.:

PAGE NO.:

weights

ML2

PAGE NO.:

$$\text{For Layer 2} \quad \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{22} \\ b_{32} \end{bmatrix}$$

$$\text{OP} \equiv \sigma \left(\begin{bmatrix} w_{11}^2 o_{11} + w_{12}^2 o_{12} + w_{31}^2 o_{13} + b_{11} \\ w_{11}^2 o_{21} + w_{12}^2 o_{22} + w_{31}^2 o_{23} + b_{22} \end{bmatrix} \right) = \begin{bmatrix} o_{11} \\ o_{22} \end{bmatrix} \leftarrow u^{[2]}$$

For

$$\text{Layer 3} \quad \begin{bmatrix} w_{11}^3 \\ w_{21}^3 \end{bmatrix} \begin{bmatrix} o_{11} \\ o_{22} \end{bmatrix} + [b_{31}]$$

$$\text{output} \equiv \sigma \left(\begin{bmatrix} w_{11}^3 o_{11} + w_{21}^3 o_{22} + b_{31} \end{bmatrix} \right) = \hat{y}_i$$

\uparrow
 $u^{[3]}$

$$u^{[1]} = \sigma(u^{[0]} w^{[1]} + b^{[1]})$$

$$u^{[2]} = \sigma(u^{[1]} w^{[2]} + b^{[2]})$$

$$u^{[3]} = \sigma(u^{[2]} w^{[3]} + b^{[3]})$$

Final expression for output $u^{[3]}$

$$u^{[3]} = \sigma \left(\sigma \left(\sigma(u^{[0]} w^{[1]} + b^{[1]}) \right) w^{[2]} + b^{[2]} \right) w^{[3]} + b^{[3]}$$

SUMMARY OF NOTES

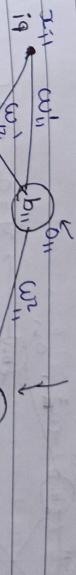
PAGE NO.



Scanned with OKEN Scanner

BACK PROPAGATION

For regression Problem : activation fun = linear



$$\text{MSE} = \frac{1}{2} (y - \hat{y})^2$$

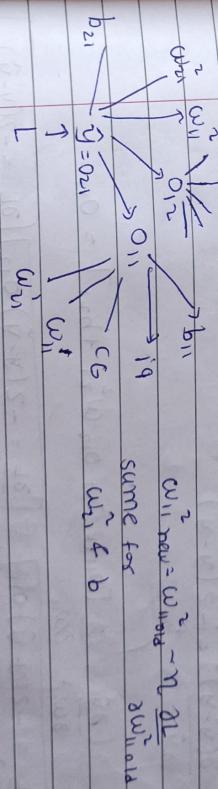
Let initial $w_1 = 1, b_1 = 0$

6 weight + 3 bias = 9 trainable parameters

- loop
 i) select point (row)
 ii) predict op (forward prop) (dot product)
 epochs \times row
 iii) choose loss fun (MSE here)
 iv) update weights and biases using gradient descent

$$w_{11} = w_{11} - \eta \frac{\partial L}{\partial w_{11}}$$

$$b_{11} = b_{11} - \eta \frac{\partial L}{\partial b_{11}}$$



We have to calculate

$$\frac{\partial L}{\partial \theta} \text{ for all } g \text{ parameters}$$

$$\theta = w_{11}, w_{12}, w_{21}, w_{11}^2, w_{12}^2, w_{11} w_{12}, b_{11}, b_{12}, b_{21}$$

}

$$\hat{y} = w_{11}^2 o_{11} + w_{11}^2 o_{12} + b_{11}$$

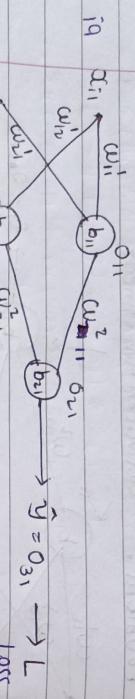
$$\hat{y} = w_{11}^2 [w_{11}^1 x_{11} + w_{11}^1 x_{12} + b_{11}] + w_{12}^2 [w_{12}^1 x_{11} + w_{12}^1 x_{12} + b_{12}]$$

only x_{11}, x_{12} is const
other g are variable

PAGE NO.:

we calculated $\frac{\partial L}{\partial y}$ already

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



$$L(w_{11}, w_{12}, w_{11}', w_{12}', b_{11}, b_{12}, o_{11}, o_{12})$$

now we need

$$\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_{11}}, \frac{\partial L}{\partial w_{11}'}, \frac{\partial L}{\partial w_{12}'}, \frac{\partial L}{\partial b_{12}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}} \rightarrow y \rightarrow w_{11}^2$$

$$\frac{\partial(y - \hat{y})}{\partial y} = 1 - 2(y - \hat{y})$$

$$\frac{\partial y}{\partial w_{11}} = \frac{\partial}{\partial w_{11}} [o_{11} w_{11}^2 + o_{12} w_{12}^2 + b_{11}] = o_{11}$$

$$\frac{\partial L}{\partial w_{11}} = -2(y - \hat{y}) o_{11} \quad \frac{\partial L}{\partial w_{12}} = -2(y - \hat{y})$$

Finally we get

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}} \quad \text{Sum for } w_{11}' \in o_{11}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{12}} \frac{\partial o_{12}}{\partial w_{12}} \quad \text{Sum for } w_{12}' \in o_{12}$$

$$\frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{11}} \frac{\partial o_{11}}{\partial b_{11}}$$

$$\frac{\partial L}{\partial b_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{12}} \frac{\partial o_{12}}{\partial b_{12}}$$

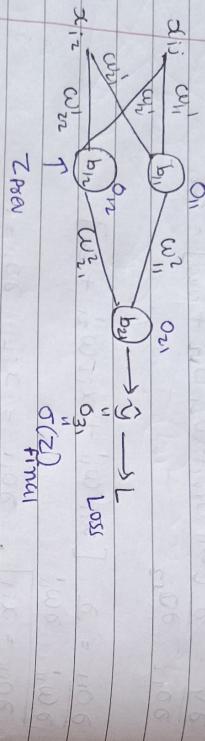
\Rightarrow Hence we calculate minima for g diff parameters

PAGE NO.:



Scanned with OKEN Scanner

For classification problem :
 $L = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$
activation function is sigmoid



must models Z

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{12}}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b_1}$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{21}}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b_2}$$

$$\text{summe for } \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_1}$$

$$L \rightarrow \hat{y} \rightarrow z_f \rightarrow o_{12} \rightarrow z_p \rightarrow w_{12}'$$

$$\frac{\partial L}{\partial w_{11}} = -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial w_{12}} = -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial b_1} = -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12})$$

$$\frac{\partial L}{\partial b_2} = -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12})$$

already calculated

PAGE NO.:

$$L \rightarrow \hat{y} \rightarrow z_f \rightarrow o_{12} \rightarrow z_p \rightarrow w_{12}'$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{11}}$$

$$= -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{12}}$$

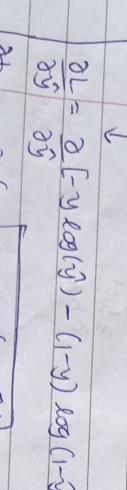
$$= -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b_1}$$

$$= -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12})$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b_2}$$

$$= -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12})$$



$$\frac{\partial L}{\partial o_{11}} = -(y - \hat{y}) \sigma'(z) \frac{\partial L}{\partial o_{12}}$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}} \frac{\partial w_{11}}{\partial z_1}$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{12}} \frac{\partial w_{12}}{\partial z_2}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial o_{11}} \frac{\partial o_{11}}{\partial b_1}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial o_{12}} \frac{\partial o_{12}}{\partial b_2}$$

$$\left[\begin{array}{l} \frac{\partial L}{\partial w_{11}} = -(y - \hat{y}) o_{11} \\ \frac{\partial L}{\partial w_{12}} = -(y - \hat{y}) o_{12} \end{array} \right]$$

$$\left[\begin{array}{l} \frac{\partial L}{\partial b_1} = -(y - \hat{y}) \\ \frac{\partial L}{\partial b_2} = -(y - \hat{y}) \end{array} \right]$$

$$\frac{\partial L}{\partial w_{11}} = -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12}) x_{12}$$

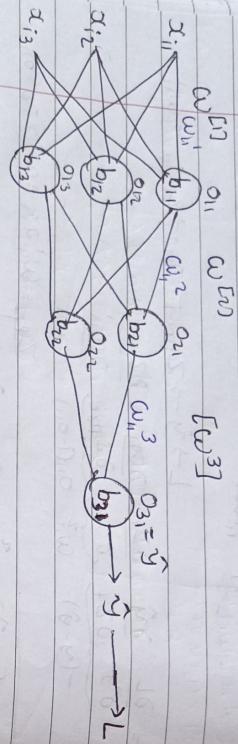
$$\frac{\partial L}{\partial w_{12}} = -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial b_1} = -(y - \hat{y}) w_{11}^2 o_{12} (1 - o_{12})$$

$$\frac{\partial L}{\partial b_2} = -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12})$$

multi
layer
perceptron

For 2 layers \Rightarrow MLP memorization



$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^2}$$

$$\text{For } \frac{\partial L}{\partial w_{11}^1} \\ L \rightarrow \hat{y} \rightarrow o_{11} \rightarrow w_{11}^1$$

$$L \rightarrow y \rightarrow o_{21} \rightarrow o_{11}$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial \hat{y}}{\partial y} \left[\frac{\partial y}{\partial o_{21}} \frac{\partial o_{21}}{\partial o_{11}} + \frac{\partial \hat{y}}{\partial o_{11}} \frac{\partial o_{11}}{\partial w_{11}^1} \right]$$

dark lines \rightarrow direct connection without
hidden layer (1 to 1 connections)

white lines \rightarrow hidden layer (1 to many connections)

$(0, 1) \text{ or } 0.5 \text{ or } 0.5 - 1.0 \rightarrow 1.0$

sigmoid function

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

≈ 0.5

How to improve a neural network

Hyperparameters

PAGE NO.: _____

PAGE NO.: _____

- 1) Fine tuning neural network
↳ Hyper parameters →
 - Batch size
 - Optimizer
 - Activation function

- 2) By solving problems:
 → Vanishing / Exploding gradient

- Not enough data
 → slow training

- overfitting

Fine tuning hyperparameter

- i) => more HL with few neurons is better than few HL with more neuron (⇒ Representation Learning)
⇒ extend HL until overfitting occurs

Neuron/layer

- ⇒ nodes should be sufficient so it capture primitive features

Batch size

- iii) Batch size
 → LR scheduling

Smaller

- (8 to 32)
↳ depends on GPU Ram

- ⇒ give generalize model
 → training fast

- ⇒ better result
 → LR varies small to high
(warming up LR)

Epoch value

- ⇒ Early stopping : is mechanism which is intelligent to understand when weight updation should be stop

- ii) Not enough data
 → we transfer learning
 → in transfer learning we use model which made by others in our model
 → unsupervised pretraining

- 2) i) vanishing/exploding gradient
 → weigh initialization
 → activation function
 → batch normalization
 → gradient clipping, for exploding gradient

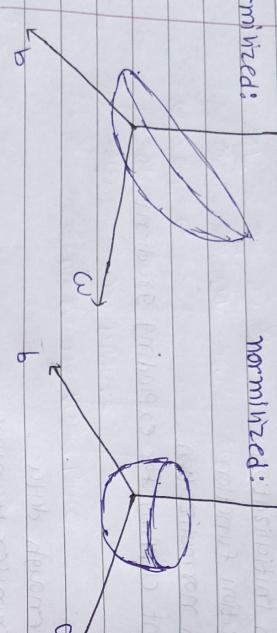
★ Feature Scaling

→ Standardization → J^* cost function
minimize J^* cost function

Unnormalized:



Normalized:



b) Data distribution

c) $\text{Weight distribution}$

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

case 1 zero initialization $W=0, b=0$

for new

$$C_{11} = \mu_{11}(0, 2_{11})$$

$$Z_{11} = W_1^T x + W_1^T x + b_{11} = 0$$

$$Z_{12} = W_1^T x + W_2^T x + b_{12} = 0$$

since $Z_{11} Z_{12} = 0$ so

$$C_{11} = C_{12} = 0$$

⇒ by adding penalty term we can reduce values of weight and avoid overfitting

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

more often regularization L become

$$L = L' + \frac{\lambda}{2} \sum_i \|W_i\|^2$$

$$\frac{\partial L}{\partial W_{\text{old}}} = \frac{\partial L'}{\partial W_{\text{old}}} + \frac{\lambda}{2} (\text{new})$$

+ new

my new

remain when we interested others zero

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

\downarrow

\downarrow

★ always to some overfitting

→ adding more rows \rightarrow stability Reducing complexity \rightarrow data augmentation, random noise, drop out, early stopping, regularization

- ↳ early stopping
- ↳ regularization

★ weight initialization

⇒ wrong weight initialization

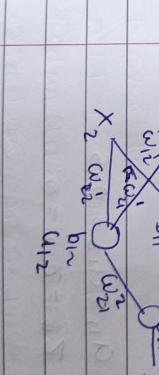
curse

gradient problem

→ vanishing GP

→ exploding GP

→ slow convergence



case 1 zero initialization $W=0, b=0$

$$Z_{11} = W_1^T x + W_1^T x + b_{11} = 0$$

$$Z_{12} = W_1^T x + W_2^T x + b_{12} = 0$$

Very high = under fitting

low = overfitting

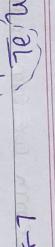
since $Z_{11} Z_{12} = 0$ so

$$C_{11} = C_{12} = 0$$

$$\left(\sum_{i=1}^K \|W_i\|^{1/n} = \sum_{L=1}^L \sum_{i=1}^n \|W_{ij}\|^{1/n} \right)$$

$$\begin{aligned} \text{Lasso} & \Rightarrow \text{Regularization hyper param} \\ L_2 & \rightarrow \text{Lasso} = L_2 + \frac{\lambda}{2n} \sum_{i=1}^n \|W_i\|^2 \\ \text{ridge} & \rightarrow \text{Lasso} = L_2 + \frac{\lambda}{2n} \sum_{i=1}^n \|W_i\|^2 \\ L_1 & \rightarrow \text{Lasso} = L_2 + \frac{\lambda}{2n} \sum_{i=1}^n \|W_i\| \end{aligned}$$

$$w_{ii}^1 = w_{ii} - \eta \frac{\partial L}{\partial w_{ii}^1}$$



$$L \rightarrow \hat{y} \rightarrow o_n \rightarrow w_{ii}^1$$

$$\text{For } \frac{\partial L}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{ii}^1} \frac{\partial a_{ii}^1}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial a_{ii}^1}{\partial w_{ii}^1}$$

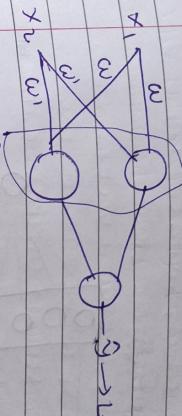
$$a_{ii} = a_{ii} = 0$$

$$\frac{\partial L}{\partial w_{ii}^1} = 0$$

no update

& no training take place

behave us 1 single neuron
(If there are 1000 neurons then
also that act as single)



for tanh

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

$$a_{ii} = a_{ii} - \overline{a_{ii}} \quad \text{since } \overline{a_{ii}} = 0$$

$$\frac{a_{ii} = 0}{a_{ii} + a_{ii} = 0}$$

sum as zero

for sigmoid

$$\sigma = \frac{1}{1+e^{-z}}$$

$$\boxed{a_{ii} = a_{ii} = 0.5}$$

$$\frac{a_{ii} = 0.5}{a_{ii} + a_{ii} = 1}$$

\Rightarrow act sum as sigmoid with zero initialize

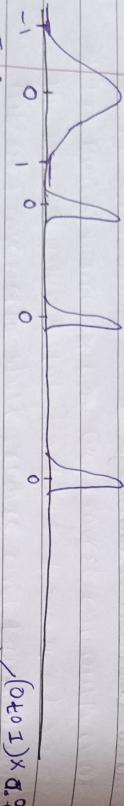
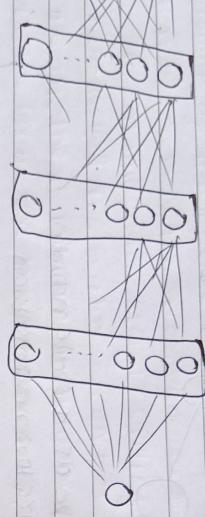
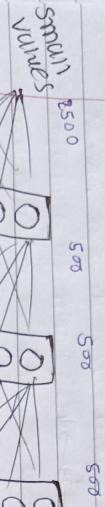
$$\frac{\partial L}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{ii}^1} \frac{\partial a_{ii}^1}{\partial z_{ii}^1} \frac{\partial z_{ii}^1}{\partial w_{ii}^1}$$

$$\frac{\partial L}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{ii}^1} \frac{\partial a_{ii}^1}{\partial z_{ii}^1} \frac{\partial z_{ii}^1}{\partial w_{ii}^1} \leftarrow x_1$$

$$\frac{\partial L}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{ii}^1} \frac{\partial a_{ii}^1}{\partial z_{ii}^1} \frac{\partial z_{ii}^1}{\partial w_{ii}^1} \leftarrow x_2$$

$$\frac{\partial L}{\partial w_{ii}^1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_{ii}^1} \frac{\partial a_{ii}^1}{\partial z_{ii}^1} \frac{\partial z_{ii}^1}{\partial w_{ii}^1} \leftarrow x_2$$

Case 3 Random initialization



Input gaussian dist
let bias = 0, np.random.rand(500, 500) * 0.01

and since weights are very small and
 $\sigma(\text{input}) \approx -1 < x_i < 1$ so

$$Z = \sum w_i x_i \text{ is small number}$$

and activation of all node very close to zero

For activation fun = tanh Z, sigmoid

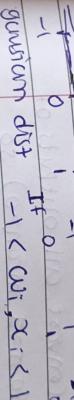
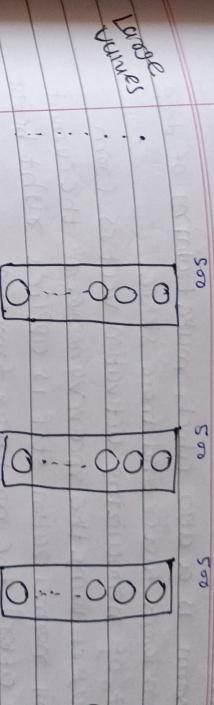
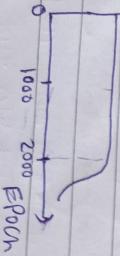
(for deep network), ReLU

activation is close to zero so during back propagation

gradient is very small so $w_{\text{new}} \approx w_{\text{old}}$ and vanishing gradient occurs

\Rightarrow In case of ReLU due to small gradient vanishing

not occur but it take more time and converge at very high epoch fun



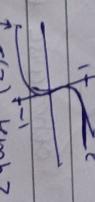
then

$$100 \leq \sum w_i x_i \leq 500$$

\Rightarrow due to high value of Z for activation fun like sigmoid

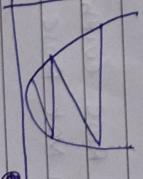
and tanh Z value converges so training will be slow

or vanishing gradient can be also occurs



very high and we get unstable

\Rightarrow For ReLU If Z is high then
activation function also high
due to this $w = \frac{\partial L}{\partial w}$



due to high activation function

$\frac{\partial L}{\partial w}$ high and we get this graph

\Rightarrow to avoid this kind of problem we uniform distribution,

Xavier / Gort. and we init with constant 0.01

MIN. J of loss 200 drop rate 0.001 or 0.01

Batch Normalization

PAGE NO. _____
PRESENTED BY _____

⇒ Batch Norm is applied with mini batch GD
⇒ applied layer by layer (can apply on any number of layers)

⇒ Batch Norm is used which make training of deep neural network faster and more stable.

⇒ It consists of normalizing activation vectors from hidden layers using the mean variance of the current batch. This normalization step is applied right before (or right after) the non-linear function.

⇒ Change in the distribution of network activations

due to the change in network parameters during training is called internal covariate shift.

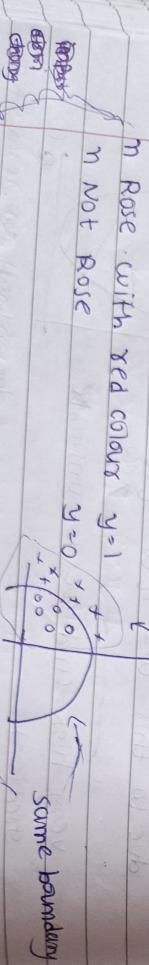
⇒ In Batch normalization we convert all outputs of HL into Gaussian distribution ($\sigma=1, \text{mean}=0$) so internal covariate shift reduce and training will be fast and stable.

⇒ If we don't use Batch Norm then LR should be low

→ initialization of parameters need to done carefully (like weights)

otherwise train is very unstable. so we use Batch Norm.

covariate shift



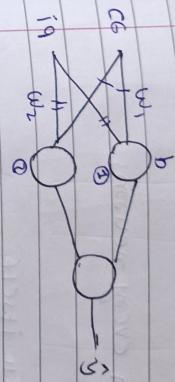
n Rose with red colour $y=1$
n Not Rose $y=0$

n Rose with diff colours $y=1$
n Not rose $y=0$

⇒ Input distribution and OLP

distribution is diff but OLP is same and input reduction

⇒ OLP sum (given by above C)
IS is called covariate shift and we need to train



$$z_{ii} = w_1 c_i + w_2 i + b$$

For batch norm 2 ways
 $z_{ii} \rightarrow z_{ii}^N \rightarrow g(z_{ii}^N) \rightarrow a_{ii}$

$$z_{ii} \rightarrow g(z_{ii}) \rightarrow a_{ii} \rightarrow a_{ii}^N$$

Let batch = n

$$(u_{ii}, z_{ii}) \cdot (z_{ii}, z_{ii}) = (u_{ii}, z_{ii})$$

↑
input weight
↓
bias
 $(u_{ii}, z_{ii}) + (1, z_{ii})$

We input n rows so get n activation on ① and n on ②

$$\text{For } ① \quad u_B = \frac{1}{m} \sum_{i=1}^m z_{ii} \quad \text{activation for input row}$$

For ② $u_B = \frac{1}{m} \sum_{i=1}^m z_{ii}$ These batch = n
so n input and n of one there for ②

$$\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_{ii} - u_B)^2} \quad \text{for every activation}$$

$$z_{ii}^N = \frac{z_{ii} - u_B}{\sigma_B}$$

$$\text{For } ② \quad u_B = \frac{1}{m} \sum_{i=1}^m z_{ii}$$

$$\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_{ii} - u_B)^2}$$

$\sigma_B + \epsilon \leftarrow \text{+ve small effect}$
to avoid zero's effect

$$y = \gamma z_{ii}^N + \beta \quad \text{initial value of } \gamma=1, \beta=0 \text{ in first row}$$

⇒ Here we done with $z_{ii} \rightarrow z_{ii}^N$

OLP sum (given by above C)
IS is called covariate shift and we need to train



denoted mean μ and standard deviation σ of input data
Normalizing the input data is called **batch normalization**

PAGE NO.:

during batch normalization every neuron has its own μ and σ
Parameter 2 are learnable γ_B and β are non learnable

PAGE NO.:

$$\text{every neuron has its own } \gamma_B + \beta$$

$$z_{ii} \xrightarrow[\substack{\mu_i=0 \\ \sigma_i=1}]{} z_{ii}^N \xrightarrow{} z_{ii}^{BN} \xrightarrow{} g(z_{ii}^{BN}) \rightarrow a_{ii}$$

\Rightarrow This process is done with every neurons

$$\text{If } r = \sigma + \epsilon, \beta = \alpha$$

$$\text{then } z_{ii}^N \rightarrow z_{ii}^{BN} = z_{ii} \text{ itself}$$

we give flexibility by giving $r \in B$. If neural network has no need of Normalization then r, β are set accordingly

\rightarrow Advantage: \rightarrow Stable so we can choose hyper para in large range
 \rightarrow training faster (we can set high LR)

\rightarrow has regularization effect



to calculate for every batch which is dependent on batch itself and changing of batch value also change activation also change and due to change on activation

system get some randomness and overfitting reduce slightly.

It has very much effect so we need to use regularization also.

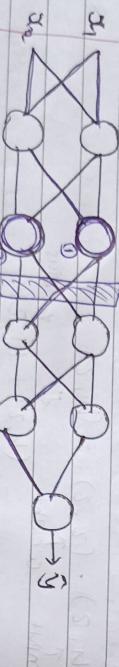
\rightarrow Reduce weight initialization impact

slow converge



weights are slow because cost fun is uniform

Cost fun without normalization



\hookrightarrow $\text{parameters } r, \beta$ too $\neq 0 \neq 0$

during back propagation r and β also get updated by loss

$$r = r - \eta \frac{\partial L}{\partial r}, \beta = \beta - \eta \frac{\partial L}{\partial \beta}$$

\Rightarrow Here we can decide batch for training

or testing where we have only one input. In this type for case we use exponentially weighted Avg (EWMA)

\Rightarrow for test data we maintain exponentially moving avg using $\sigma_1, \sigma_2, \sigma_3, \dots, \sigma_n$ for n batch

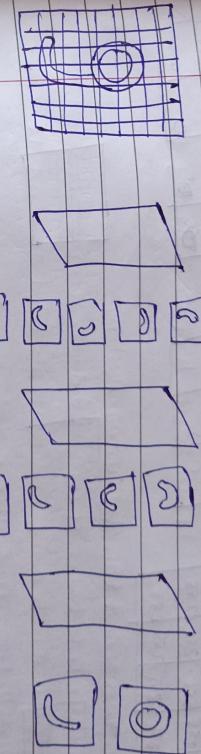
and we final use σ for further calculation

CNN (convolutional) Neural Network

CNN = matrix multiplication

→ ANN is not works in image data because of
i) high computational cost

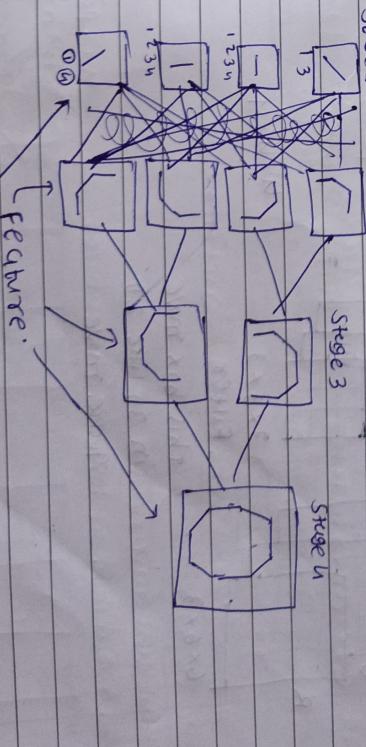
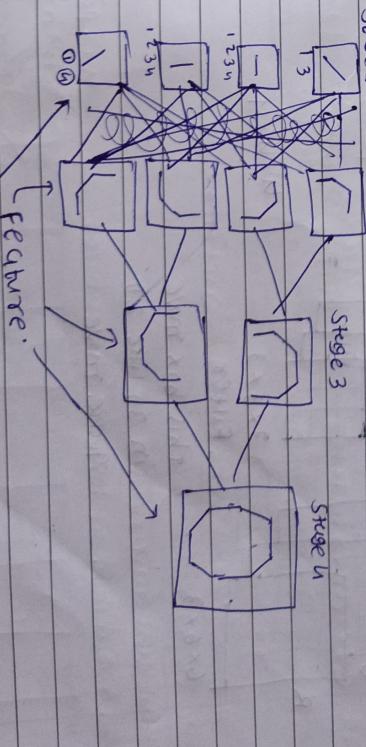
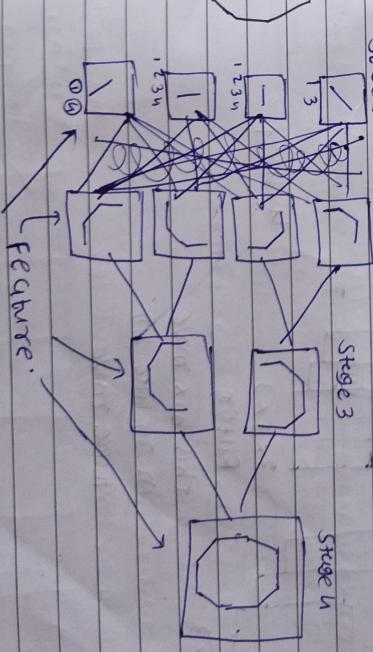
ii) overfitting
iii) loss of info into like spatial arrangement of pixels



Complex Feature

⇒ CNN extract features from given data as we go ahead
It detect complex features

⇒ There are two types of cell
Simple cell: detect edges
Complex cell: detect complex features from simple cell



$$0(-1) + 0(0)(-1) + 0(0)(-1) + (0)(0)(0) + (0)(0)(0) + (0)(0)(0) + (0)(0)(0)$$

PAGE NO.: _____

PAGE NO.: _____

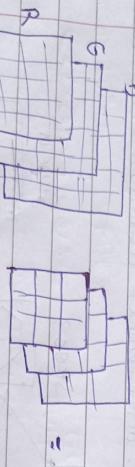
Convolution operation (edge detection)

horizontal filter / kernel		
3x3		
0	0	0
0	0	0
0	0	0
-1	-1	-1
0	0	0
256	256	256
256	256	256
256	256	256

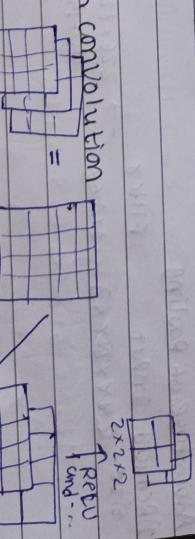


- ⇒ In CNN we use filters like weights in CNN
- ⇒ Values inside filter is compute from back propagation

working with RGB Images



Problem with convolution



$$\left[\frac{m+f-1}{s} + 1 \right] \times \left[\frac{n+f-1}{s} + 1 \right] \quad (\text{if padding } P = 0)$$

$$\left[\frac{m+2p-f}{s} + 1 \right] \times \left[\frac{n+2p-f}{s} + 1 \right] \quad (\text{if padding } P > 0)$$

- ⇒ Padding and strides
- ⇒ In Padding we added boundary with value zero
- ⇒ Strides decide that how filter is moving (one column should be skip)
- ⇒ General values in padding is zero.

Image $m \times n \Rightarrow (m-f+1) \times (n-f+1)$

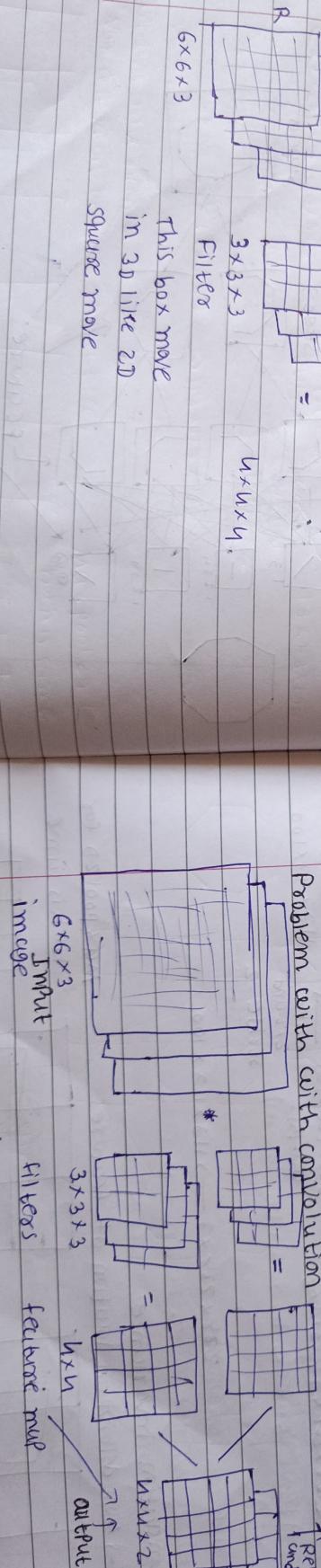
Filter $f \times f \times 3 \Rightarrow \left[\frac{m-f}{s} + 1 \right] \times \left[\frac{n-f}{s} + 1 \right] \times 3$ (if stride is s)

(L · J floor operation)

$$\left[\frac{m+2p-f}{s} + 1 \right] \times \left[\frac{n+2p-f}{s} + 1 \right] \times 3 \quad (\text{if padding } P & \text{ Strides } s)$$

Why strides required

- If we need high level feature
- Computing time low



- ⇒ We need to reduce feature map otherwise size of output is very big. This can be done by increasing stride value

→ There is majority 2 problems of memory issue

② Translation Variance

Translation Variance

→ After applying convolution layer feature will be the by convolution variance and is called translation variance.

Translation variance → Translation Pooling

→ Translation variance is more sensitive to image rotation and translation than it is to noise.

W) No Need of training

→ In Keras there are 3 types of pooling

1) Max Pooling

2) Avg Pooling ← avg of numbers

3) Global Pooling ← max of all numbers

global
max
avg

disadvantage of Pooling

→ can't be used in image segmentation task

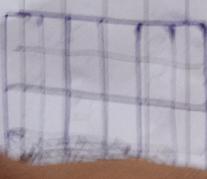
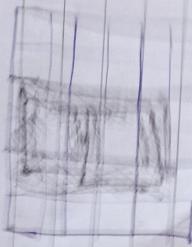
→ loss lose information $5 \times 5 \rightarrow 2 \times 2$

lose 75% information

$$228 \times 228 \times 3 \rightarrow 112 \times$$

Translation Variance

→ location of features



PAGE NO.: _____

PAGE NO.: _____

→ There are majorly 2 problems of memory issue
 1) Translation variance

2) Translation Variance

3) using (max Pooling) Pooling - Enhanced Features
 (only for max Pooling)

Translation variance

⇒ After applying convolution layer feature which is that by CNN is dependent on location this is called translation variance

⇒ we can avoid this by using Pooling

Pooling is actually a way to down sample feature map so features become independent of location and it called

Translation invariance

Ex: $\begin{matrix} 3 & 1 & 1 & 3 \\ 1 & 5 & 0 & 2 \\ 1 & n & 2 & 1 \\ 1 & 2 & n \end{matrix}$ here low level details are eliminated

After PELU $\Rightarrow \begin{matrix} 5 & 3 \\ 7 & n \end{matrix}$

Type: max pooling ⇒ In this process first we

size: 2×2

stride = 2 apply PELU on output feature map

Advantage of Pooling

1) Reduced Input Filter $\rightarrow 226 \times 226 \times 100$ Pooling with $113 \times 113 \times 100$

Size $228 \times 228 \times 3$ $3 \times 3 \times 3$ Output with $113 \times 113 \times 100$

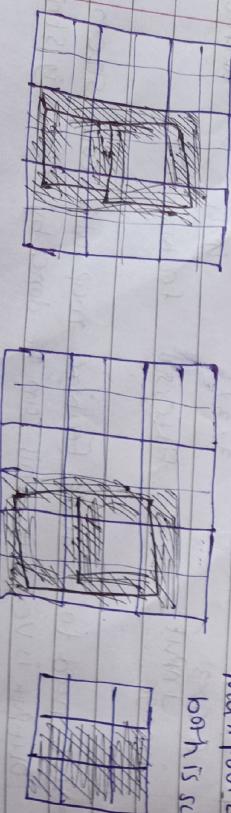
Stride = 2

$228 \times 228 \times 3 \Rightarrow 113 \times 113 \times 100$

2) Translation invariance

⇒ Not depend on location of features

Max pool 2D for both is sume

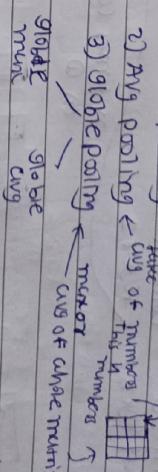


Disadvantage of Pooling

⇒ can't be used in image segmentation task

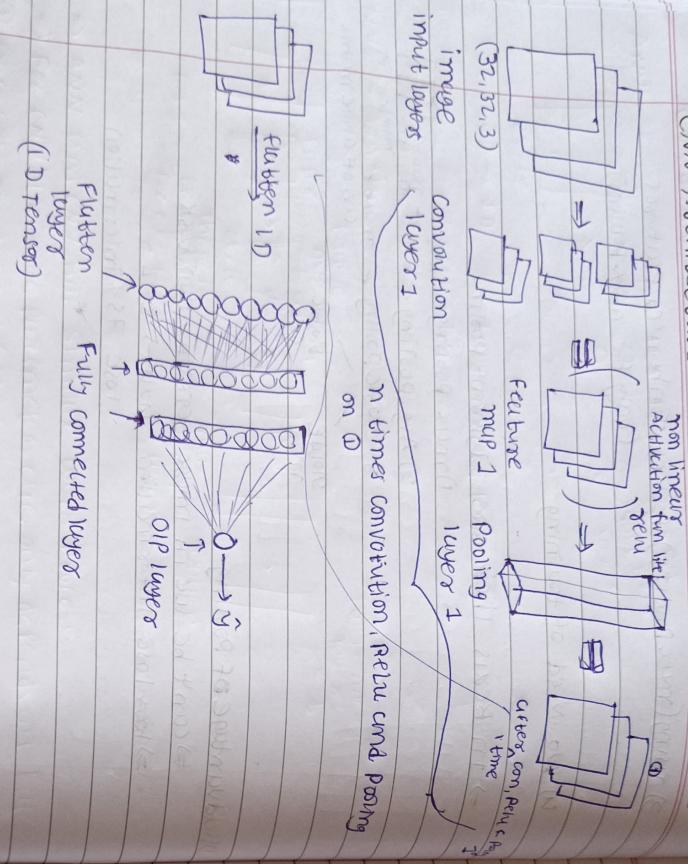
⇒ loses information $n \times n \rightarrow 2 \times 2$

lose 75% information



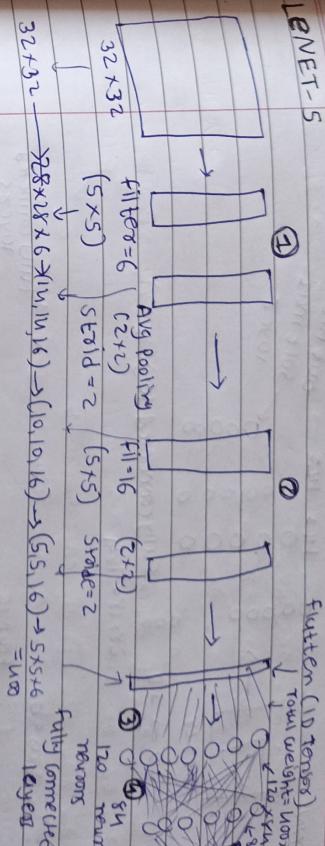
CNN Architecture

PAGE NO.: _____



PAGE NO.: 089
TOPIC: CNN

LENET-5

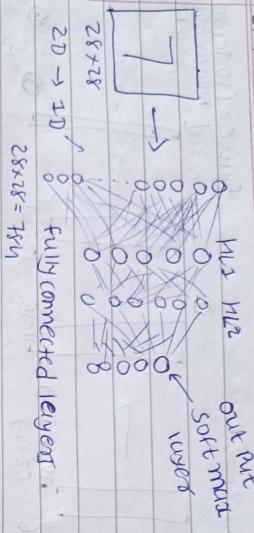


CNN Architecture

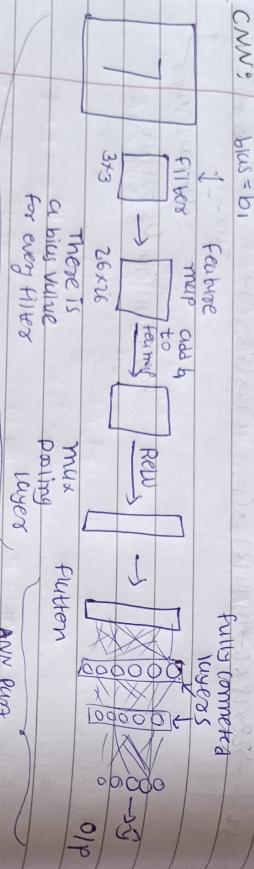
- 1) LENET-5 made by Yann Le Cun
- 2) Alex NET
- 3) GoogleLe NET
- 4) VGGNET
- 5) RESNET
- 6) Inception

CNN & ANN

ANN:



ANN:



\Rightarrow

ANN and CNN has similar work process.

size and number of filters

For Ex (3, 3, 3)

Learnable parameters is every filter have own bias value

Input (1080, 1080, 3)

$$= (3 \times 3 \times 3) \times 50 + 50 = 270(50) + 50$$

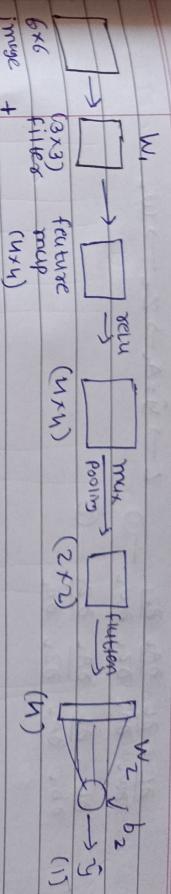
$$= 1350 + 50 = 1400$$

\Rightarrow while ANN learnable parameter is depends on size of input learnable parameters are

$$= 1080 \times 1080 \times 3$$

Backpropagation in CNN

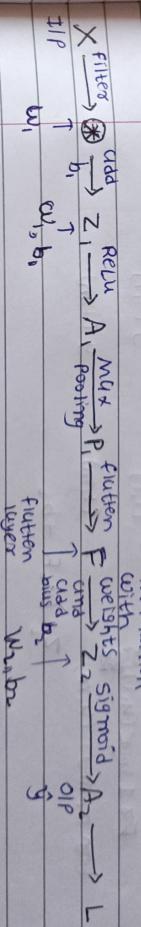
W1



Learnable parameters
 $W_1 = (3, 3) \quad W_2 = (1, n) = 10 + 5$
 $b_1 = (1, 1) \quad b_2 = (1, 1) = 15$ trainable parameters

\Rightarrow After forward propagation we will calculate loss

Logical flow



Forward propagation

$$Z_1 = \text{conv}(X, W_1) + b_1$$

$$A_1 = \text{ReLU}(Z_1)$$

$$P_1 = \text{max_pool}(A_1)$$

$$F = \text{fuction}(P_1)$$

$$Z_2 = F \cdot W_2 + b_2$$

$$A_2 = \sigma(Z_2)$$

Gradient descent

$$W_1 = W_1 - \eta \frac{\partial L}{\partial W_1}$$

$$b_1 = b_1 - \eta \frac{\partial L}{\partial b_1}$$

$$W_2 = W_2 - \eta \frac{\partial L}{\partial W_2}$$

$$b_2 = b_2 - \eta \frac{\partial L}{\partial b_2}$$

W₁ in training phase → W₁ → Z₁ → A₁ → P → F → 0 → Z₂ → A₂ → L

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad \text{--- (1)}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial b_1} \quad \text{--- (2)}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial F} \frac{\partial F}{\partial R} \frac{\partial R}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \quad \text{--- (3)}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial F} \frac{\partial F}{\partial R} \frac{\partial R}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial F} \frac{\partial F}{\partial R} \frac{\partial R}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

Flatten layers
ANN part

$$\text{For 1 image } (u_{1,1}) \rightarrow A_{1,1} \rightarrow L \rightarrow y_{(1,1)}$$

for m images (u_{m,1}) → (1,m) → y (1,m)

$$z_2 = w_2 F + b_2$$

$$A_2 = \sigma(z_2)$$

$$\text{For } \frac{\partial L}{\partial u_1} = \frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial z_2} \frac{\partial z_2}{\partial w_1}$$

for single image

$$\frac{\partial L}{\partial u_1} = \frac{w_2 - y_i}{\sigma(u_1)} \cdot u_1$$

for single instance
we take u₁

$$dA_2 = \sigma(z_2) \cdot [1 - \sigma(z_2)] = u_2(1-u_2)$$

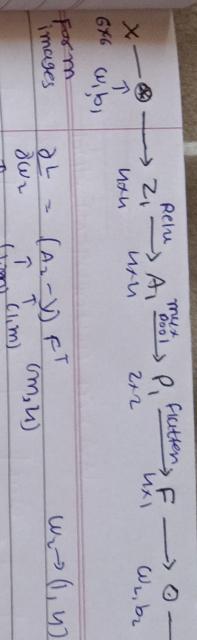
$$\frac{\partial z_2}{\partial u_2} = F$$

$$\frac{\partial z_2}{\partial w_2} = 1 \quad \frac{\partial L}{\partial w_2} = (u_2 - y_i) F$$

replace u₂ by An matrix and y_i to y matrix

$$\frac{\partial L}{\partial w_2} = (A_2 - y) F^T$$

$$\frac{\partial L}{\partial b_2} = (A_2 - y) F^T$$



$$\text{for (3) \& (4) } z_1 = \text{conv}(x, w) + b_1$$

$$A_1 = \text{ReLU}(z_1)$$

$$F = \text{Flatten}(P)$$

$$A_2 = \sigma(z_2)$$

$$L = \frac{1}{m} \sum_i (-y_i \log(A_2) - (1-y_i) \log(1-A_2))$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial F} \frac{\partial F}{\partial R} \frac{\partial R}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial F} \frac{\partial F}{\partial R} \frac{\partial R}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

gradient need to calculate
(A₂-y)

$$\frac{\partial z_1}{\partial F} = w_1, \quad \frac{\partial F}{\partial P}, \quad \text{no trainable parameter}$$

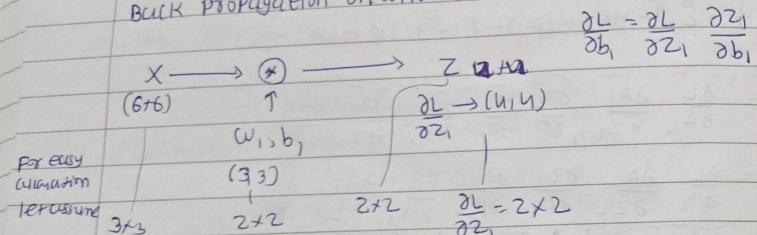
reshape(P, shape)

$$\frac{\partial L}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial F} = (A_2 - y) w_1 \cdot \text{reshape}(P, \text{shape})$$

$$\frac{\partial A}{\partial} \frac{\partial L}{\partial A_i} = \begin{cases} \frac{\partial L}{\partial p_{i,y}} & \text{if } A_{mi} \text{ is the max element} \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial A}{\partial z_1} = \begin{cases} 1, & \text{if } z_{1,y} > 0 \\ 0, & \text{if } z_{1,y} \leq 0 \end{cases}$$

BACK PROPAGATION ON CONVOLUTION:



$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \quad w_1 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \quad z_{xi} = z_1 = \begin{bmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \end{bmatrix}$$

$X \otimes w_1 + b$

Convolving $X \otimes w_1$ can be done by adding bias (we get)

$$z_{11} = x_{11}w_{11} + x_{12}w_{12} + x_{13}w_{21} + x_{21}w_{22} + b_1$$

$$z_{12} = x_{11}w_{11} + x_{13}w_{12} + w_{21}w_{21} + x_{22}w_{22} + b_1$$

$$z_{21} = x_{21}w_{11} + x_{22}w_{12} + x_{31}w_{21} + x_{32}w_{22} + b_1$$

$$z_{22} = x_{22}w_{11} + x_{23}w_{12} + x_{31}w_{21} + x_{33}w_{22} + b_1$$

$$\frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

~~z1~~

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_{11}} \frac{\partial z_{11}}{\partial b_1} + \frac{\partial L}{\partial z_{12}} \frac{\partial z_{12}}{\partial b_1} + \frac{\partial L}{\partial z_{21}} \frac{\partial z_{21}}{\partial b_1} + \frac{\partial L}{\partial z_{22}} \frac{\partial z_{22}}{\partial b_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_{11}} + \frac{\partial L}{\partial z_{12}} + \frac{\partial L}{\partial z_{21}} + \frac{\partial L}{\partial z_{22}} = \text{sum} \left(\frac{\partial L}{\partial z_i} \right)$$

$$\frac{\partial L}{\partial b_1} = \text{sum} \left(\frac{\partial L}{\partial z_i} \right) \quad \leftarrow \text{scalar}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_{11}} \frac{\partial z_{11}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z_{11}} \frac{\partial z_{11}}{\partial w_{11}} + \frac{\partial L}{\partial z_{12}} \frac{\partial z_{12}}{\partial w_{11}} + \frac{\partial L}{\partial z_{21}} \frac{\partial z_{21}}{\partial w_{11}} + \frac{\partial L}{\partial z_{22}} \frac{\partial z_{22}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{12}} = w_{11} \rightarrow w_{21}$$

replace.

$$\text{sum for } w_{11} \quad \frac{\partial L}{\partial w_{ij}} \quad i=1,2 \quad j=1,2$$

From eqns of $\frac{\partial z_{xy}}{\partial w_{ij}}$ we can find $\frac{\partial z_{xy}}{\partial w_{ij}} \quad x=1,2, y=1,2$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z_{11}} x_{11} + \frac{\partial L}{\partial z_{12}} x_{12} + \frac{\partial L}{\partial z_{21}} x_{21} + \frac{\partial L}{\partial z_{22}} x_{22}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial z_{11}} x_{12} + \frac{\partial L}{\partial z_{12}} x_{13} + \frac{\partial L}{\partial z_{21}} x_{22} + \frac{\partial L}{\partial z_{22}} x_{23}$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial z_{11}} x_{21} + \frac{\partial L}{\partial z_{12}} x_{22} + \frac{\partial L}{\partial z_{21}} x_{31} + \frac{\partial L}{\partial z_{22}} x_{32}$$

$$\frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial z_{11}} x_{22} + \frac{\partial L}{\partial z_{12}} x_{23} + \frac{\partial L}{\partial z_{21}} x_{32} + \frac{\partial L}{\partial z_{22}} x_{33}$$



we can simplify this into convolution operation

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \quad \frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

$$\frac{\partial L}{\partial z_1} = \text{conv}(x, \frac{\partial L}{\partial z_1})$$

$$\frac{\partial L}{\partial z_1} = \text{conv}(x, \frac{\partial L}{\partial z_1})$$

?

$$\frac{\partial L}{\partial z_1} = \text{sum}(\frac{\partial L}{\partial z_1})$$

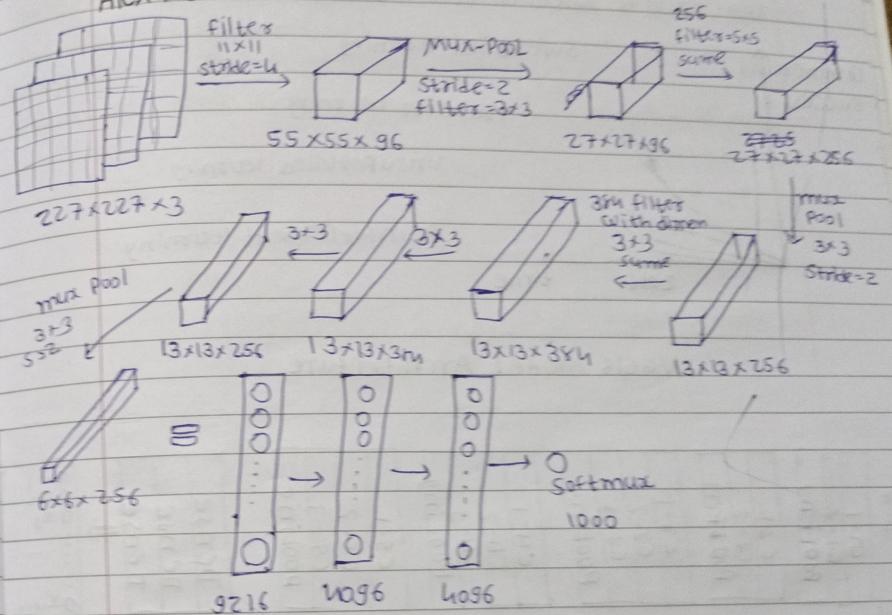
pretrained models in CNN

⇒ why

→ Data hungry

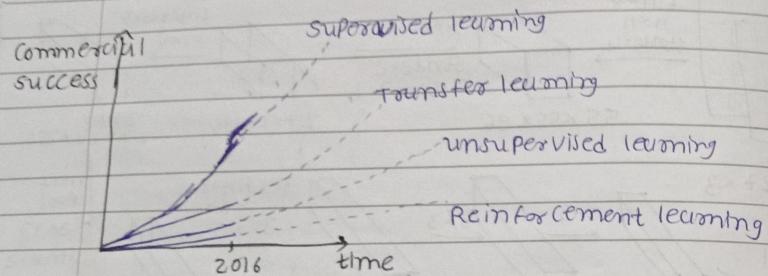
→ Time

Alex Net

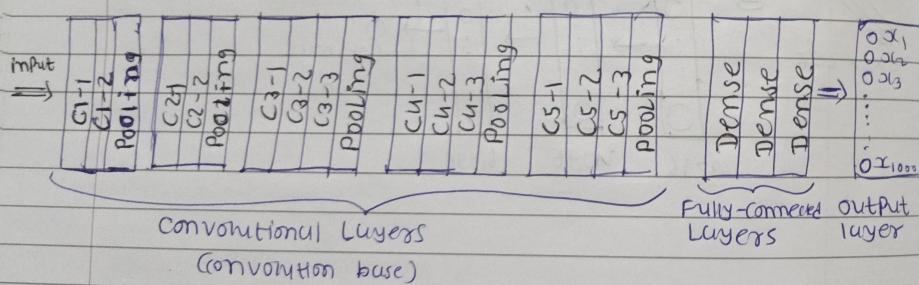


Transfer Learning

⇒ Transfer Learning is a research problem in machine learning that focuses on applying knowledge gained while solving one task to a related task.
For ex. Knowledge gained while learning to recognize car could be used in recognition of truck.



VGG16 Model Architecture



- ⇒ Convolution Layers find relation between pixel in 2D matrix and fully connected part do classification
- ⇒ In this we freeze convolution base and only train fully connected layers a

ways of transfer learning

Feature extraction
⇒ In this Dense layers is replace by our dense layers and training of convolutional layers freeze training occurs only on our dense layers

Fine tuning
⇒ In fine tuning we also train some last convolution layers

RNN (Recurrent neural network)

⇒ RNN is type of sequential model to work on sequential data

Ex prediction of animals: sequence of input like shape of nose, eyes, ears etc.

sequential data: text (hi, my name is vatsal), time series data, speech (~m/m~)

why use RNN?

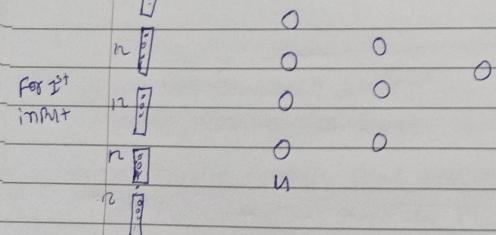
Input	Output
1 Hi my name is vatsal	0
2 I love ramux	0
3 Indicating the match	1

12 unique words \downarrow size 12
 hi → [0 1 0 0 ... 0]
 my → [0 0 1 0 0 ... 0]
 name → [0 0 0 1 0 0 ... 0]

vectorization

ANN Text classification

12 inputs $W = 2 \times 12$ weights



60

⇒ Here input size is varies but for ANN there should be const input size

For 2nd input 36×4

3rd 48×4

So for ANN we use zero padding but it is not efficient

Forward Propagation in RNNs

Review

x_1^{in} movie was good

Sentiment

1

x_2^{in} movie was bad

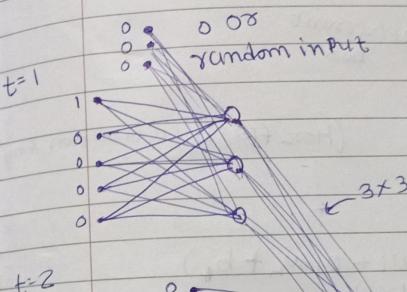
0

x_3^{in} movie was not good

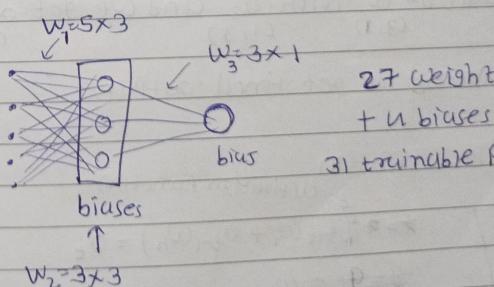
0

[1 0 0 0 0] [0 1 0 0 0] [0 0 1 0 0] [0 0 0 1 0]
 movie was good bad

[0 0 0 1]
 not



t=2



Scanned with OKEN Scanner

review

x_{11}	x_{12}	x_{13}
x_{21}	x_{22}	x_{23}
x_{31}	x_{32}	x_{33}
x_{34}		

sentiment

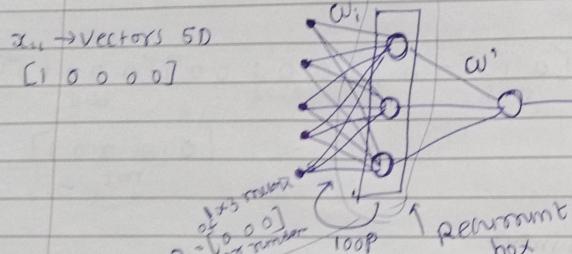
1

0

0

$x_{11} \rightarrow$ Vectors 5D

[1 0 0 0 0]



$$\begin{aligned}
 t=1 & \quad x_{11} \xrightarrow{w_i} o_1 = f(x_{11}w_i + b_1) \quad (\text{Here } f(x) \text{ is activation fun}) \\
 t=2 & \quad x_{12} \xrightarrow{w_i} o_2 = f(x_{12}w_i + o_1w_h + b_2) \\
 t=3 & \quad x_{13} \xrightarrow{w_i} o_3 = f(x_{13}w_i + o_2w_h + b_3)
 \end{aligned}$$

Now, dot Product of w^T with o_3 and we get output

(3, 1)

(1, 3)

(1, 1)

and after sigmoid we get final output

⇒ Simplified Representation

$$\begin{aligned}
 & \text{activation function: } g \\
 & \star = \star f(x_{it}w_i + o_{t-1}w_h) = o_t \\
 & g = g(o_{t-1}w_h) \\
 & \text{If } o_t \text{ is last OP} \\
 & \text{then we pass it through activation fun} \\
 & g = g(o_{t-1}w_h)
 \end{aligned}$$

i = row

t = time step

Types of RNN

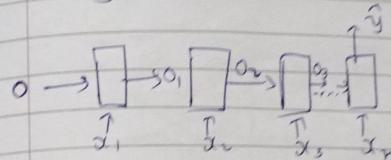
i) many to one

input → sequence

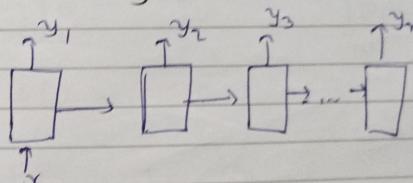
output → non sequence

sentiment analysis

Rating Prediction



ii) one to many



Ex describe image,
music generation

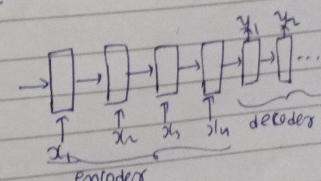
iii) many to many

Input → sequence

Output → sequence

Ex pos tagging
sum length many to many
variable length " " "
machin translation

Eng to Hindi



Scanned with OKEN Scanner

iv) one to one (ANN, NN)

~~or~~ Input = non sequential

Output = " "

Eg: image classification

This is simple neural network



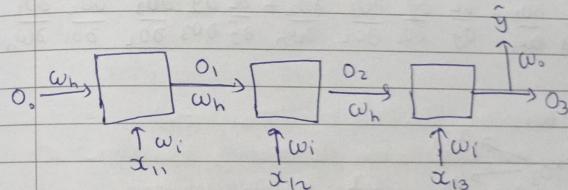
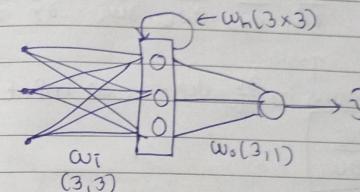
Back Propagation in RNN

Text

cat mat rat $x_1[1 \ 0 \ 0] [0 \ 1 \ 0] [0 \ 0 \ 1]$

rat cat mat $x_2[0 \ 0 \ 1] [0 \ 0 \ 0] [0 \ 1 \ 0]$

mat rat cat $x_3[0 \ 1 \ 0] [0 \ 1 \ 0] [1 \ 0 \ 0]$



$$O_1 = f(x_{11}w_i + O_0w_h)$$

$$O_2 = f(x_{12}w_i + O_1w_h)$$

$$O_3 = f(x_{13}w_i + O_2w_h)$$

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

$$L = -y_i \log \hat{y}_i - (1-y_i) \log(1-\hat{y}_i)$$

minimiz loss using gradient descent

$$w_i, w_h, w_o \quad w_i = w_i - \eta \frac{\partial L}{\partial w_i}$$

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

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



$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_i}$$

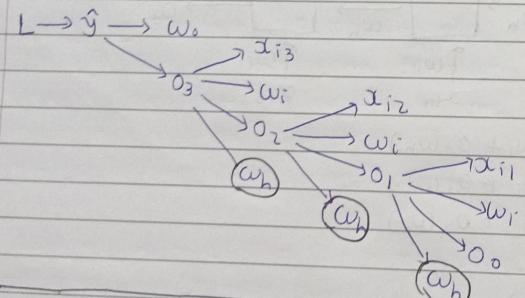


$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_i} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_i}$$

Total terms in $\frac{\partial L}{\partial w_i}$ depends on input

$$\frac{\partial L}{\partial w_i} = \sum_{j=1}^{t=3} \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_j} \frac{\partial o_j}{\partial w_i}$$

$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial w_h} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_h} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_h}$$



$$\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_j} \frac{\partial o_j}{\partial w_h}$$

problems with RNN

- ⇒ i) Problem of long term dependency (vanishing gradient problem)
- ii) unstable gradients (exploding gradient)

i) Problem of long Term dependency (Vanishing)

input out

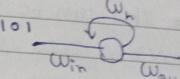
1 0 1 1

0 0 1 0

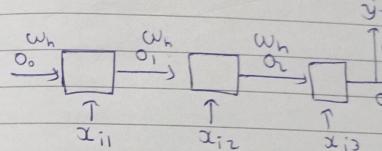
0 0 0 0

1 1 1 1

input hidden



3 time step



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

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial w_{in}} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_{in}} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}}$$

↑
short term dependency

↑
long term dependency

$$\begin{aligned} \text{Long Term dependency} & \quad \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{100}} \frac{\partial o_{100}}{\partial o_{99}} \dots \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}} \\ \text{look like} & = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{100}} \prod_{t=2}^{100} \left(\frac{\partial o_t}{\partial o_{t-1}} \right) \frac{\partial o_1}{\partial w_{in}} \end{aligned}$$

$$o_1 = \tanh(\alpha_{i1} w_{in} + o_0 w_h)$$

$$o_2 = \tanh(\alpha_{i2} w_{in} + o_1 w_h)$$

$$o_t = \tanh(\alpha_{it} w_{in} + o_{t-1} w_h)$$

$$\frac{\partial o_t}{\partial o_{t-1}} = \tanh'(\alpha_{it} w_{in} + o_{t-1} w_h) w_h$$

$$\frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{in}} \left[\sum_{t=2}^{100} [tanh'(a_t(w_{in} + w_t \cdot h_t)) \cdot h_t] \right] \frac{\partial a_t}{\partial w_{in}}$$

0 to 1

so above term ≈ 0

$\frac{\partial L}{\partial w_{in}}$ only short term contribution is consider.

- solution
 \Rightarrow This can be solved by using different activation function ReLU, leaky ReLU
 \rightarrow better weight initialization
 \rightarrow skip RNNs
 \rightarrow LSTM

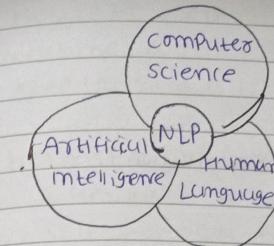
ii) unstable Training (Exploding gradients)

- \rightarrow In this problem long term dependency become very large so gradient and weights become infinite and model does not train.
- \Rightarrow If LR is large then same problem occurs

- solution:
 \rightarrow gradient clipping (Set max value of gradient)
 \rightarrow controlled Learning Rate (LR)
 \rightarrow LSTM

Natural Language Processing (NLP)

Introduction to NLP



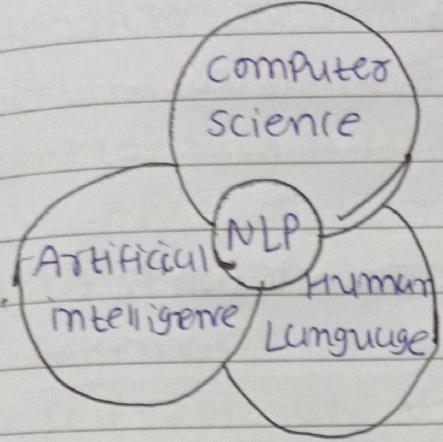
\Rightarrow In neuropsychology, linguistics and the philosophy of language, a natural language or ordinary language is any language that has evolved naturally in humans through use and repetition without conscious planning or premeditation. Natural languages can take different forms, such as speech or signing. They are distinguished from constructed and formal languages such as those used to program computers or to study logic.

Real world applications

- \Rightarrow Contextual Advertisements
- \Rightarrow Email clients - Spam filtering, smart reply
- \Rightarrow Social media - removing adult content, opinion mining
- \Rightarrow search engines
- \Rightarrow chatbots



Introduction to NLP



⇒ In neuropsychology, linguistics and the philosophy of language. ~~u~~ natural language or ordinary language is any language that has evolved naturally in humans through use and repetition without conscious planning or premeditation. Natural languages can take different forms, such as speech or signing. They are distinguished from constructed and formal languages such as those used to program computers or to study logic.

Real world applications

- ⇒ Contextual Advertisements
- ⇒ Email clients - spam filtering, smart reply
- ⇒ Social media - removing adult content, opinion mining
- ⇒ search engines
- ⇒ chatbots

common NLP Tasks

- ⇒ Text / Document classification
- ⇒ sentiment Analysis
- ⇒ informative Retrieval
- ⇒ Part of speech Tagging
- ⇒ Language detection & Machine Translation
- ⇒ conversational Agents
- ⇒ Knowledge graph & QA systems
- ⇒ Text summarization
- ⇒ Topic modeling
- ⇒ Text generation
- ⇒ spell checking and grammar correction
- ⇒ Text parsing
- ⇒ speech to text

Approaches to NLP

- i) Heuristic methods
- ii) machine Learning Based methods
- iii) Deep learning Based methods

i) Heuristic approaches (GOALS)

A Heuristic is any approach to problem solving or self-discovery that employs a practical method that is not guaranteed to be optimal perfect or rational but is nevertheless sufficient for reaching an immediate short-term goal or approximation.

Examples

- ⇒ Regular Expressions
- ⇒ curated
- ⇒ open mind common sense

ii) Machine learning approaches

⇒ advantage of this approaches is when some time its not possible to make rules by human so in this kind of scenario ML approaches used because in ML rules are made by machine itself.

⇒ Algorithms used

- Naive Bayes
- Logistic regression
- SVM
- LDA
- Hidden markov models



iii) Deep learning approaches

→ ML approaches loss sequential information so
DL approaches is very useful

→ Architectures used

RNN

LSTM (long short term memory)

CRF/CNN

Transformers

Autoencoders

* Challenges in NLP

- Ambiguity (I saw the boy on the beach with my binoculars)
- Contextual words (I ran to the store because we ran out of milk)
- Colloquialisms and slang (piece of cake)
- Synonyms
- Long, sarcasm and tonal difference
- spelling Errors
- Creativity (poems, dialogue, scripts)
- Diversity

What is NLP Pipeline

⇒ NLP pipeline is a set of steps followed to build an end-to-end NLP software.

⇒ NLP software consists of the following steps:

- Data Acquisition
- Text Preparation
 - Text cleanup
 - Basic Preprocessing
 - Advance Preprocessing

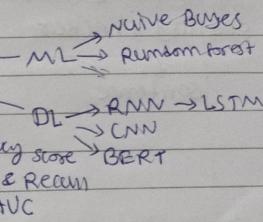
→ Feature engineering

→ modelling

- model building
- Evaluation

→ Deployment

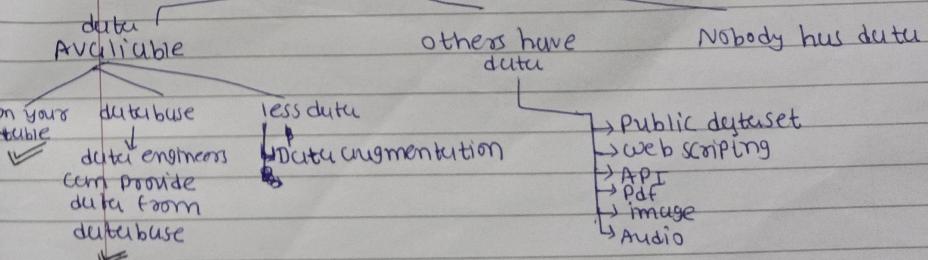
- Deployment
- monitoring
- Model update



Points to remember

- It's not universal
- Deep learning pipelines use slightly different
- Pipeline is non-linear

i) Data Acquisition



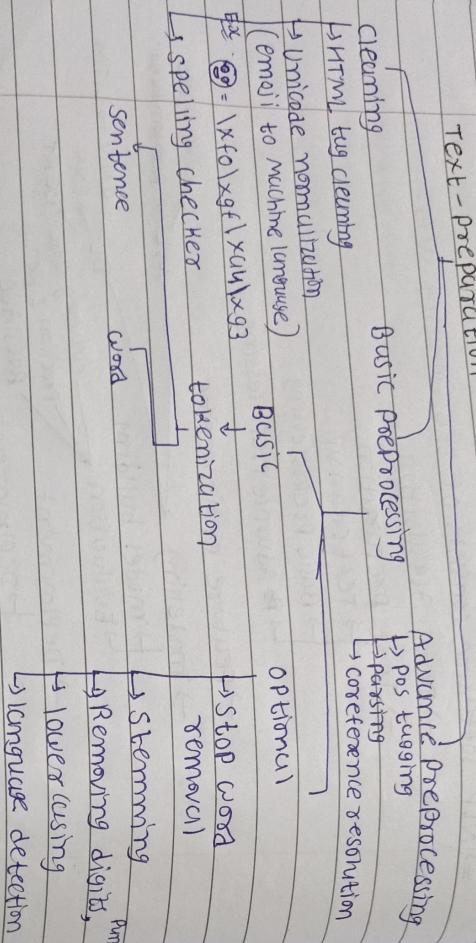
iii) Feature Engineering

→ Feature convert to numbers

Eg. review text sentiment

ii) text-preprocessing

Text - Preparation



	positive words	negative words	neutral	sentiment
5	2	6	0	
3	1	6	1	
1				



⇒ Feature engineering by

ML algo

- base on domain knowledge need to do feature engi.
- and create feature
- Advantage: can understand how result come
- Disadv: need to have knowledge about domain

ML algo

- ML algo do feature engi it self
- Adv: without domain knowledge can make model
- disadv: can not know how actually work

iv) Modelling

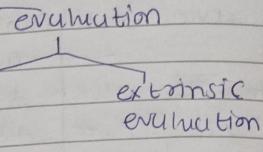
- modeling
- ↳ Heuristic
- ↳ ML Algo
- ↳ DL
- ↳ Cloud API

which approach should be used
is depend on 2 things
↳ Amount of data
↳ Nature of Problem

v) Deployment

- Deployment
- ↳ API (microservice)
- ↳ Chatbot

monitoring update



Feature Engineering

Common terms

corpus: combination of all words of data set which also contain repeated word

vocabulary: contain unique words of dataset

Document: each review is called Document

word: word

* One Hot Encoding

In this ~~Technique~~ Technique Text is convert into vector of dimension $V = \text{vocabulary}$

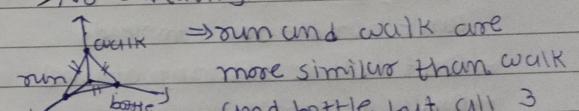
Problems with OHE → sparsity (contain large number of 0)

→ no fixed size

→ OOV (out of vocabulary)

If input is not in vocabulary

→ No capturing of semantic



⇒ so OHE cannot capture actual meaning of the sentence

* Bag of words

$V=5$

- | | people | watch | campus | write | comment |
|----|--------|-------|--------|-------|---------|
| D1 | 1 | 1 | 1 | 0 | 0 |
| D2 | 0 | 1 | 2 | 0 | 0 |
| D3 | 1 | 0 | 0 | 1 | 1 |
| D4 | 0 | 0 | 1 | 1 | 1 |

	people	watch	campus	write	comment
D1	1	1	1	0	0
D2	0	1	2	0	0
D3	1	0	0	1	1
D4	0	0	1	1	1

Adv: able to capture semantic of the sentence
easy to implement

Disadv: dimension of vocabulary is huge so it slow the algo.
out of vocabulary

\Rightarrow Bag of words note that how many time specific word
use occurs not note order of words

n dimensional space

people

→ single between two vector
denote similarity between them

→ Shallow → feature extraction

Count Vectorizer

D_1

D_2

D_3

D_4

D_5

D_6

D_7

D_8

D_9

D_{10}

D_{11}

D_{12}

D_{13}

D_{14}

D_{15}

D_{16}

D_{17}

D_{18}

D_{19}

D_{20}

D_{21}

D_{22}

D_{23}

D_{24}

D_{25}

D_{26}

D_{27}

D_{28}

D_{29}

D_{30}

D_{31}

D_{32}

D_{33}

- \Rightarrow If input contain words which not in vocabulary then Bow ignore those words
- Advantage: simple and intuitive fix size can handle OOV for some extent capture semantic little bit
- Disadv: sparsity can't handle OOV completely can't capture ordering of word can't do semantic well

For Ex: This is a very good movie

This is not a very good movie

\rightarrow Bow denote these two sentence are very related but they are completely opposite.

Adv: used in information retrieval

Disadv: sparsity

OOV dimension high

$$IDF(t) = \log_e \frac{\text{Number of occurrence of term } t \text{ in document}}{\text{Total number of doc with term } t \text{ in them}}$$

↑ minimize IDF using loge

When loge(1) so $\log_e(1+1) = 2$

$Tf \leftarrow$ particular term how important for

That Document

$Tf \times IDF$

* Custom Features

For Ex Number of +ve reviews number of -ve reviews positive reviews, number of words negative

	watch	campus	people	write	campus + write
1	1	0	0	0	0
0	1	1	0	0	0
0	0	0	1	1	0
0	0	0	1	0	1

deep learning base technique

* word2vec

→ word embedding: In NLP word embedding is a term for the representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning. Ex: BOW, Tf-IDF Techniques

embedding types

frequent base
BOW
TF-IDF
Glove
(matrix factorization)

Prediction base
word2vec

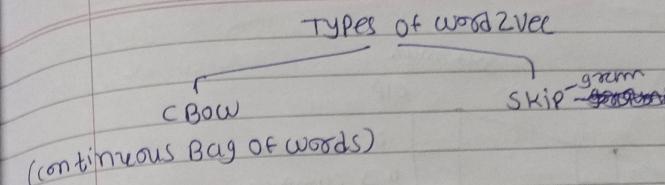
Adv: semantic meaning
low dimension vector
dense vector (Non zero values in matrix) (Not sparse)

→ we can use pretrain model and self-train model
→ we will use the pretrained weights of word2vec that was trained on Google news corpus containing 3 billion words. This model consists of 300-dimensional vectors for 3 millions words and phrases.

Intuition

	King	Queen	Mum	Woman	Monkey
Gender	1	0	1	0	1
Wealth	1	1	0.3	0.3	0
Power	1	0.7	0.2	0.2	0
Weight	0.8	0.9	0.6	0.5	0.3
Speak	1	1	1	1	0

$$\text{King} - \text{mum} + \text{woman} \approx \text{queen}$$



⇒ C-BOW

Fake problem
dummy problem
context word

Target / context word

watch (campus X) for, dutch, science

→ our target is to convert above sentence into 3D vector

target ← for window size 5
context word

training training of deep learning model

X (IP)

watch, for
campusX, dutch
for, science

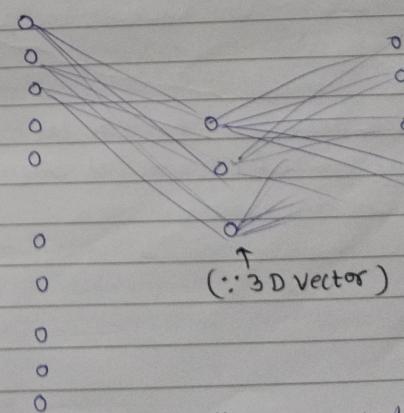
Y (OP)

campusX
for
dutch

window size = 3

watch [1 0 0 0]
campus [0 1 0 0]

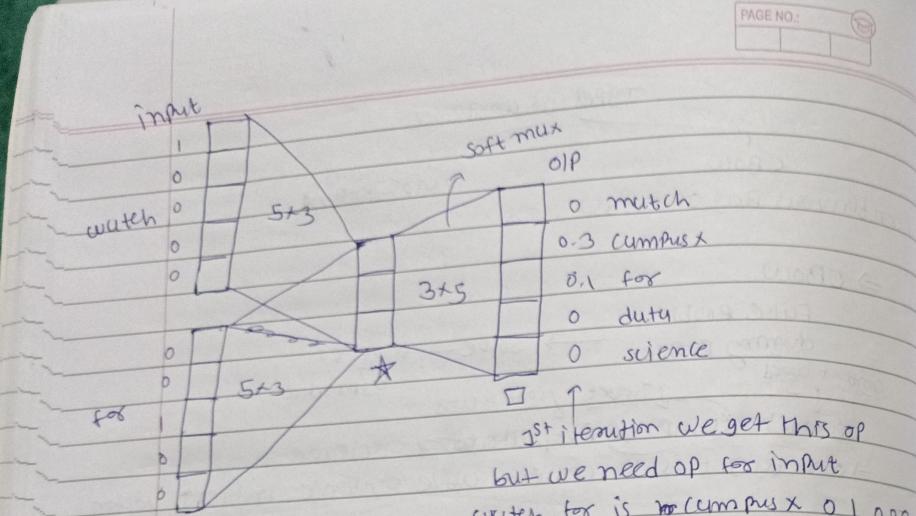
Creating neural network for above training data



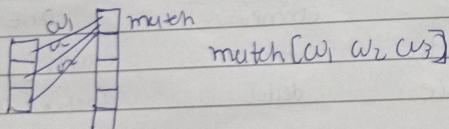
Fully connected



Scanned with OKEN Scanner



⇒ Vector representation of any word is given by three weights ~~from~~ coming from \star to \square



small data = CBow gives best result
large data = Skipgram

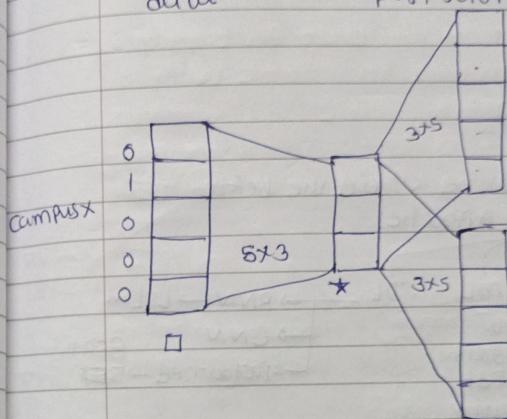
⇒ Skipgram

x (inputs)

campusx → watch, for
for → campusx, duty
duty → For, science

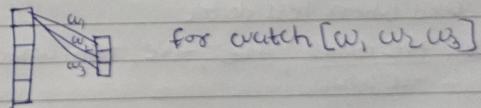
y (outputs)

window size = 3



⇒ for input campusx OIP should be watch, for in every iteration updates weights and minimize loss

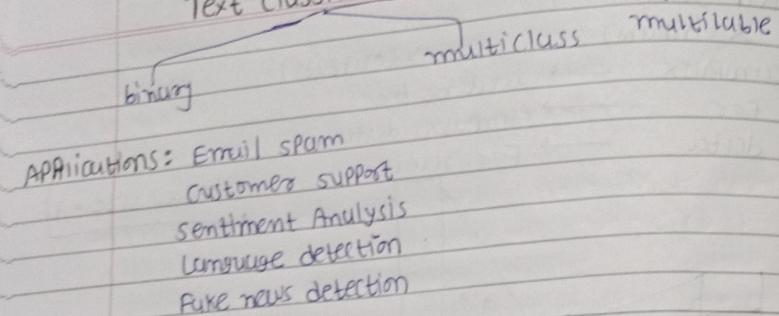
⇒ Vector Form of word is ~~comes~~ corresponding weight of \star to \square



⇒ to improve wordvec quality

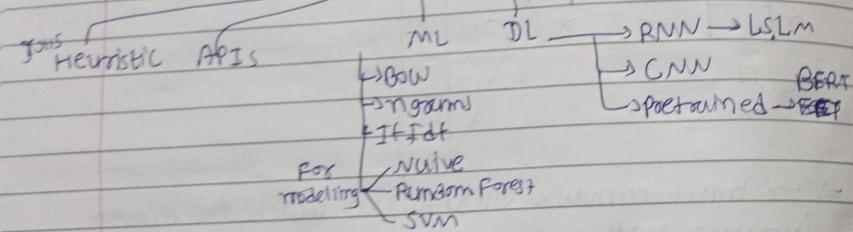
- ↑ increases training data
- ↑ increases dimension of vector
- ↑ increase window size

Text classification



→ Here we use same pipeline as we use before.

Different Approaches



i) Heuristic Approach

Used when less data

Ex: +ve words, -ve, -ve, total words
-ve words

ii) APIs

⇒ Amazon, Google have built in solution using ~~their~~ their

API and can ~~do~~ do classification of text.

disAdv: Those company cost money for this task

iii) Machine learning Approaches

→ using BoW

→ Preprocessing

Naive
RF

Compare these two models

Tf-Idf

word2vec

Pretained
word2vec

on your
data

use only if you
have enough data

use only if more than
80% Vocabulary match

Practical advice: ensemble technique

heuristic features

Deep learning

make sure you have balance dataset

Some As many Project as you can



Part of speech Tagging in NLP (POS Tagging)

Plane of attack
 1) What is POS tagging
 2) Applications
 3) Spacy → Code demo
 4) HMM → Viterbi algo
 POS tagging is a preprocessing step

- ⇒ In simple words, we can say that POS tagging is a task of labelling each word in sentence with its appropriate part of speech.
- ⇒ In traditional grammar, a part of speech or part-of-speech is a category of words that have similar grammatical properties.

Why not tell someone?
 adverb adverb verb noun punctuation mark,
 sentence closer

- Application:
- 1) Name entity recognition
 - 2) Question Answering system
 - 3) Word sense disambiguation
 - 4) Chatbox

How POS tagging works

→ First of all we have some dataset which are labeled by human itself, and using this we will train our HMM model

Noun Word Model
 Nitish loves campusx
 Nitish google campusx
 Nitish google campusx
 Ankita google campusx
 Ankita loves will
 will loves google
 will will google campusx

PAGE NO.:

PAGE NO.:

⇒ Emission Probability

word	N	M	V
Nitish	2/10	0	0
loves	0	0	3/5
campusx	3/10	0	0
google	1/10	0	2/5
will	2/10	1/2	0
Ankita	2/10	0	0
cum	0	1/2	0

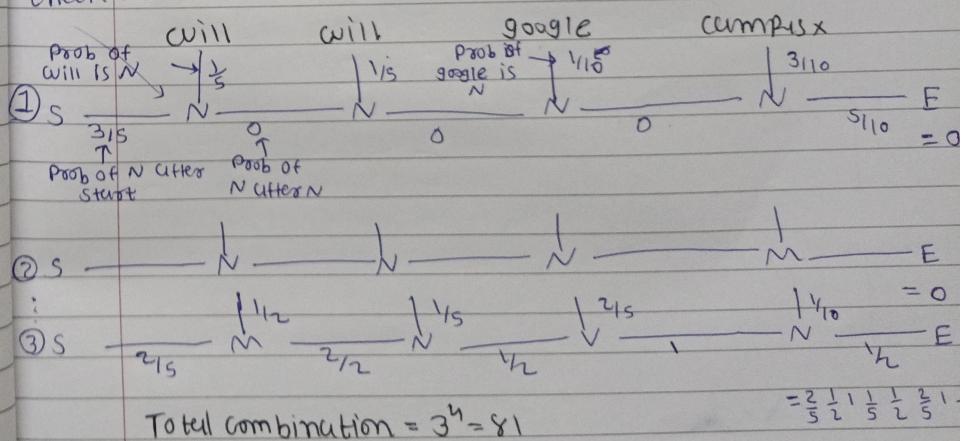
← given N, M or V

Probability of word

⇒ Transition probability (What is prob of next.)

	N	M	V	End
Start	$\frac{3}{10}$ $N \rightarrow N$	$\frac{2}{5}$ $N \rightarrow M$	$\frac{0}{5}$ $N \rightarrow V$	$\frac{0}{5}$ $N \rightarrow End$
N	0	0	$\frac{5}{10}$ $M \rightarrow V$	$\frac{5}{10}$ $N \rightarrow E$
M	$\frac{2}{10}$ $V \rightarrow N$	0	0	$\frac{0}{5}$ $M \rightarrow E$
V	$\frac{5}{10}$ $V \rightarrow M$	0	0	$\frac{0}{5}$ $V \rightarrow E$

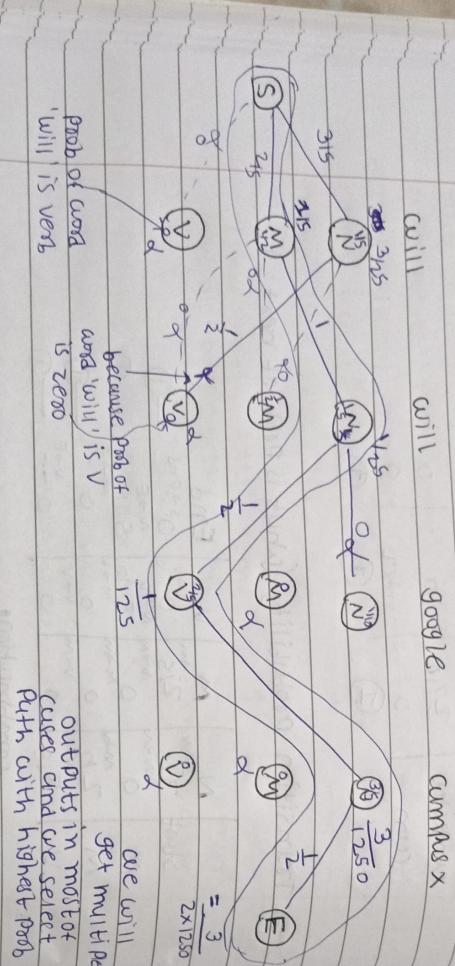
combination
Check for every combination



Scanned with OKEN Scanner

will will google cumpusx
 N N N N
 S V V V V
 M M M M

⇒ Viterbi Algorithm



because prob of word 'will' is V
 'will' is verb
 outputs in most of cues and we select Path with highest prob