

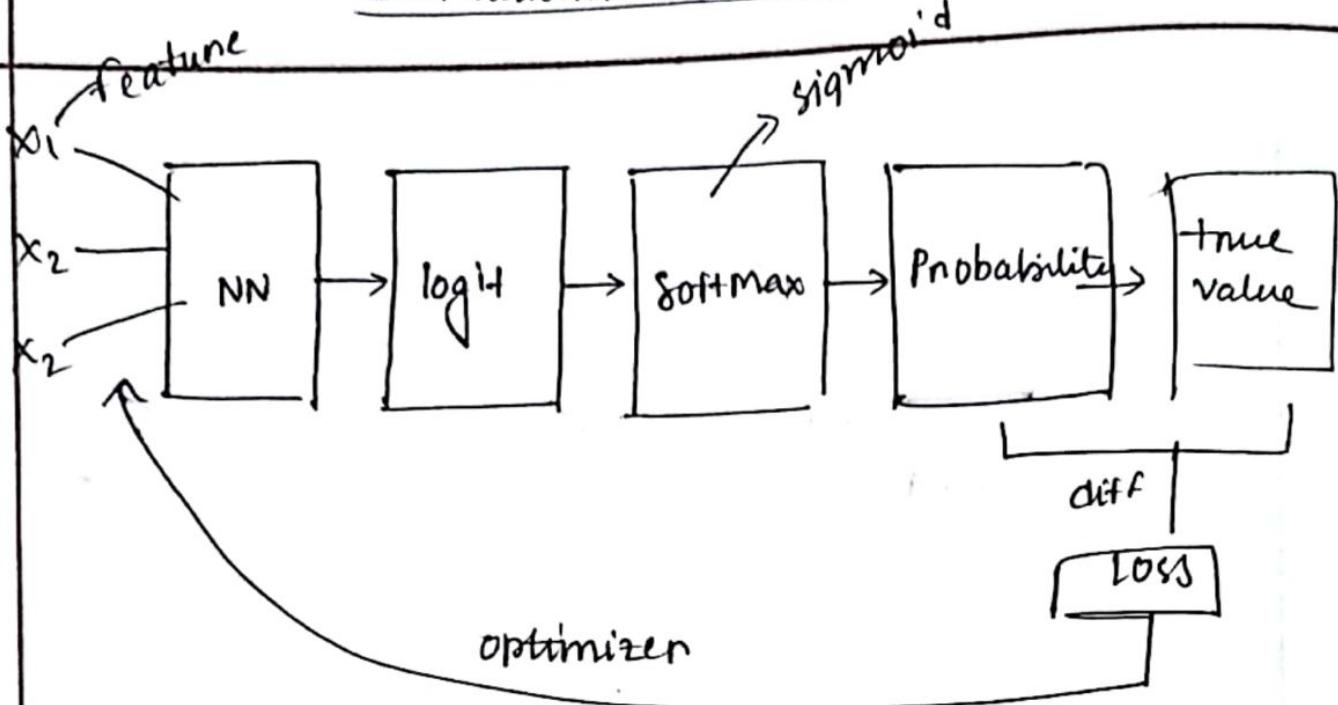
CSE: 4237

Topics:

1. Introduction to Neural Network
2. Neural Network Basics
3. Cross Entropy Loss
4. Shallow Neural Network
5. Hyperparameters and parameters
6. CNN
7. Parameter Calculation
8. Word Embedding
9. Natural Language Processing and Word Embedding
10. RNN, LSTM (incomplete)
11. HMM (incomplete)

- Please do not rely on my calculations—there is a good chance they are wrong.
- If you don't understand any word, there is a good chance that I don't understand/remember it either.

Introduction to Neural Networks



Softmax multiclass classification

- Starts with a dataset that contains n input features.
- The input features are fed into a neural network. The neural network can consist of several layers, including, input, hidden, output
- A forward propagation is performed propagating the input features through the neural network.

$$y = \text{softmax}(K_1 w_1 + K_2 w_2 + \dots + K_n w_n + b)$$

After forward pass, the neural network produces output values called logits

- Then the logits become input to activation function sigmoid / softmax. If multiclass softmax gives probability of each class. Sigmoid fnc. pass into regular linear regression (with logistic regression 20).

Difference between machine learning and deep learning

ML

→ works on small amount of data

→ dependent on low-end machine

→ divides tasks into subtask,
solves individually and finally combines result

→ takes less time to train

→ testing time may increase

Supervised learning

We are given a dataset and we already know what our correct output should look like

→ We have some input (x) and we learn a function to map output (y)

DL

→ works on large amount of data

→ dependent on high-end machine

→ solves problem end-to-end

→ long time to train

→ less time to test data

Structured data

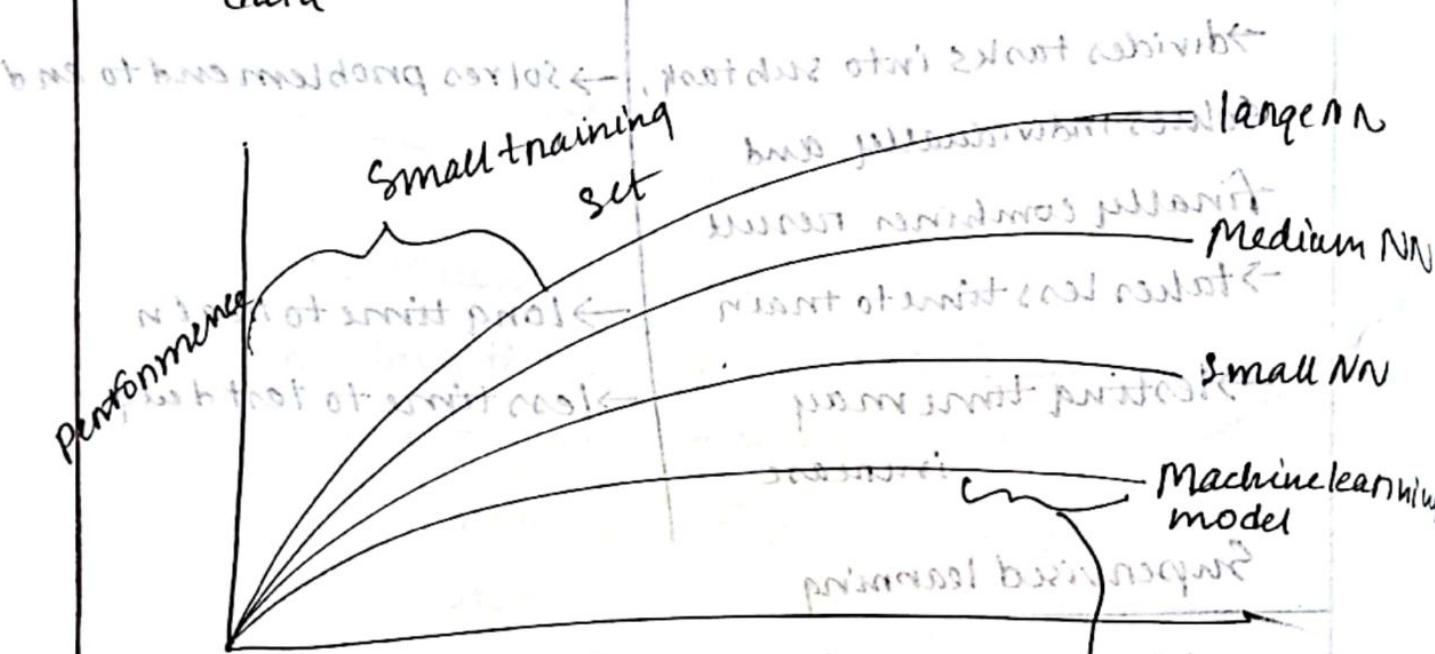
→ databases and tables

JM

Unstructured data

→ images, video, audio, text

→ harden to computer to make use of this kind of data

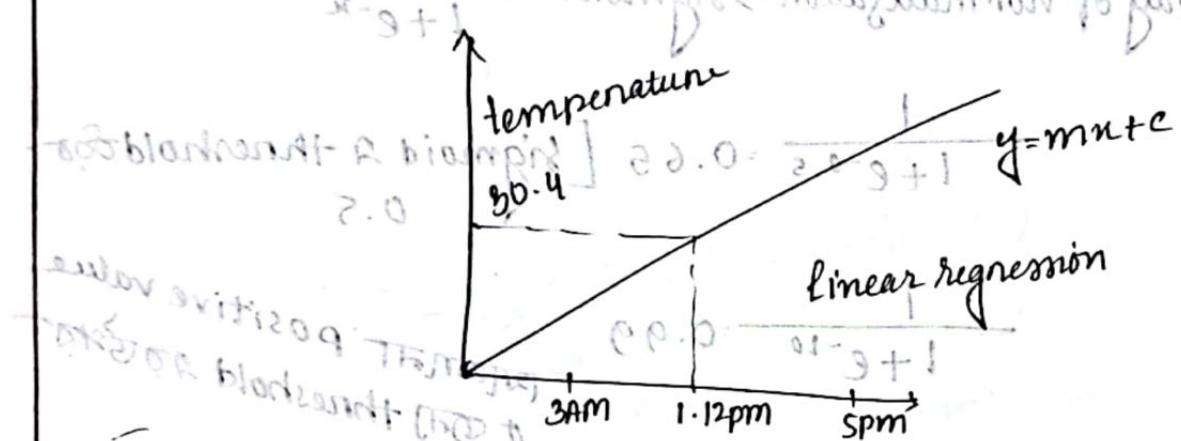


→ few small phrasos such as the graph is very simple
Amount of data
(labeled) → data is not performance

Single Neural Networks

→ Given data about the size of houses, you want to fit a function that will predict the price of the houses.

→ It's linear regression because the price as a function of size is a continuous output



Regression → continuous value (or output)

Classification → discrete value (or output)

Regression > threshold value → classification

This is logistic regression.

Neural network unit - logistic regression

Suppose here threshold is 25.

$y > 25 \rightarrow \text{hot}$

$y < 25 \rightarrow \text{cold}$

Logistic Regression: If we pass linear regression through sigmoid function, we get logistic regression, which is a classification.

$y = \text{min}(x, 1 - x)$, or normalized $0 \leq y \leq 1$.

way of normalization \rightarrow sigmoid $= \frac{1}{1 + e^{-x}}$

$$\text{Sigmoid} \quad \frac{1}{1 + e^{-0.5}} = 0.65 \quad [\text{Sigmoid } \approx \text{threshold } 0.5]$$

$$\text{Nonlinear point} \quad \frac{1}{1 + e^{-10}} = 0.99 \quad [\text{constant positive value } \approx 1 \text{ at threshold } 0.5]$$

$$\text{Nonlinear function} \quad \frac{1}{1 + e^{0.5}} = 0.37 \quad [\text{constant negative value } \approx 0 \text{ at threshold } 0.5]$$

changes between 0 to 1.

$$\rightarrow \text{Sigmoid function} \quad \frac{1}{1 + e^{100}} \approx 0 \quad [\text{sigmoid is a non-linear function}]$$

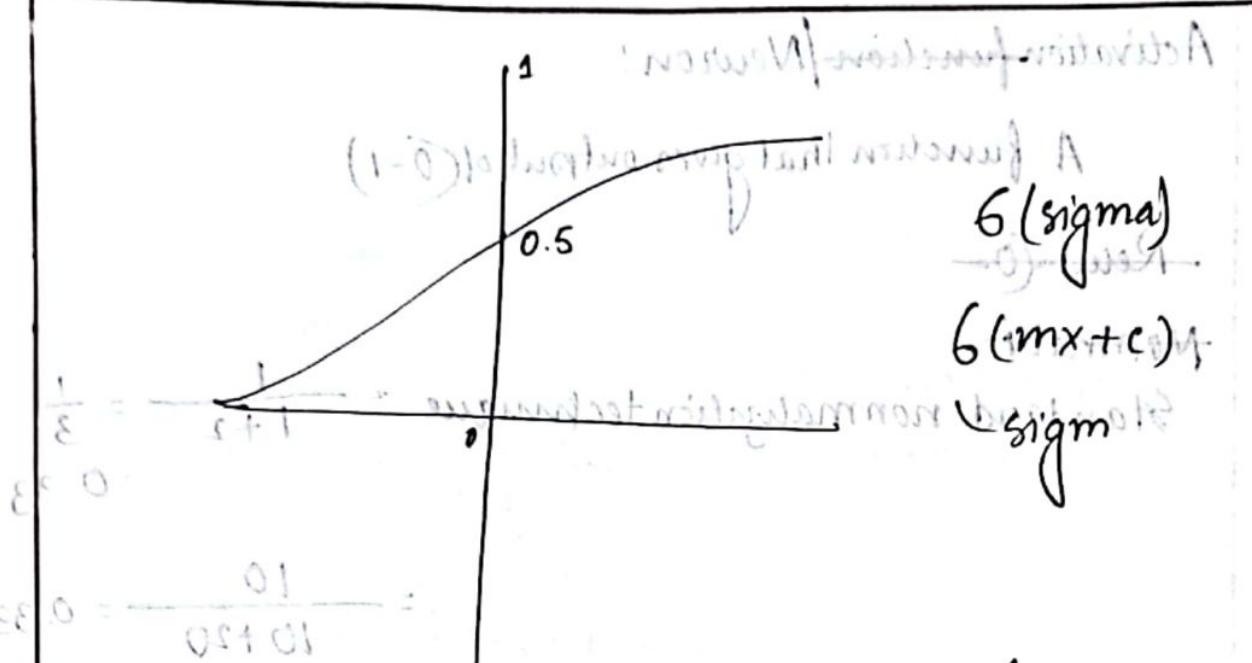
non-linear function \rightarrow non-linear function

non-linear sigmoid \rightarrow linear function

$y = \text{sigmoid}(x) \approx 0.5$

task \rightarrow 2 classes

class \rightarrow 2 sets



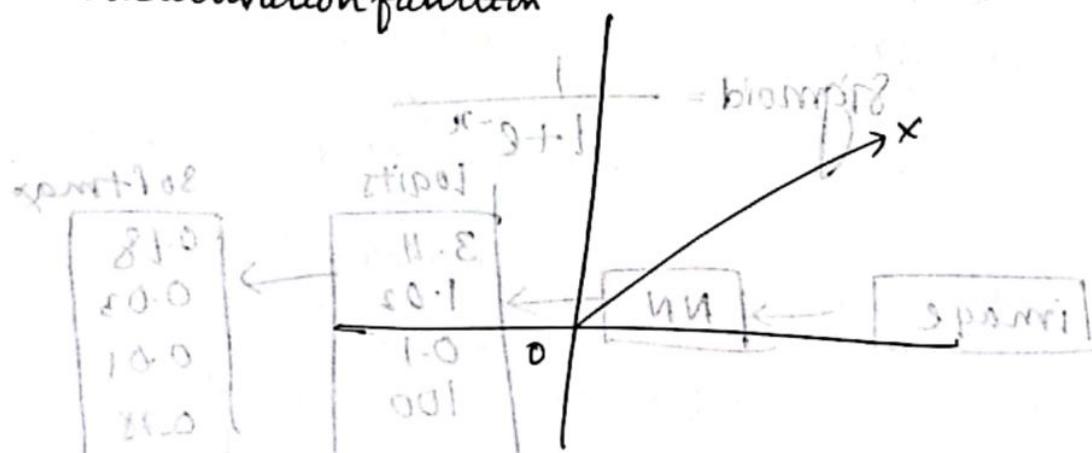
$$SE_0 = \frac{0.1}{10.450} =$$

$m+e \rightarrow$ **Sigmoid** (Activation function)

Activate or deactivate AS THE NEED

sigmoid function based on value

Rectified Linear Unit
ReLU activation function



input गरिए तो output गरिए तो negative
एवं value रो शून्य हो

Activation function/Neuron:

A function that gives output of (0-1)

~~(sigmoid)~~
ReLU (0)

~~(softmax)~~

Standard normalization technique

$$= \frac{1}{1+2} = \frac{1}{3}$$
$$= 0.33$$

$$= \frac{10}{10+20} = 0.33$$

(softmax)

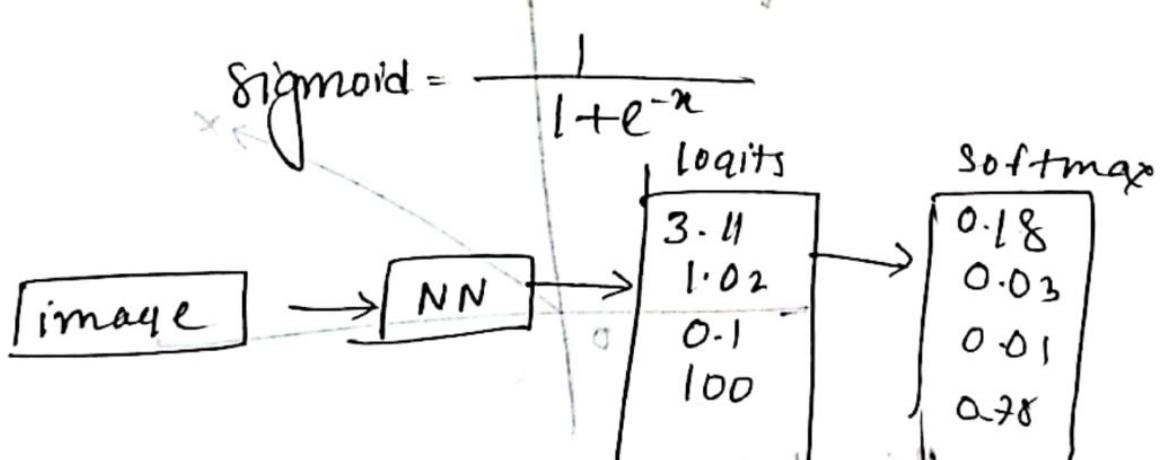
$$\text{softmax} = \frac{e^{x_i}}{\sum^n e^{x_i}}$$

softmax = $\frac{e^{x_i}}{\sum^n e^{x_i}}$

softmax = $\frac{e^{x_i}}{\sum^n e^{x_i}}$

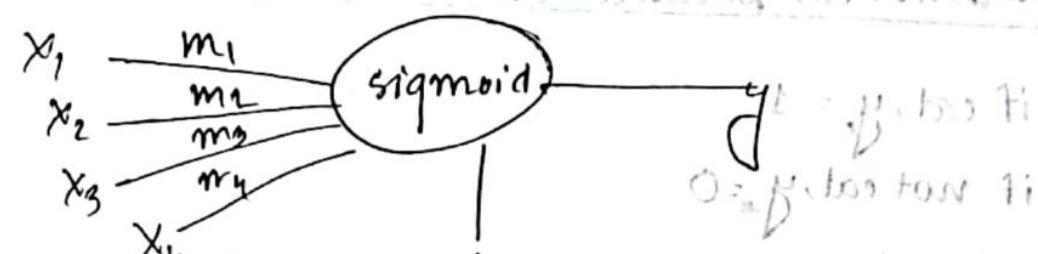
softmax uses standard normalization

softmax uses standard normalization



softmax = $\frac{e^{x_i}}{\sum^n e^{x_i}}$

To create feature vectors, the pixel intensity values are unrolled and shaped for each color.



dimension = input + 1

$$\hookrightarrow \text{feature } 0 = 50.0 - 1 =$$

$\#(n+1)$ dimension for sigmoid input and 1 for output

$$y = f(x) = \sigma(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + c)$$

Image feature

Image \Rightarrow pixels \Rightarrow features, \Rightarrow dimension

$$24 \times 24 \text{ Image} \Rightarrow \text{feature} = 24 \times 24 \times 3 = 1728 \rightarrow \text{RGB}(3 \text{ frames})$$

Neural network \Rightarrow input \Rightarrow feature matrix

(or 1D row vector)

$$64 \times 64 \times 3 = 12288$$

for 64x64 image \Rightarrow 1D \Rightarrow $(1, 64 \times 64)$, $(1, 64 \times 64)$

1 row, 12288 columns

transpose \Rightarrow 1 column, 12288 rows

neural network \Rightarrow data "flatten" \Rightarrow 144 matrix

Neural Network Basics

Cat & not cat prediction - classification

if cat, $y_t = 1$

if not cat, $y_t = 0$

after first iteration $y_i = 0.02$ (was supposed to be 1)

$$\text{loss} = \text{target} - y_i + \text{regularization term}$$

$$= 1 - 0.02 = 0.98$$

want to decrease loss function by gradient descent (1+ε)th

Optimization

Logistic regression

Logistic regression is a learning algorithm used in supervised learning problem when the output y are all ~~are~~ either zero or one

→ The goal of logistic regression is to minimize the error between its predictions and training data

→ Given an image represented by feature vector x , the algorithm will evaluate the probability of a cat being in the picture

Given x , $y = P(y=1|x)$, where $0 \leq y \leq 1$

The parameters

used to start

The input feature vector: $u \in \mathbb{R}^{n_u}$ where n_u is the number of features.

The training label $y \in \{0, 1\}$

parameters | The weights: $w \in \mathbb{R}^{n_w}$ where n_w is the number of features
the bias $b \in \mathbb{R}$

The output $y = \sigma(w^T u + b)$

weight vector

transposed

Sigmoid function $s = \sigma(w^T u + b) = \sigma(z) = \frac{1}{1 + e^{-z}}$

maps continuous variables into probabilities

(0 or 1). therefore suitable for binary classification

output is either 0 or 1. with probability of 0.5

0.5

pd to choose 0 or 1 for further info. Normalized pd

→ make decision based on threshold

Role of bias

understanding SFT

- The bias helps the activation function to be shifted to left or right, to better fit the data.
- changes in weights affect the steepness of the sigmoid curve, whilst the bias offsets it, shifting the curve entirely so it fits better.
- Bias can only influence the output values, it does not interact with the actual input data.

Bias is a measure of how easy it is to get a node to fire

- → for a node with high bias, the output tends to be intrinsically high, with small positive weights and inputs producing large positive output (near to 1)
- Bias can be negative, leading to sigmoid output near 0
- If bias is too small, the output will be decided by the value of weights and input alone

If z is a large positive number, then $\sigma(z) = 1$

If z is a small or large negative number then $\sigma(z) = 0$

If $z=0$, then $\sigma(z)=0.5$.

Cost function

To train the parameters w and b , we need to define a cost function.

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b) \text{ whenever } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$$

Given $(x^1, y^1), (x^n, y^n)$ we want $\hat{y}^{(i)} \approx y^{(i)}$

not violating $\hat{y}^{(i)} \approx y^{(i)}$

predicted value \approx actual value

Loss function

loss function measures the discrepancy between the prediction $\hat{y}^{(i)}$ and the desired output $y^{(i)}$

→ In other words, the loss function computes the error for single training sample.

Error loss function

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

Here is an extra $1/2$ in the equation.

It is because when we take the derivative of cost function, that is used in updating the parameter during gradient descent, that 2 in the power get cancelled with the $1/2$ multiplier.

Is squared error function a good choice for?

→ The square error function (commonly used function for linear regression) is not very suitable for logistic regression.

In case of logistic regression, the hypothesis/prediction is non-linear (sigmoid function), which makes the square error function to be non-convex.

(→ On the other hand, logarithmic function is a convex function for which there is no local optima, so gradient descent works well.)

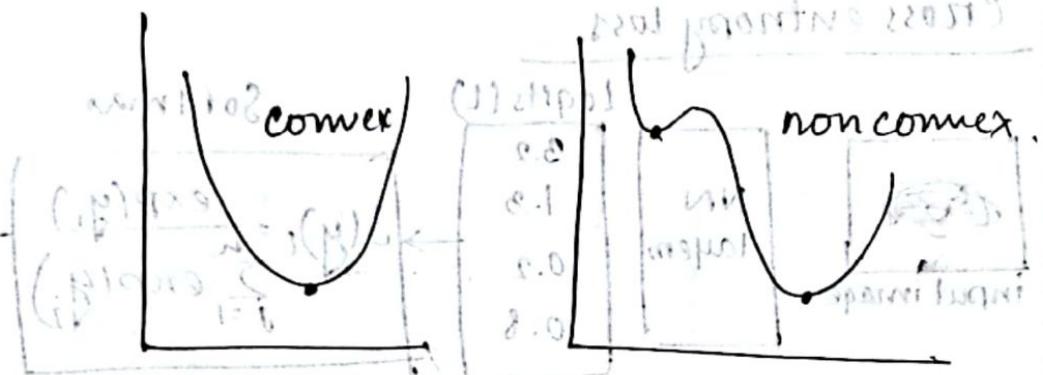
two
variables

poly

tos

separ

polys



If you are doing binary classification, squared error function generally also penalizes examples that are correctly classified but are still near decision boundary - thus creating a 'margin' where misclassifications are not allowed.

Gradient descent waste a lot of time getting prediction very close to {0,1}

Is it necessary for loss function to be convex?

Convexity is a desirable property for loss function in optimization problems. A convex loss function guarantees the existence of a unique global minimum. It simplifies the optimization by ensuring that the solution is not trapped in local optima, making it easier to find best possible solution.

Slide: Logistic Regression / Neural Network Basics

Get to VC plane Gradient Descent when working on FC

Logistic regression: cross entropy loss

$$L(\hat{y}^{(i)}, y^{(i)}) = -\left(y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})\right)$$

If $y^{(i)}=1$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$ where $\log(\hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 1

If $y^{(i)}=0$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1-\hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 0.

Cost function

Cost function is the average of loss function of the entire training data set. We are going to find the parameters w and b that minimises the overall cost function.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

So far we know

$$y = \sigma(w^T n + b), \sigma^{(2)} = \frac{1}{1+e^{-x}}$$

$$\begin{aligned} J(w, b) &= -\frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})] \end{aligned}$$

Gradient descent

Our cost function is a convex.

→ first we initialize w, b to 0, 0 or random values in the

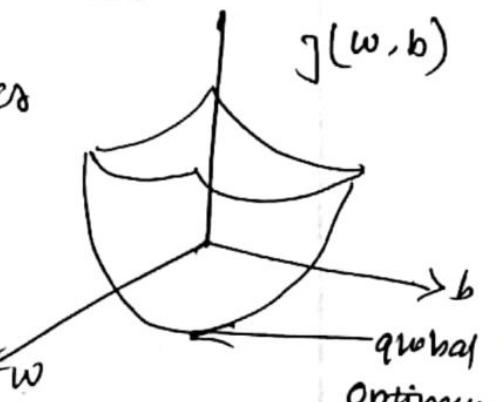
convex function and try to improve the values to reach minimum value.

we want parameters that minimizes

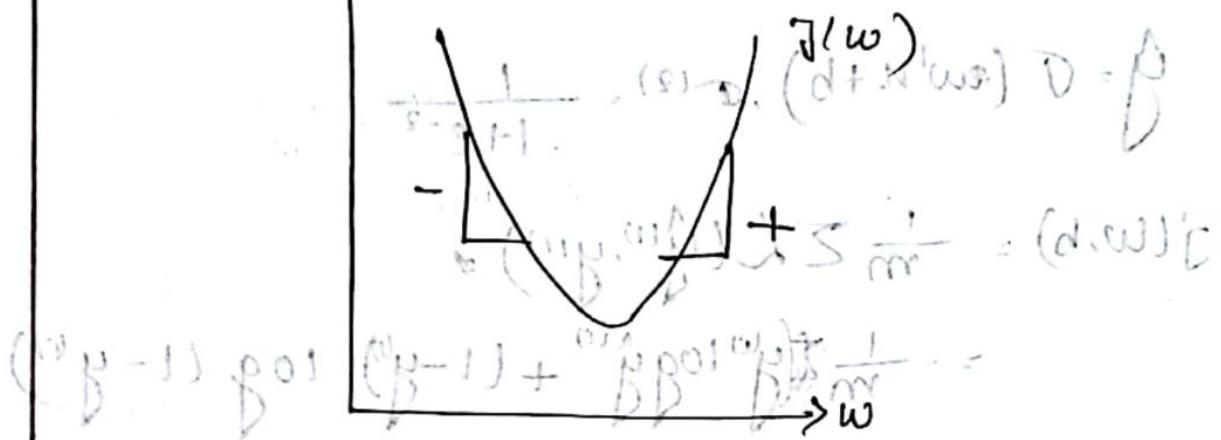
$$J(w, b)$$

Gradient starts at the initial point and take a step in the steepest downhill direction after each iteration.

It will try to reach global optimum.



Gradient Descent



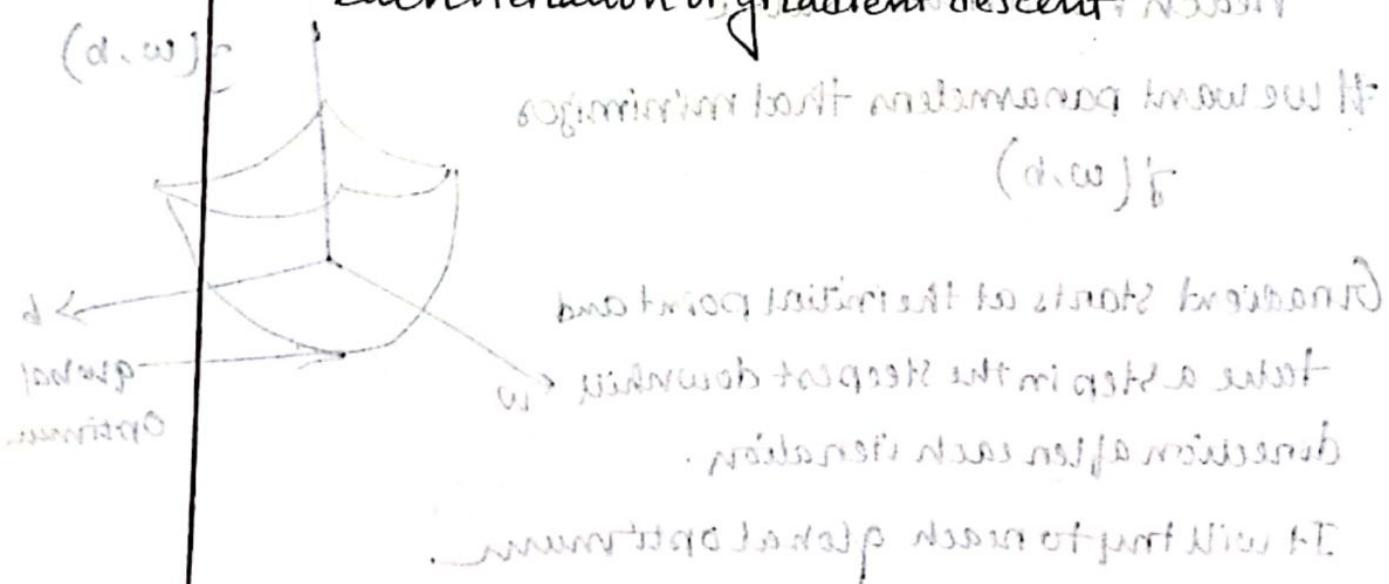
$$w := w - \alpha \frac{dJ(w)}{dw}$$

Model Weights

Learning Rate

Learning rate

Set up a new problem to 0.0 of classification error
of set α learning rate, how biggen steps we choose at
each iteration of gradient descent



Actual update rule: w, b updated in $J(w, b)$

We want parameters w, b to minimize $J(w, b)$

→ to find a point b we take $\frac{\partial J(w, b)}{\partial w}$ to minimize $J(w, b)$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

Same with w ! Weights of change μ : 0 until α

derivative instant function rate of change of $f(a)$

graph which is not too flat between two points

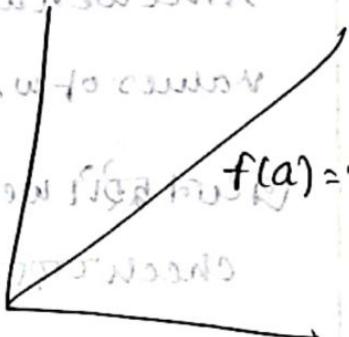
a

function with a lot of steps and jumps, we have to calculate

$\rightarrow 100$

$$a = 2, f(a) = 6$$

$$a = 2.001, f(a) = 6.003$$



If we shift a by 0.001, $f(a)$

shifts by 3 times 0.001

$$a = 5, f(a) = 15$$

$$a = 5.001, f(a) = 15.003$$

Hence slope is 3

$$\frac{d f(a)}{da} = 3$$

The mean rate of change is 3

Do we actually need gradient descent? for large datasets

Let's say, we have 1 weight. To find the ideal weight
to minimize our cost, we need to try a bunch of
values for w , (1000 values) (about 1000 iterations)

$$\text{cost}(w) = \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

→ it takes 0.04 seconds to check 1000 different
weights values on our neural network on my machine

Since we have computed the cost for a wide range
values of w , we can just pick one with the smallest
cost

Get this weight from among 1000 values

check matrix (1000x1000) 1000 points
1000 points

$$z = (A)x - b$$

$$100 z = 100 A x - 100 b$$

\therefore equal matrix

$$z = \frac{(A)x - b}{100}$$

↓ 1000 points for start weight

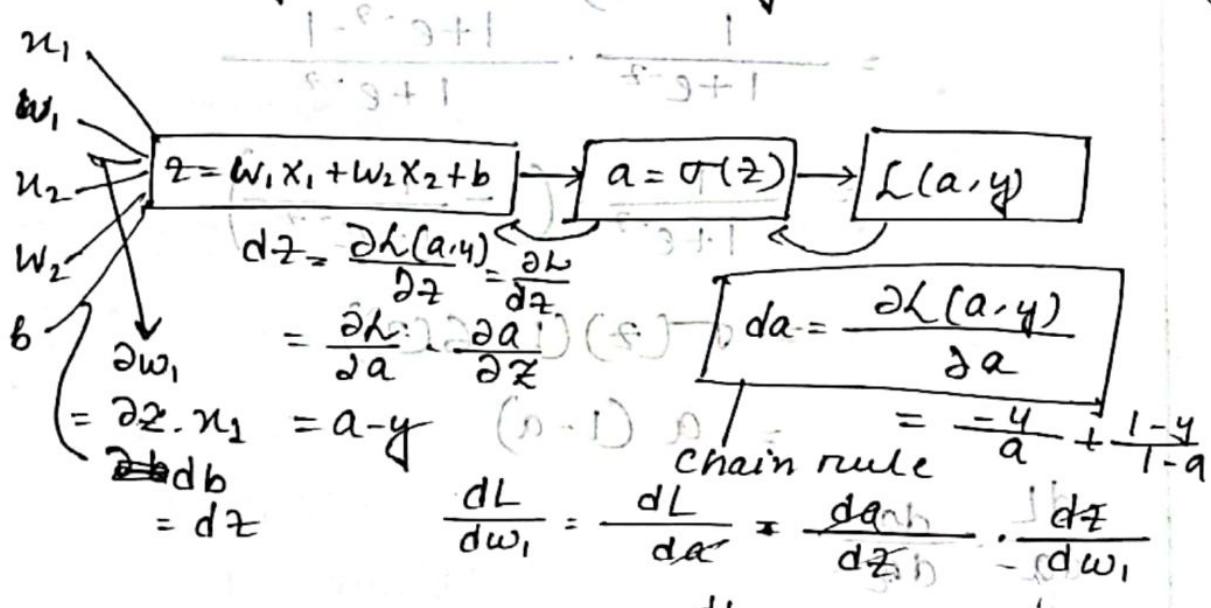
Logistic Regression: Forward Propagation

$$z = w^T x + b$$

$$y = a = \sigma(z) \rightarrow \text{Prediction / Sigmoid output}$$

$$L(a, y) = -(y \log a + (1-y) \log(1-a)) \rightarrow \text{cross entropy loss}$$

computing loss of a single training example



$$L(a, y) = -y \log a - (1-y) \log(1-a)$$

$$\frac{dL}{da} = \frac{d}{da} (-y \log a - (1-y) \log(1-a))$$

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

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

$$\frac{da}{dz} = -\frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right)$$

$$= -\frac{d}{dz} \left(\cancel{1} + e^{-z} \right)^{-1}$$

$$= -1 \times (1+e^{-z})^{-2} (-e^{-z})$$

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

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

$$\boxed{(p, q) \lambda} \xrightarrow{\frac{(p, q) p = p}{1+e^{-z}}} \boxed{\frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}} \right)} = f$$

$$\boxed{\frac{(p, q) \lambda G}{pG}} \xrightarrow{\cancel{pG}} \boxed{f(z)(1-f(z))}$$

$$\frac{dL}{da} = \frac{da}{dz} \cdot \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-z}}$$

$$= \left(\frac{-y}{a} + \frac{1-y}{1-a} \right) \times \frac{1}{1-a}$$

$$= \frac{-y(1-a) + a(1-y)}{a(1-a)} \times \frac{1}{1-a}$$

$$= \frac{(-y(1-a) + a(1-y))}{a(1-a)} \times \frac{1}{1-a} = \frac{(-y(1-a) + a(1-y))}{a(1-a)} \times \frac{1}{1-a}$$

$$= \frac{(-y + ay + a - ay)}{a(1-a)} \times \frac{1}{1-a} = \frac{a - y}{a(1-a)}$$

$$= a - y$$

$$z = w_0 x_0 + w_1 x_1 + b$$

$$\frac{\partial z}{\partial w_1}$$

$$\frac{\partial z}{\partial w_1} = n_1$$

$$\frac{dL}{dw_1} = \frac{dL}{da} \times \frac{da}{dz} \times \frac{dz}{dw_1}$$

$$\text{Let } a = (a - y)x_1$$

$$\frac{dL}{dw_2} = \frac{dL}{da} \times \frac{da}{dz} \times \frac{dz}{dw_2}$$

Logistic regression, gradient descent from examples

$$J=0, dw_1=0, dw_2=0, db=0, w_0=0, w_1=0, w_2=0, b=0$$

$$\text{for } i=1 \text{ to } m$$

$$z = w_0 x_0(i) + w_1 x_1(i) + b$$

$$a(i) = \text{sigmoid}(z)$$

$$J += y(i) \log a(i) + (1-y(i)) \log(1-a(i))$$

Backward pass: $\Delta = \frac{1}{m} \sum_{i=1}^m$

$$d(z) / dz(i) = a(i) - y(i)$$

$$dw_1^+ = d(z) / dz(i) \cdot x_1(i)$$

$$dw_2^+ = -d(z) / dz(i) \cdot x_2(i)$$

$$db^+ = d(z) / dz(i)$$

$$J/m$$

$$dw_1 / m$$

$$dw_2 / m$$

\sum $\xrightarrow{[1] \text{ (p)}}$ No of layer
 $\xrightarrow{[1]} \text{No of sample}$
 $\xrightarrow{[1]} \text{No of node}$

$$d + ex + w_1 + w_2 + b = 0$$

Gradient descent

$$w_1 = w_1 - \alpha x dw_1$$

$$w_2 = w_2 - \alpha x dw_2$$

$$b = b - \alpha x db$$

Here, we have done one step of gradient descent.
we need to repeat it multiple times in order to take multiple steps of gradient descent.

In this code, we implement two for loops, which is not very efficient. Solution = vectorization technique

Modified LR Gradient Descent

$$\theta = 0, [dw_1=0 \quad dw_2=0] \quad db=0, w_1=0, w_2=0, b=0$$

$$dw = np.zeros((nx, 1)) \quad (nx, w = 8)$$

In backward pass

(el bias = 1)

previously modified

$$dw_1 += \frac{dZ(i)xw_1(i)}{m} \quad dw_1 = h(i)x dZ(i) \quad m=15$$

$$dw_2 += \frac{dZ(i)xw_2(i)}{m} \quad dw_2 = h(i)x dZ(i) \quad m=15$$

$$(i) \in X \times (i) \in b = \frac{1}{m} \sum b$$

$$(i) \in X \times (i) \in b = \frac{1}{m} \sum b$$

$$(i) \in b = \frac{1}{m} \sum b$$

$$m = 15$$

$$m = 15 b$$

$$m = 15 b$$

In Gradient Descent,

previously

$$\begin{aligned} d\omega_1 &= d\omega_1/m \\ d\omega_2 &= d\omega_2/m \end{aligned}$$

modified

$$d\omega = d\omega/m$$

$$(d\omega) = \begin{bmatrix} d\omega_1 \\ d\omega_2 \\ \vdots \\ d\omega_m \end{bmatrix} = A$$

$$(d\omega) = B - C = B - (A)^T \cdot (A) \cdot B = B - A^T \cdot A \cdot B$$

Vectorizing logistic regression (forward)

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \\ n^{(1)} & n^{(2)} & \dots n^{(m)} \\ \vdots & \vdots & \vdots \end{bmatrix} \text{ while } R^{n \times m}$$

$$[\omega^T] \begin{bmatrix} n^{(1)} & n^{(2)} & \dots n^{(m)} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{aligned} X &= n \times m \\ \omega^T X X^T &= (1 \times n) \cdot (n \times m) \\ &= (1 \times m) \end{aligned}$$

$$\begin{aligned} z^{(1)} &= \omega^T n^{(1)} + b \\ a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

1st training example

$$\begin{aligned} z^{(2)} &= \omega^T n^{(2)} + b \\ a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

2nd training example

$$\begin{aligned} z^{(3)} &= \omega^T n^{(3)} + b \\ a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

3rd training example

$$[z^{(1)} \ z^{(2)} \ z^{(3)} \dots z^{(m)}] = W^T X + [b \dots b] = [w^T n^{(1)} + b \quad w^T n^{(2)} + b \dots \quad w^T n^{(m)} + b]$$

$\hookrightarrow 1 \times m$

$$z = \text{np.dot}(w.T, x) + b$$

+ b
broadcasting rule
 $\rightarrow (1, 1)$

$w \otimes b = wb$

$$A = [a^{(1)} \ a^{(2)} \dots \ a^{(m)}] = \sigma(z)$$

$$dz = \frac{dL}{da} \cdot \frac{da}{dz} \quad dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

$$dz^{(3)} = a^{(3)} - y^{(3)}$$

~~principal diagonal~~

$$dz = [dz^{(1)} \ dz^{(2)} \dots \ dz^{(m)}]$$

$$A = [a^{(1)} \ a^{(2)} \dots \ a^{(m)}]$$

$$Y = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}]$$

$$(m \times n) \cdot dz = A - Y = [(a^{(1)} - y^{(1)}) \ (a^{(2)} - y^{(2)}) \ \dots \ (a^{(m)} - y^{(m)})]$$

with $m < n$

$$dF^T w^T v = F^T | \quad dF^T w^T v = F^T |$$

$$(F^T)^T v = F v$$

permutation

$$dF^T w^T v = F^T | \quad dF^T w^T v = F^T |$$

$$(F^T)^T v = F v$$

permutes

$$dF^T w^T v = F^T | \quad dF^T w^T v = F^T |$$

$$(F^T)^T v = F v$$

permutes

$$dF^T w^T v = dF^T w^T v = [d \ \ dF^T w^T v \ \ \dots \ \ dF^T w^T v]$$

$$[d \ \ dF^T w^T v \ \ \dots \ \ dF^T w^T v]$$

$m \times n < 1$

Gradient computation

previously

$$dw = 0$$

$$dw^+ = X^{(1)} dZ^{(1)}$$

$$dw^+ = X^{(2)} dZ^{(2)}$$

$$dw^+ = X^{(m)} dZ^{(m)}$$

$$dw^l = m$$

$$db = 0$$

$$db^+ = dZ^{(1)}$$

$$db^+ = dZ^{(2)}$$

$$db^+ = dZ^{(m)}$$

$$db^l = m$$

modified

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)}$$

$$db = \frac{1}{m} np \cdot \text{sum}(dZ)$$

$$dw = \frac{1}{m} \sum_{i=1}^m dZ^{(i)} \cdot X$$

$$dw = \frac{1}{m} \begin{bmatrix} 1 & \dots & 1 \\ n^{(1)} & n^{(2)} & \dots & n^{(m)} \end{bmatrix} \begin{bmatrix} dZ^{(1)} \\ dZ^{(2)} \\ \vdots \\ dZ^{(m)} \end{bmatrix}$$

$$dw = \frac{1}{m} [X^{(1)} dZ^{(1)} + X^{(2)} dZ^{(2)} + \dots + X^{(m)} dZ^{(m)}]$$

Single iteration of gradient descent

$$Z = w^T X + b$$

$$Z = np \cdot \text{dot}(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - y$$

$$dw = \frac{1}{m} X^T dZ$$

$$db = \frac{1}{m} np \cdot \text{sum}(dZ)$$

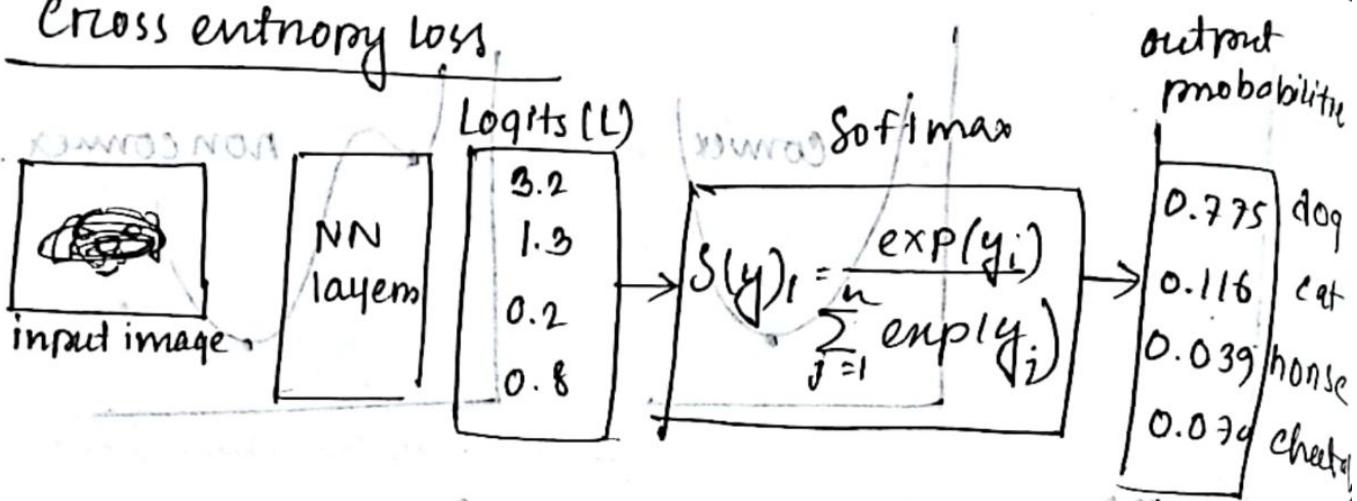
$$w := w - \alpha dw$$

$$b := b - \alpha db$$

[B gradient update]

slide: cross entropy loss

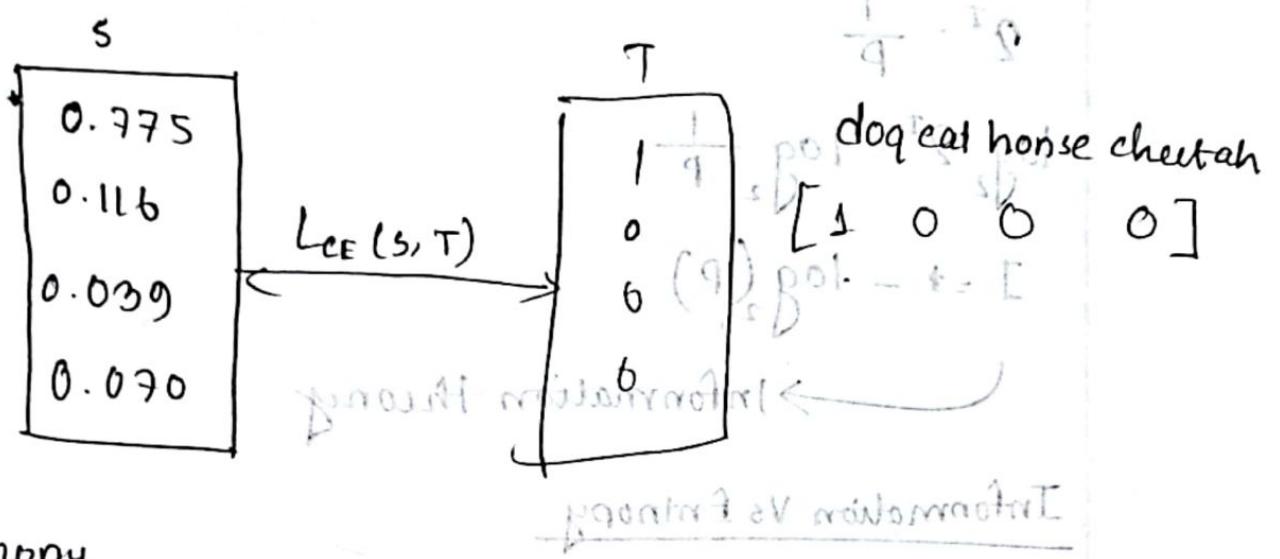
Cross entropy loss



Logits

The vector of raw (non-normalised) predictions that a classification model generates, which is ordinary ordinarily then passed to a normalisation function.

If the model is solving a multi-class classification problem, logits typically become an input to the softmax function. The softmax function then generates a vector of (normalised) probabilities with one value for each possible class.

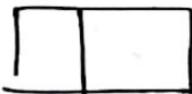


Entropy

Entropy of a non-random variable x is the level of uncertainty inherent in the variable possible outcome.

$$H(X) = \begin{cases} \int_{\Omega} p(x) \log p(x), & \text{if } x \text{ is continuous} \\ -\sum_x p(x) \log p(x), & \text{if } x \text{ is discrete.} \end{cases}$$

$$p = 1/2$$



$$\text{number of box} = 2^{\text{info}}$$

$$\text{number of ques/info} = 1$$

W	O	R	S

$$\text{number of box} = 4$$

$$\text{number of ques} = 2$$

$$\text{number of box} = 8$$

$$\text{number of ques} = 3$$

$$2^{\text{info}} = \text{number of box}$$

$$\text{number of box} = 2^I$$

$$P = \frac{1}{\text{number of box}}$$

$$\text{number of box} = \frac{1}{P}$$



$$2^I = \frac{1}{P}$$

$$\log_2(2^I) = \log_2 \frac{1}{P}$$

$$I = -\log_2(P)$$

→ Information theory

Information Vs Entropy

Information provides a way to quantify the amount of surprise for an event measured in bits.

Entropy provides a way to measure of the average amount of information needed to represent an event drawn from a probability distribution for a random variable.

What is information?

Information is measured in bits. When measuring information, we are measuring the number of bits or the number of "yes/no" questions we would have to ask to reach the same conclusion.

$$\frac{1}{2} = \text{1 bit of information}$$

Reason of negative sign

cross-entropy loss function

negative binomial distribution and sigmoid function

sigmoid function is a probability quantity

$$\log(p(x)) < 0 \text{ for all } p(x) \in (0, 1).$$

Cross entropy loss

→ Cross entropy loss is also called logarithmic loss, log loss, logistic loss.

→ Each predicted class probability is compared to the actual class desired output of 0 or 1.

→ A score loss is calculated that penalizes the probability based on how far it is from the actual expected value

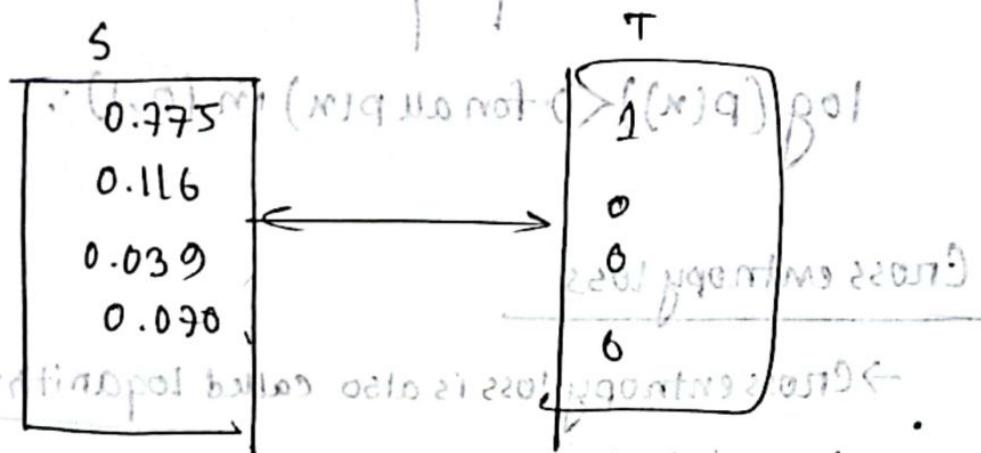
$$f f d e . 0 =$$

Cross entropy loss

prefix writing for to avoid it

$$L_{CE} = -\sum_{i=1}^n t_i \log(p_i), \text{ for } n \text{ classes,}$$

where t_i is the truth label and p_i is the softmax probability for i^{th} class



$$L_{CE} = \sum_{i=1}^n T_i \log(S_i)$$

$$L = -[1 \log_2(0.775) + 0 \log_2(0.116) + 0 \log_2(0.039)]$$

plus denominator with log base 2 of 1 - 0.775 = 0.225

$$= -\log_2(0.775)$$

$$= 0.3673$$

Binary cross entropy loss

loss av 2201

$$\text{Binary cross entropy loss} L = - \sum_{i=1}^n t_i \log(p_i)$$

$$\rightarrow \text{simplified int} = [-t \log(p) + (1-t) \log(1-p)]$$

When using log base 2, the unit of entropy is bits,

whereas with natural log the unit is nats.

• t is probability of class 1, p is predicted probability of class 1

• t is probability of class 1, p is predicted probability of class 1

Categorical crossentropy loss / Sparse categorical crossentropy

Both categorical crossentropy and sparse categorical crossentropy have the same loss function. The only difference between the two is how truth labels are defined.

Categorical crossentropy is used when true labels are one-hot encoded. $[1, 0, 0]$

In sparse categorical crossentropy, truth labels are integer encoded. $[1, 2, 3]$

Loss vs cost

not priorities was priority

The loss function computes the error for a single training example, while the cost is the average of the loss function of the entire training set.

$$L = -\frac{1}{N} \left[\sum_{j=1}^N [t_j \log p_j + (1-t_j) \log (1-p_j)] \right]$$

2 free terms are left

positive losses are positive terms loss op 1

losses for plus will increase if real score with small predictions

losses will decrease if real score is close to target

losses will increase if real score is far from target

$$[0, 0, 1] \rightarrow \text{loss 0.5}$$

$$[0, 0, 1] \rightarrow \text{loss 0.5}$$

WORDLE

The goal of wordle is to guess 5 letter words, with 6 guesses.

After answering, each letter turns green, amber or grey

green \Rightarrow correct spot

Amber = \varnothing letter is in word, but not in correct spot

Grey - not in the word.

5757 five letter words in english.

At beginning = there are $\log_2(5757)$ or 12.49 bits of hidden info.

As you play, you \varnothing acquire more information about the hidden word.

The goal is to create an algorithm that in each guess chooses the word that adds the most info.

(and reduces uncertainty the most)

16

MIDDLE

- the initial guess is 'tires'
- there are only 12 words in English language that matches this criteria
- We have reduced the uncertainty from 12.5 bits to 3.6 bits
- Our guess makes us gain 8.9 bits of info
- after the second guess "onate", the only possible word is great, but we still don't know if it's

On the other hand, if we enter the word 'fuzzy', there are 62% ^{5 letter} words in English that do not contain 'F' 'U' 'Z' 'Y'; so 62% of time we will end up with a word that's gray:
all gray is the outcome with the least amount of information gained.

If we put 'taras' as the initial word, in only 2% of the times we will get a wordle all grey.

93). of the time we will know at least one letter:

$$((B-1) \text{ pot } (B-1) + ((B-1) \text{ pot } B)) - ((B-1) \text{ pot } B)$$

fuzzy reduces our possible word list from

$$((B-1) \text{ pot } \text{ another } ((B-1) \text{ pot } B)) - ((B-1) \text{ pot } B)$$

to at least one letter

taries reduces it to 78.

$$((B-1) \text{ taries } ((B-1) \text{ pot } B)) - ((B-1) \text{ pot } B)$$

to at least one letter

~~Quiz 10~~

Quiz 10

and the remaining 42% of experiments is not covered by

but it is proposed that the other 50% of the time

at minimum half of the words is eliminated with

minimum fuzzy matching

$$((B-1) \text{ pot } B) \sum_{i=1}^m \frac{1}{m} = ((B-1) \text{ pot } B) \sum_{i=1}^n \frac{1}{m} = (dust)$$

$$((B-1) \text{ pot } (B-1) +$$

Slide: Shallow Neural Networks

Activation functions

activation functions can be used in hidden layers and output layers.

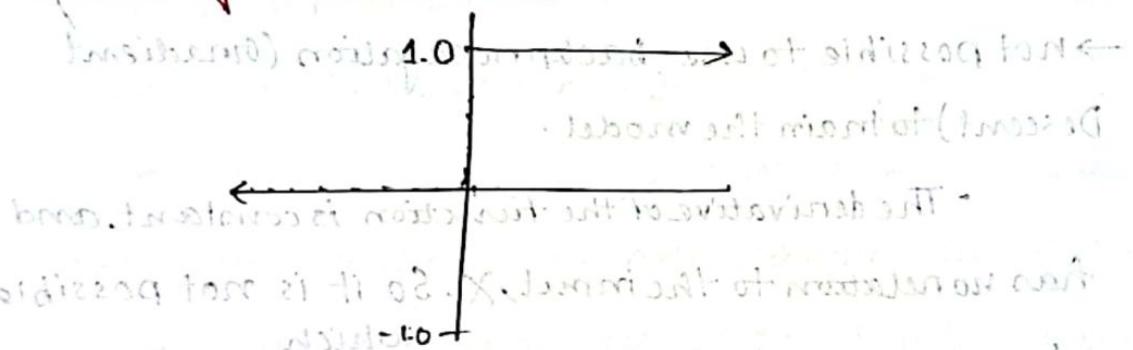
$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

Binary step function



threshold based activation function

- If the input value is above or below a certain threshold, the neuron is activated and sends exactly the same signal to the next layer.

Problem

⇒ does not allow multi-value outputs.

⇒ It is not recommended to use in hidden layers

because it does not represent derivative learning value

Types of activation function

sigmoid activation

Linear activation function.

linear boundary

→ form $A = cx$

⇒ It takes input multiplied by weights of each neuron and creates an output signal proportional to the input.

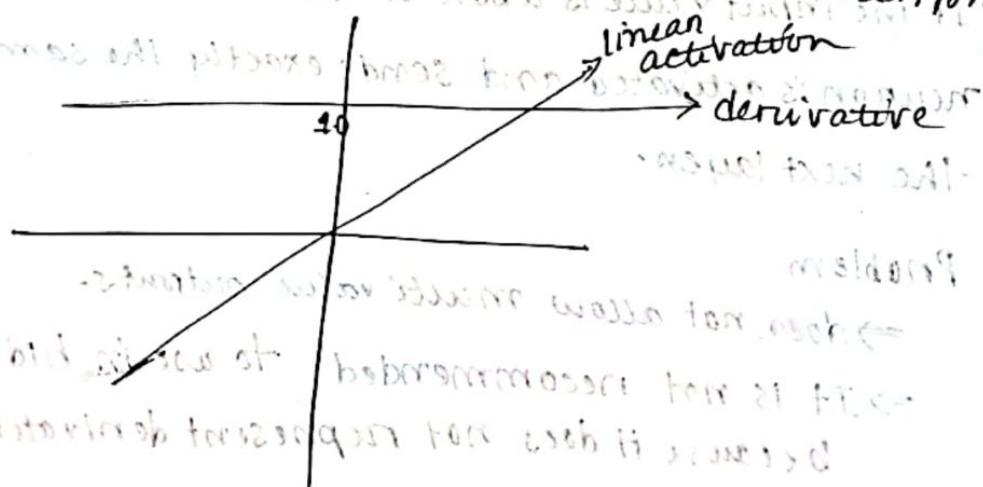
Linear function is better than step function because it allows multiple outputs, not just yes and no.

Problems:

→ not possible to use backpropagation (Gradient Descent) to train the model.

The derivative of the function is constant, and has no relation to the input, x . So it is not possible to go back and understand which weights in the

In the input neurons can provide a better prediction of the output. All pictures come from $\frac{1}{x}$.



- When $A = cx$ is derived from x , we reach c. That is there is no relationship with x in the derivative, therefore there is no learning.

→ All layers of neural network will collapse into one -

- with linear activation function, no matter how many layers in the neural network, the last layer will be a linear function of the first layer (because a linear combination of linear functions is still a linear function). So a linear activation function turns the neural network, on even into a deep neural network into just one layer.

→ A linear acti. A neural network with ~~an~~ activation function on without any activation function is simply a linear regression model. It has limited power and ability to handle complexity varying parameters of input data.

. with linear activation function or no activation function

$$z^{[1]} = W^{[1]}X + b^{[1]}$$
$$a^{[1]} = z^{[1]}$$

Q: Why do we need non-linear activation function?

→ If linear activation function or no activation function is used, the output signal becomes a simple linear function.

→ Linear functions are only single grade polynomials.

⇒ A non activated neural network will act as a linear regression with limited learning power.

→ We want our neural network to learn non-linear states. Because, we will give it complex real world information such as images, videos, texts, and sounds which are non linear and have high dimensionality to process. To learn to train neural network.

→ Multilayered deep neural network can learn meaningful features from data.

⇒ It allows backpropagation because they have a derivative function, which is related to the input.

→ They allow stacking of multiple layers of neurons to create a deep neural network. Multiple hidden layers of neurons are needed to learn complex datasets with high level of accuracy.

for example

$$a^{[1]} = z^{[1]} = W^{[1]}X + b^{[1]}$$

$$a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$= W^{[2]}a(W^{[1]}X + b^{[1]}) + b^{[2]}$$

$$= \underbrace{W^{[2]}W^{[1]}X}_{\text{m}} + \underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{c}$$

Neural network is outputting a linear function of inputs. Composition of two or more linear functions is also a linear function.

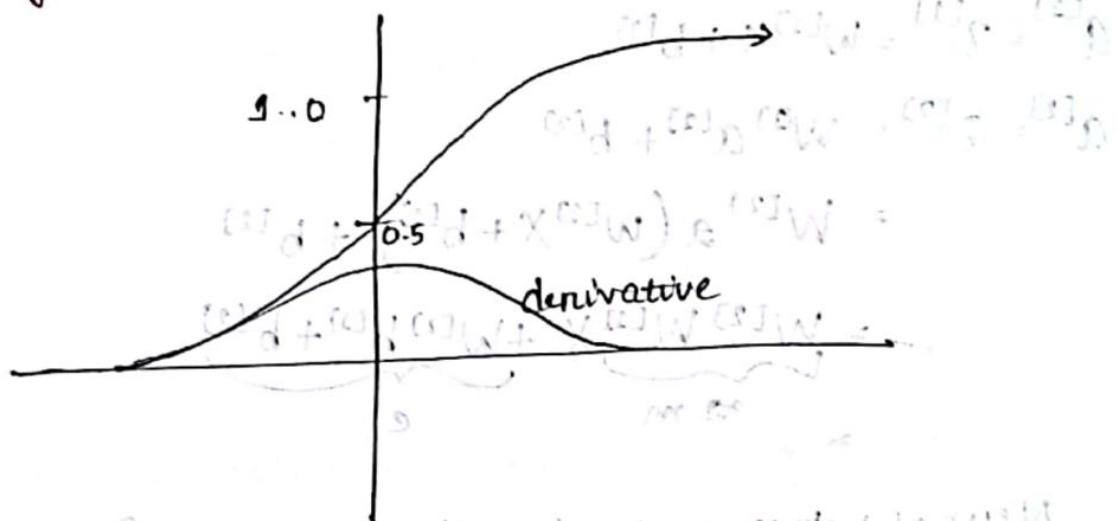
⇒ One place you can use a linear activation function when you are solving a regression problem, meaning the output is a real number. But only in the last hidden layer and not in the intermediate layers.

part of a single page, comes out to be output of (ReLU) function

function of input x is $\max(0, x)$.

For example if $x = 1$ then $\max(0, 1) = 1$. If $x = -1$ then $\max(0, -1) = 0$.

Sigmoid / Logistic function



→ ~~gradient descent is primitive & unsmooth function~~
It can be derived because it is different from

step function, which means learning can happen!

→ smooth gradient, prevents jumps in the

output values.

→ output values bound between 0 and 1.

→ if x is above 2 or below -2, tends to bring y value
(prediction) to the edge of the curve, very close to 1 or 0.

This enables clear prediction.

⇒ The sigmoid function is the most frequently used activation function, but there are alternatives

Problem:

Vanishing Gradient: for very high and very low values of X , there is almost no change to the prediction. The derivative values in these ranges are very small and converge to 0. This is called vanishing gradient and learning is minimal.
→ The network can refuse to learn further, or being too slow to reach an accurate prediction.
→ When slow learning occurs, the optimization algorithm that minimizes error can be attached to local minimum values and cannot get maximum performance from the ~~artificial~~ artificial neural network.
⇒ output is not zero centered. So output of all the neurons will be of same sign.

⇒ computationally expensive

It is caused by softmax so when it consider softmax function had polynomial nature of square function so it is difficult to bring all the numbers to the same range.

• gradient?

Hyperbolic Tangent function (tanh)

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

fundamental basis of deep learning

Linear function: $f(u) = u \quad (-\infty, \infty)$ (straight line)

Step function: $f(u) = \begin{cases} 0 & \text{for } u < 0 \\ 1 & \text{for } u \geq 0 \end{cases}$ (0-1)

Sigmoid function: $f(u) = \sigma(u) = \frac{1}{1+e^{-u}}$ (0-1)

Hyperbolic tan: $f(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$ (-1-1)

Derivative of tanh:

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

$$\frac{d}{dz} g(z) = \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2}$$

$$\frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\frac{1 - \tanh^2(z)}{1} = 1 - \tanh^2(z)$$

$$\exp(-\text{tanh}(z)) = \frac{1}{1 + \tanh^2(z)}$$

$$(XH) = \exp(-\text{tanh}(z)) = \frac{1}{1 + \tanh^2(z)}$$

$$(XH) = \exp(-\text{tanh}(z)) = 1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)^2$$

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

$$= 1 - a^2$$

ReLU

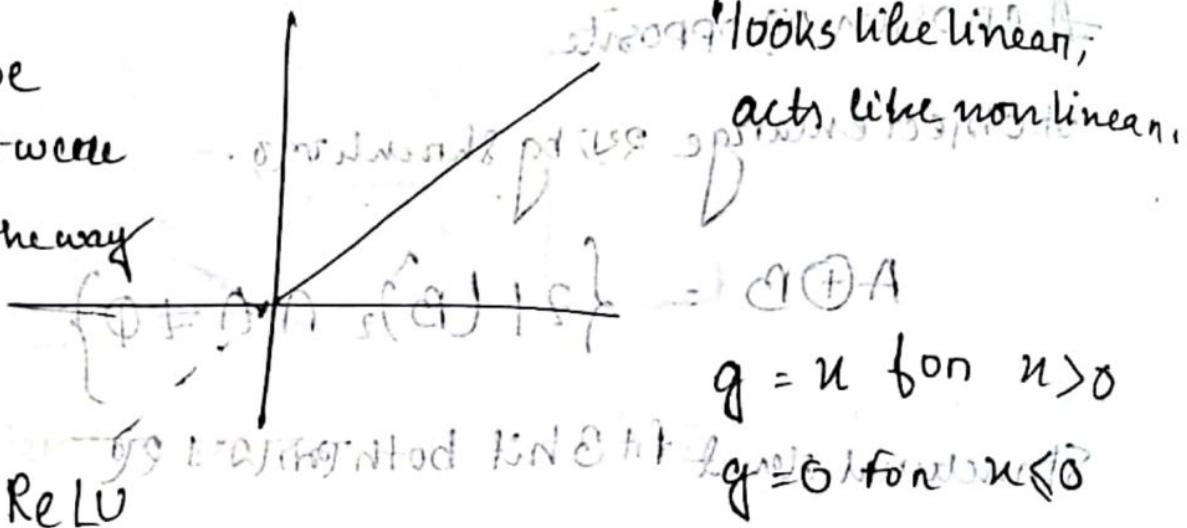
ReLU

dropout - randomly
 dropping some neurons
 of neural network
Hyperbolic tangent function

#1 linear # backpropagation

notable

It would be
linear, if it were
linear all the way



advantage

- positive → linear रूपान्तर
- loss function → $u^2/2$,
- so derivative weight
update आसानी से होती है
(for positives)

- neural networks
fast converge होती है
- dropout ज़हारा ज़हारा होता है
- negative value drop
माफ़िज़ी मार्फ़िज़ी होती है

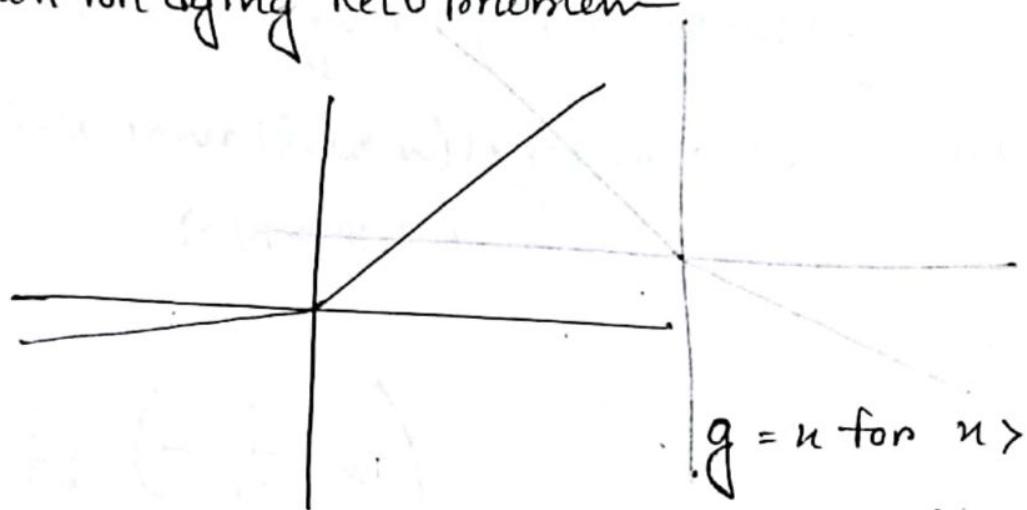
- computational
load ज़हारा ज़हारा होता है
- calculation easy

dying ReLU problem:

disadvantage

- negative → zero यूटी-एफ
- positive → linear
- negative value → ज़हारा ज़हारा होता है / derivative 0 रख दिया जाता है

Solution for dying ReLU problem



$$g = n \text{ for } n > 0$$

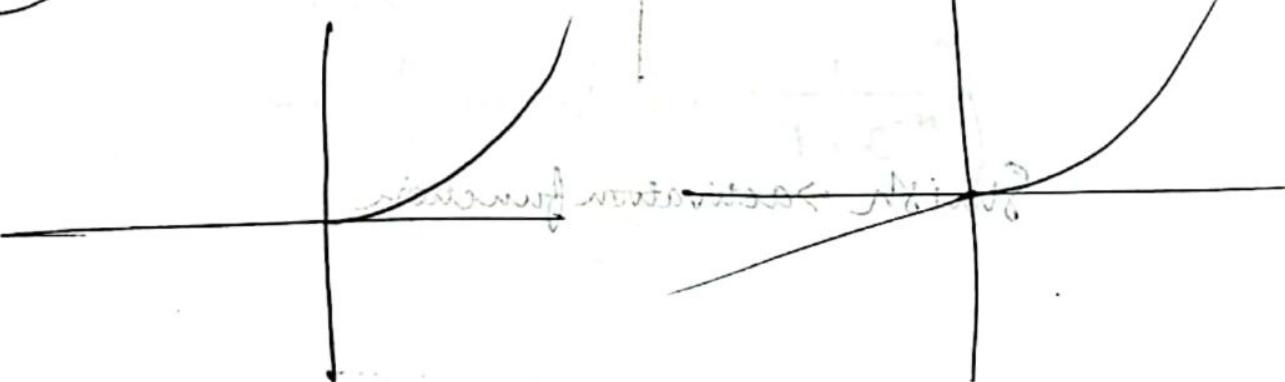
leaky ReLU activation functions

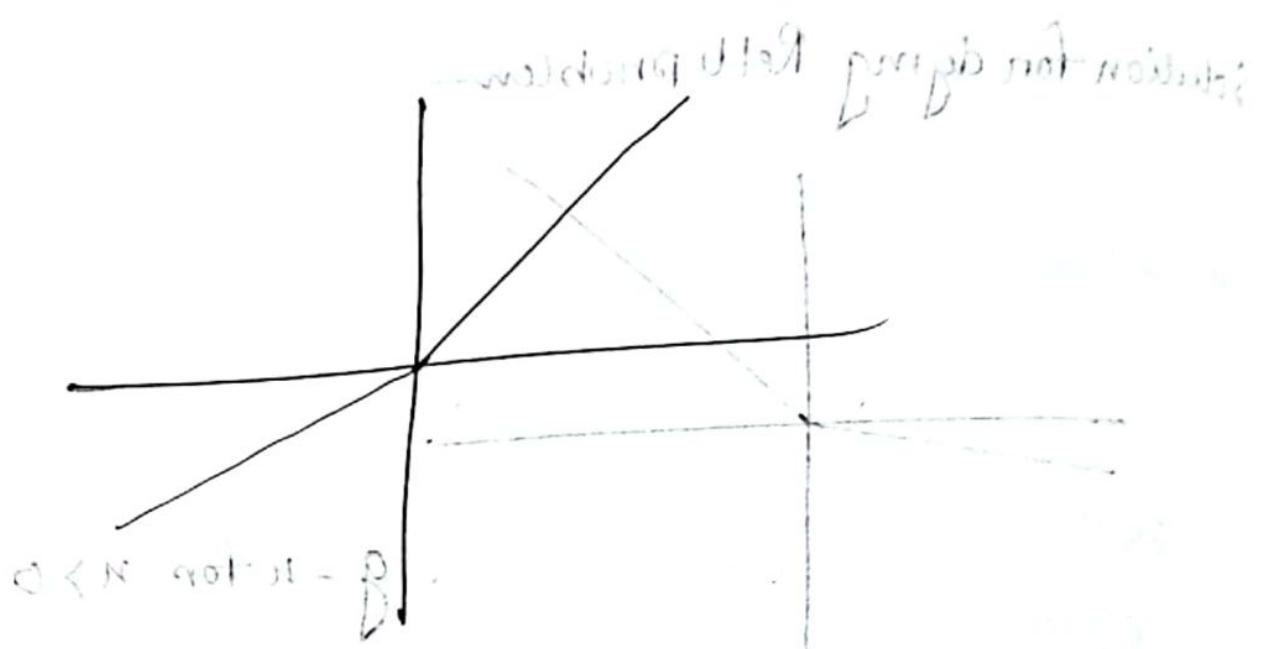
$$g = 0.1n \text{ for } n \leq 0$$

which activation function is suitable depends on data

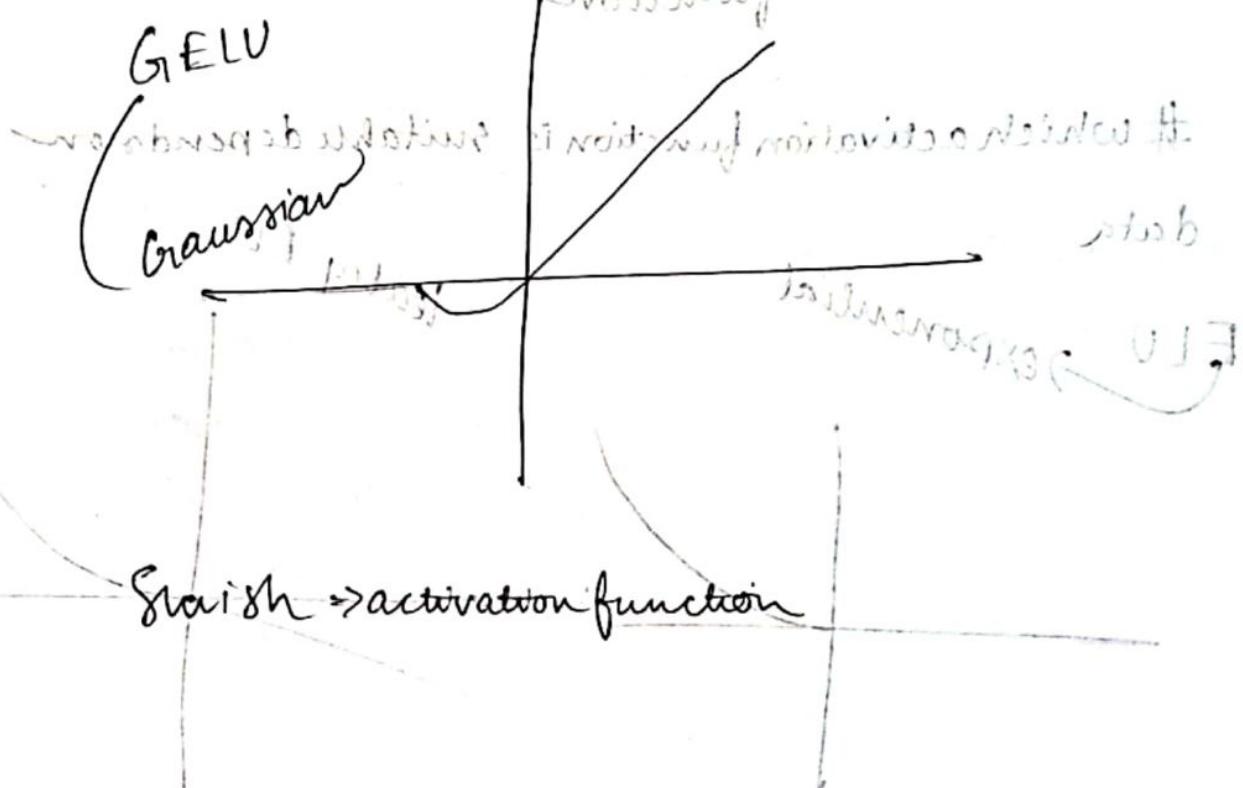
ELU exponential

leaky ELU





Parameterized ELU



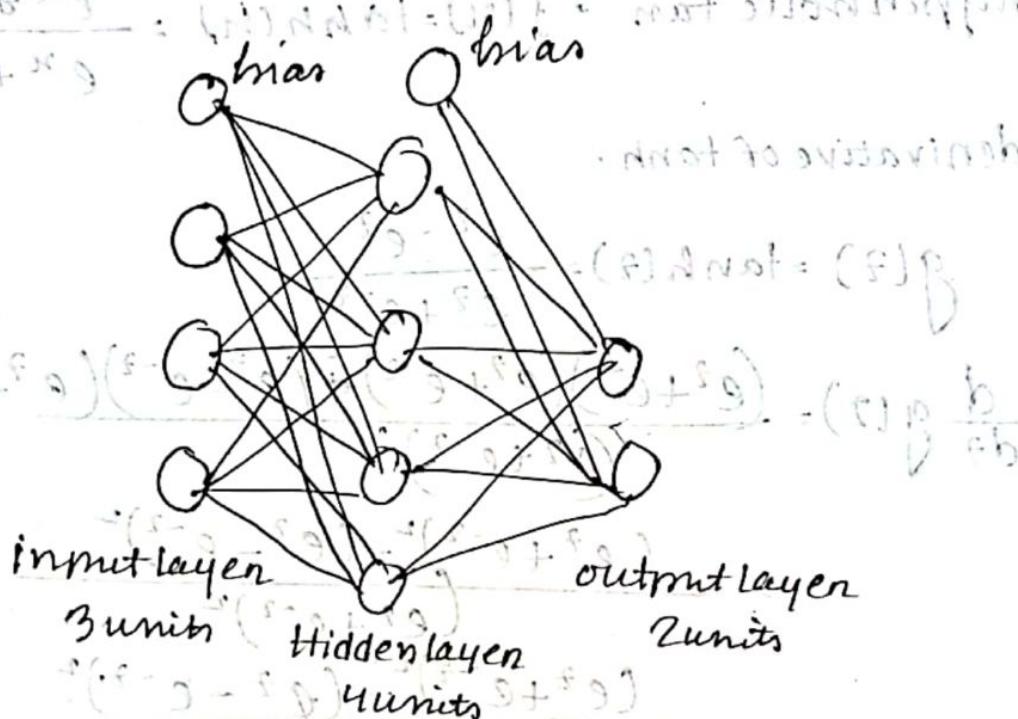
Hyperparameters and parameters

Parameters (w, b) $w \in \mathbb{R}^n$ {weights} $b \in \mathbb{R}$ {bias}

weight, bias set at training time. A set of 10^2 or 10^3 can change
 \Rightarrow $w \in \mathbb{R}^{n \times m}$ $b \in \mathbb{R}^m$

\Rightarrow not set by users $w \in \mathbb{R}^{n \times m}, b \in \mathbb{R}^m$ {constant w/ comp.}

\Rightarrow learned from data



Input layer \Rightarrow 3 node, hidden layer \Rightarrow 4 node
 \Rightarrow 12 connections. \Rightarrow connection \Rightarrow weight \Rightarrow 12
 \Rightarrow 2 layers \Rightarrow total weight = 3×4

Hidden layer \Rightarrow each node output layer \Rightarrow 2 node
 \Rightarrow 8 connections, so for 4 hidden layer = 4×2

one bias for 4 nodes in hidden layer = 4×1

one bias for 2 nodes in output layer = 2×1

$$\text{total weights} = 12 + 8 + 4 + 2 = 26$$

Hyperparameters

⇒ External to model

⇒ can not be estimated from data, set before training

⇒ determines network structure

Examples: dropout, # of hidden layers, weight initialization,

Active function

- Epoch, Batch size, optimizer, learning rate

Tuning techniques

⇒ Randomly chooses search

Randomly chooses all possible combinations of hyperparameters

⇒ Grid Search

Tests all possible combination of hyperparameters of given machine learning algorithms.

⇒ Bayesian optimiser

Based upon Bayes' rule and considers previously known knowledge to help narrow down search space of good hyperparameter combinations

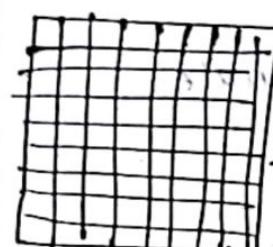
⇒ Gradient-based

A methodology to optimise several hyperparameters,

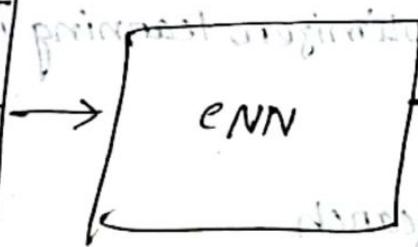
Based on the computation of the gradient of machine learning model selection criterion w.r.t hyperparameters, prominent among them is gradient backpropagation for NN.

CNN

convolutional Neural Network



200x200x3



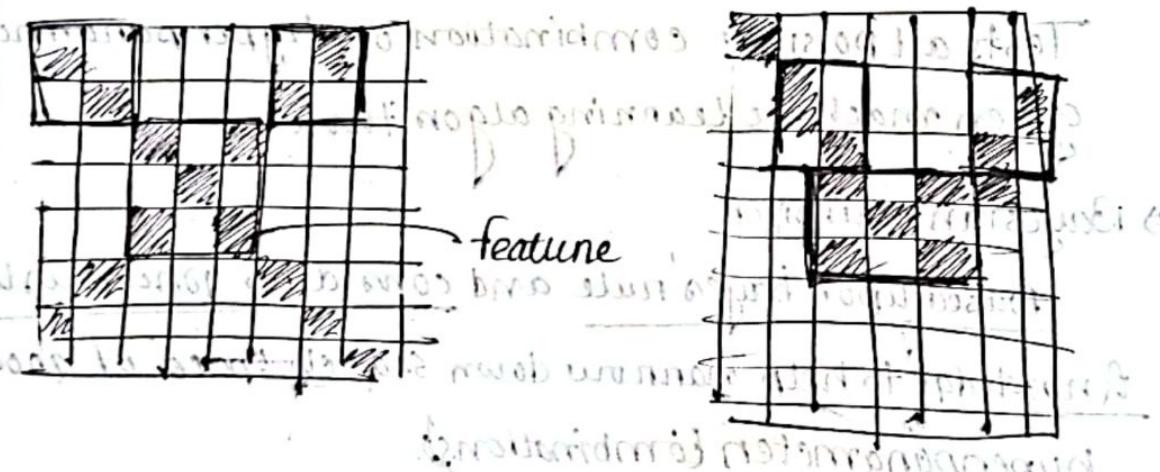
$X \in \mathbb{R}^0$

input-output

feature extraction

not meaningful

→ Feature engineering as follows:-



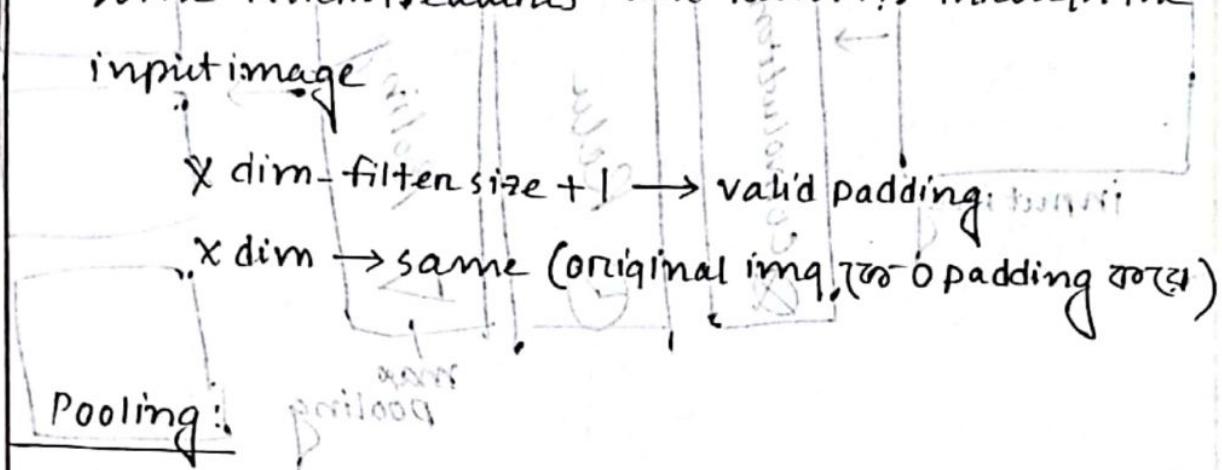
based on forward pass

CNN learns which feature signature of pixels best fits for

Fitting

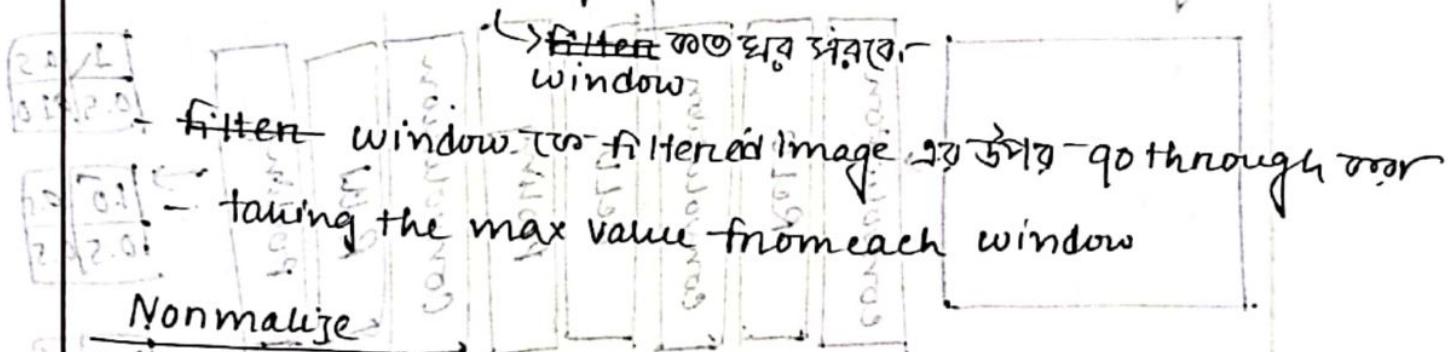
b33002 top 601301

This model does feature engineering by itself. So it finds some 'filters/features' and runs it through the input image.



Shrinking the image stack

- fit specific stride pick window



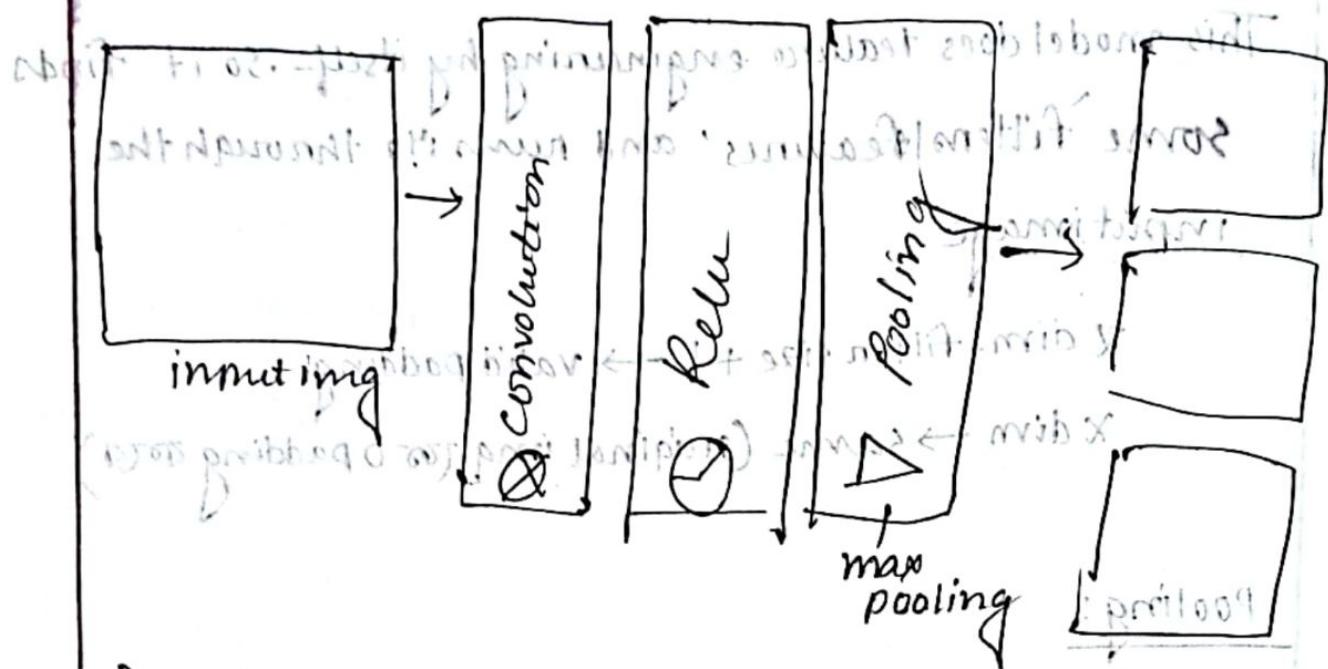
Nonnormalize

Negative value 0 (ReLU)

Negative value 0 or positive value 0 discard

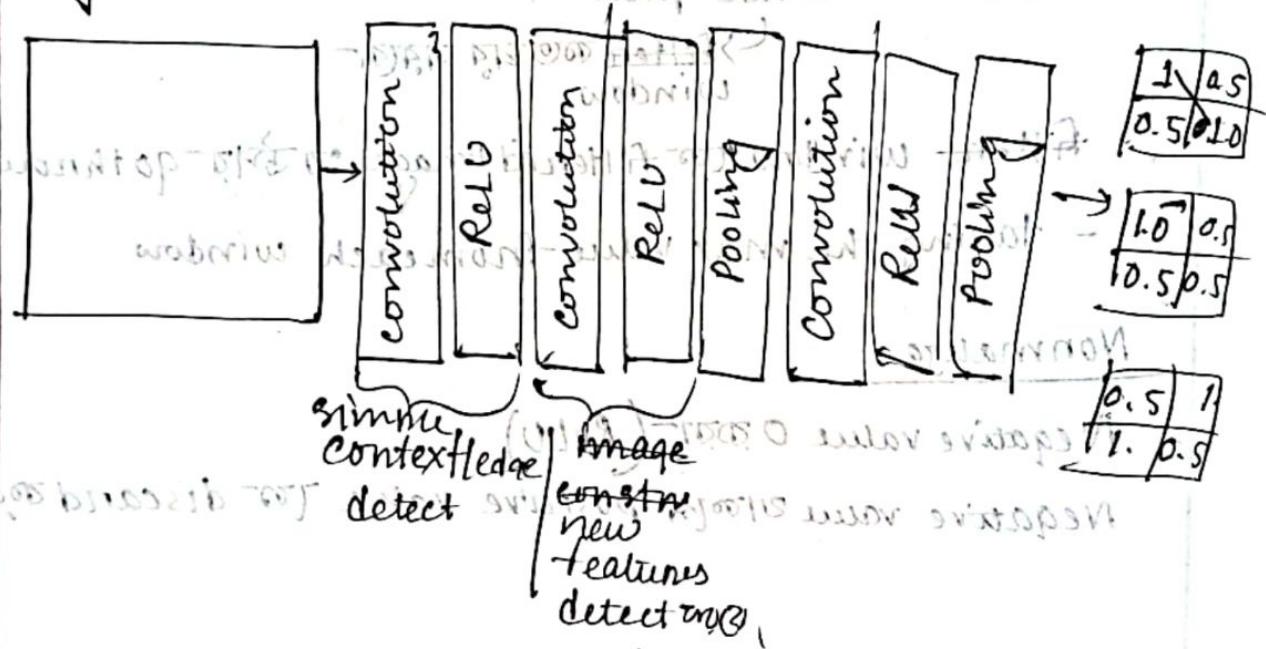
constant
bias

Layers get stacked

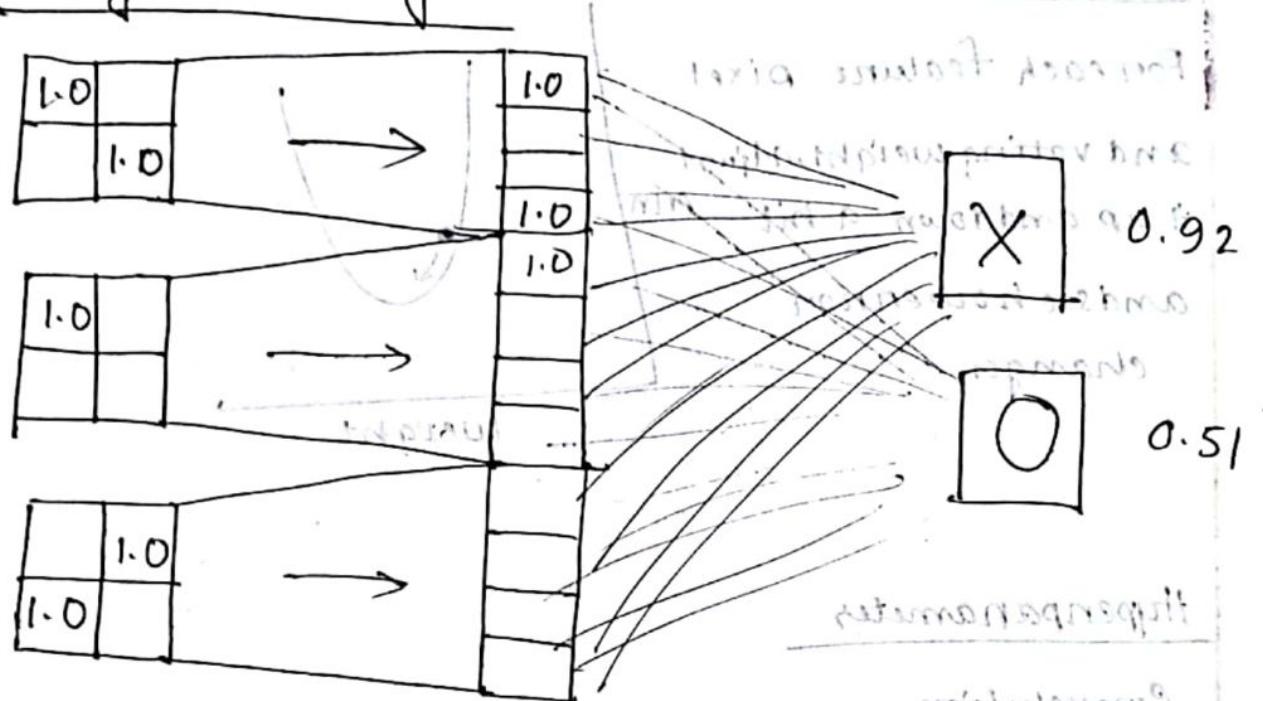


Deep stacking

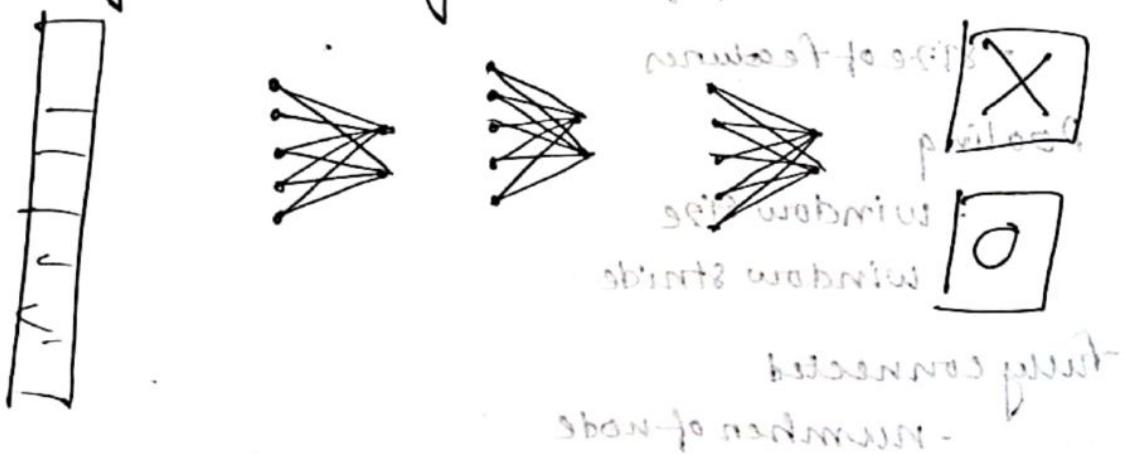
Layers can be repeated several times



Fully connected layer



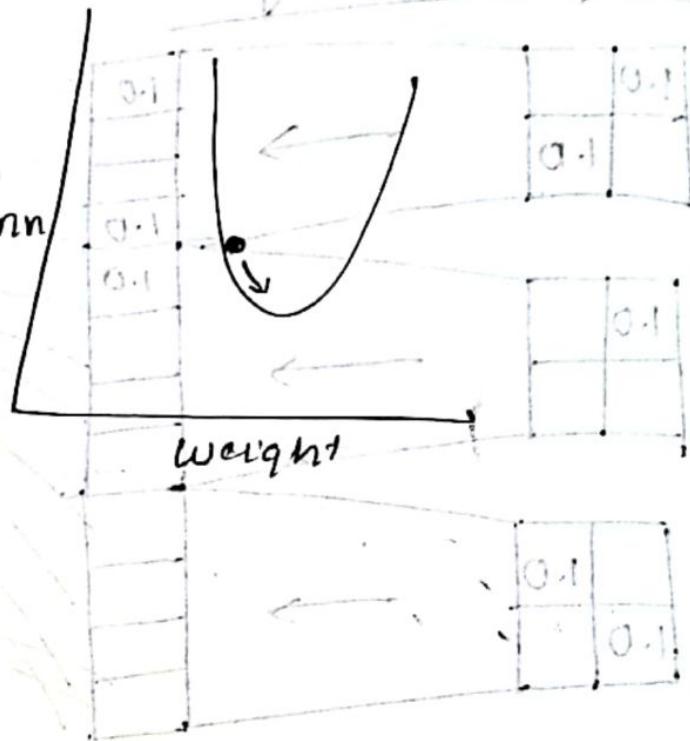
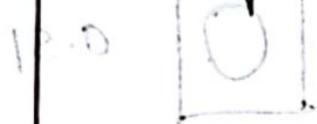
Fully connected layers can be stacked too.



Backpropagation = filter of weight learn tool.

Gradient descent

For each feature pixel and voting weight, adjust it up and down a bit and see how error changes.



Hyperparameters

Convolution

- number of features

- size of features

Pooling

- window size

- window stride

fully connected

- number of nodes

- output nodes



Limitation

What are the responsibilities of the manager?

- only captures local 'spatial' pattern in data; If the data can not be made to look like an image, conv net is less useful.

Rule of thumbs

Abbildung

If $\sum_{i=1}^n \beta_i = 1 + \Sigma - \Omega$ then good linearizations

If shuffling columns does not affect data, there won't be any loss of information.

160559

sovereign masters

No. 1

8

[[01-01-2013]]

Flute
(spinetone)

1

(1.8.8 version)

(Answers) English

e:महाराष्ट्र

1. ~~ANSWER~~

(mid. L)

(1.2.3.2028)

JOURNAL

Dr. R. M. S. Khan et al.

01. ~~Reasons for punishment~~

Q: What is your name?

Slide: Parameter calculation of NN

Parameter calculation in NN

non trainable

par = $\text{tf.Variable}(\text{tf.ones}([10, 10, 1]))$ image height
image width
channel

my_input = Input(shape=(10, 10, 1))

my_output = Convolution 2D(1, 3)(my_input)

height/width

\rightarrow 3x3 filter to start

output will become $10 - 3 + 1 = 8 \times 8$

my_output = Activation('sigmoid')(my_output)

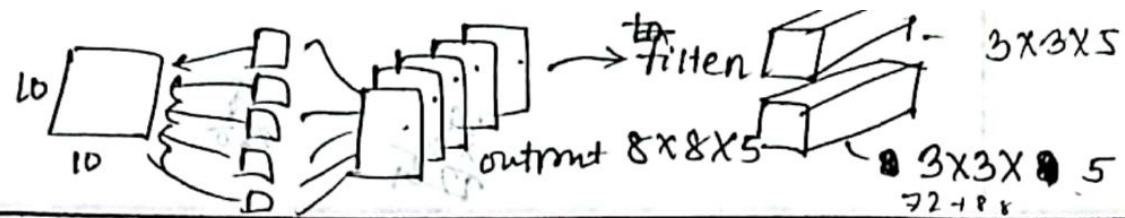
Model summary

Layer	Output shape	Param
Input-17 (Input layer)	[None, 10, 10, 1]	0
Conv2D-19 (conv2D)	(None, 8, 8, 1)	10 (3x3 filter = 9 values + 1 bias)
activation()	(None, 8, 8, 1)	0

Total param: 10

trainable param = 10

non trainable = 0



my input = Input(Shape = (10, 10, 1)) \rightarrow 5 3x3 filter

my output = Convolution2D(5,3)(my_input) kernel

my output: Activation('sigmoid')(my_output)

`myoutput = convolution2D(2,3)(myinput)`

my output = Ac - .

744 : $8 \times 82 = 656$ जिसका भाग 8 है।

$$8 \times 8 = 64$$

layer outfit shape paran

Input-19 [none,10,10,1] 0

conv2d-22 [none, 8, 8, 5] 50

activation 22 (none, 8, 8, 5) 0

com 2d-23 (none, 6, 6, 2) 46X

activation (none, 6, 6, 2) *mislabeling* (1)

www.english-test.net

[(2,001,35000)] ~~✓-kappa~~

$\{(3, 30, 35000)\}$ (333333) 201m

851 (2x2) ~~क्षेत्रीय~~ अंतर्राष्ट्रीय

g\spit-to see
-11256 - 11

ପ୍ରକାଶକ ଓ ମୁଦ୍ରଣ କମିଶନ୍

$25 = 5 \times 5$: missing factor

— — — — — *mostar*

78 (ESTIMATE) 2-Hydroxy

layer	Output shape	Param
Input-20	[None, 100, 100, 3]	0
conv2d-24	(None, 98, 98, 8) 8x[3x3x3]+1	224
activation	(None, 98, 98, 8)	0
conv2d-25	(None, 96, 96, 1) 3x3x8 + 1	73
activation	(None, 96, 96, 1)	0
1D convolution	(None, 96, 5)	480
layer	Output shape	Param
input-25	[(None, 100, 5)]	0
conv4c	[(None, 98, 8)]	128
	-filter (3x5) no of filter 8 tot param is 1 filter: $(3 \times 5) + 1 = 16$ total param = $8 \times 16 = 128$	
activation		
conv 1d-5	[(None, 96, 1)] $((3 \times 8) + 1)$	25

was ^(P.M.) (1 + 10%)

Number of parameters in conv layer would be :

$\text{dim} \Theta = ((m \times n \times d) + 1) \times k$

((shape of the width of the filter x shape of height of filter) x

- number of filters(channel) in previous layer } + 1) X number of filters.

0 ————— [Blank area]

Output shape

$$\text{Slope} = \frac{y_2 - y_1}{x_2 - x_1}$$

$$\text{Floor} \left(\frac{(w - F + 2P)}{s} + 1 \right)$$

0 ————— [PIVOT POINT]

0 _____ (F.F.E.S.J.-] Isogeny

० ————— [४३८५४] ————— ० - विजयगढ़ी

115555-10 ——— [126841]X115 ——— [1-1213]

848.806 =

Dubois II → (1-1218) → 0

0 ————— (4-218) ————— 150Tn

111555-13 (213A132) — [2148158] (2148158)

پاکستان

5 [1944-1] 1944-1945

3081 10181-3 31-1000000

$$\begin{array}{r} 26 \\ \times 16 \\ \hline 156 \\ 26 \quad 0 \\ \hline 560 \end{array}$$

$$\begin{array}{r} 28 \\ \times 28 \\ \hline 224 \\ 560 \end{array}$$

$$(3 \times 3) + 1$$

$$\begin{array}{r} (4 \times 4) \\ \times 5 \\ \hline 6 \times 5 \\ 26 \times 16 \end{array}$$

Example

layer

Conv2d

$$[-1, 16, (28, 28)] \xrightarrow{3 \times 3 \times 16} 16 \times (5 \times 5) + 1 \rightarrow 416$$

dropout

$$\times (1 + \text{random value}) \cdot [-1, 16, 28, 28] \rightarrow \text{After dropping } 0$$

ReLU

$$[-1, 16, 14, 14] \longrightarrow 0$$

$$\begin{array}{r} 28, \\ 14, \\ \hline 8, 4 \\ 16 \end{array}$$

Maxpool

$$[-1, 32, 14, 14] \xrightarrow{\text{max pooling}} 32 \times (5 \times 5 \times 16) + 1 \rightarrow 12832$$

$$16 \times (15 \times 15) + 1$$

Dropout-6

$$\left(\frac{[-1, 32, 14, 14]}{2} \right) \xrightarrow{+ (98+7-6) \text{ nodes}} 0$$

ReLU

$$[-1, 32, 14, 14] \longrightarrow 0$$

$$\text{Maxpool } [-1, 32, 7, 7] \longrightarrow 0$$

$$\text{flatten-9} \longrightarrow [-1, 1568] \longrightarrow 0$$

$$\text{Linear-10} \longrightarrow [-1, 512] \longrightarrow (1568 + 1) \times 512 = 803,328$$

$$\text{Dropout-11} \longrightarrow [-1, 512] \longrightarrow 0$$

$$\text{ReLU} \longrightarrow [-1, 512] \longrightarrow 0$$

$$\text{Linear-13} \longrightarrow [-1, 128] \longrightarrow (513 \times 128) = 65664$$

$$\text{Dropout-14} \longrightarrow [-1, 128]$$

$$\text{ReLU-15} \longrightarrow [-1, 128]$$

$$\text{Linear-16} \longrightarrow [-1, 10]$$

$$1290$$

$$\text{output shape} = \left(\frac{W_H - f + 2P}{S} + 1 \right)$$

Example: different operational point of view

1) input shape = 1024

arrow to point

Layer

1024

1024

0

FC(256)

256

$(1024 \times 256) + 256$

dropout

0

$= 262400$

FC(128)

128

$(256 \times 128) =$

fanin

128

32768

maxpool(3)

$\text{floor}\left(\frac{(128 - 3 + 2 \times 0)}{3}\right) + 1$

42

~~120~~

FC(3)

3

~~42~~

output

3

0

295294

Natural language Processing and word embedding

Bag of words

We create a vocabulary with the words that are present in our dataset.

Review). This movie is very scary and long and is slow.

Review 2: This movie is not scary and is slow.
Review 2: This movie is not scary and is slow.

Review 2: This movie is not scary and is slow. It's a trashy and dumb

Review 3: This movie is spooky and good

This movie is very scary and long not slow.

Devon 2 B02 prehistoric sp. 0 0 0

Review 3 1 1 1 0 0 0 0 0

~~word by reference to one hot encoding~~

25 ~~would be presence of one non-ecological~~
..., # presence calculated 27.0

#context এর
#horizon এর
#obj এর

vocabulary ମାତ୍ରାଙ୍କ ଶବ୍ଦଗୀ.

BOW example

word not mentioned

D-1: It was the best of times

presented

D-2: It was the worst of times

absence of word

D-3: It is the time of stupidity

absence of word

D-4: It is the age of foolishness

absence of word

It was the best of times, worst of times, is age foolishness,

D-1 1 1 1 1 1 1 0 0 0 0 0

D-2 0 0 0 0 0 0 0 0 0 0 0

D-3 0 0 0 0 0 0 0 0 0 0 0

D-4 1 0 1 0 1 0 0 0 0 1 1 1

vocabulary size 10 words; but each document has only few of them.

The result is a vector with lots of zeros, called sparse representation.

Solution 1 for BOW

Stemming Vs Lemmatization

significance word

Stemming

→ removes last few characters from a word often leading to incorrect meaning and spelling

spelling

Caring → car

0 0 0 0 0

0 0 0 0 0

1 1 1 0 0

1 1 1 0 0

need to

consider the context and convert the word to

Lemmatization

→ considers the context and converts the word to meaningful base form

bare form

caring → care

0 1 1 0 1 1 1 1 1

→ computationally expensive as it involves lookups table

0 1 0 1 0 1 0 1 1

0 1 0 1 0 1 0 1 1

Solution 2

using ngram

for example, using trigram

It was the best of time

It was the best of time

Term-frequency (TF)

$$TF_{t,d} = \frac{n_{t,d}}{\text{Total number of terms in document } d}$$

D-1: It was the best of time worst of times

D-2: It was the worst of times

D-3: It is the time of stupidity

It was the best of time worst of stupidities.

$$D-1 \quad 1 \quad 1 \quad 2 \quad 1 \quad 1 \quad 1 \quad 0 = \frac{6}{12} = 0.5 \text{ (7)}$$

$$D-2 \quad 1 \quad 1 \quad 1 \quad 0 \quad 1 \quad 0 = \frac{4}{12} = 0.33 \text{ (8)}$$

$$D-3 \quad 1 \quad 0 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 = \frac{5}{12} = 0.416 \text{ (6)}$$

Calculating TF

$$D-1 \quad \frac{1}{7} \quad \frac{1}{7} \quad \frac{2}{7} \quad \frac{1}{7} \quad \frac{1}{7} \quad \frac{1}{7} \quad 0 = \frac{5}{14} = 0.357 \text{ (7)}$$

$$D-2 \quad \frac{1}{8} \quad \frac{1}{8} \quad \frac{1}{8} \quad 0 \quad \frac{1}{8} \quad \frac{1}{8} \quad \frac{1}{8} = \frac{6}{16} = 0.375 \text{ (8)}$$

$$D-3 \quad \frac{1}{8} \quad 0 \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad \frac{1}{6} \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad \frac{1}{6}$$

$$D-3 \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad 0 \quad 0 \quad \frac{1}{6} \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad 0 \quad \frac{1}{6} \quad \frac{1}{6}$$

TF-IDF

probabilistic model

$$W_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

non-informative words

TF-IDF for each document

2 rows for probabilities of the probabilities

for two words Doc-1 prob Doc-2 prob Doc-3 prob

	0	0	0	0	0	0	0	0
It	0	0	0	0	0	0	0	0
was	0.026	0.03	0	0	0	0	0	0
the	0	0	0	0	0	0	0	0
best	0.068	0	0	0	0	0	0	0
time	0	0	0	0	0	0	0	0
worst	0	0	0.08	0	0	0	0	0
stupidity	0	0	0	0	0	0	0	0
is	0	0	0	0	0.08	0.48	0.08	0

10230

info gain

Word embedding

$$\underbrace{W \times D}_{\text{matrix}} \cdot \underbrace{D \times H}_{\text{matrix}} = W \times H$$

Word representation:

vocabulary matrix containing all words
each word one hot encoding representation
vector length = vocabulary size

1 hot representation.

size of vocabulary $[V] = 10,000$
 $V = [a, aar, \dots, \text{food}, \text{UNK}]$

Man (5391)	woman 9853	computer	food
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
1	1	0	0
0	0	0	0
0	0	0	0

meaning नाही -

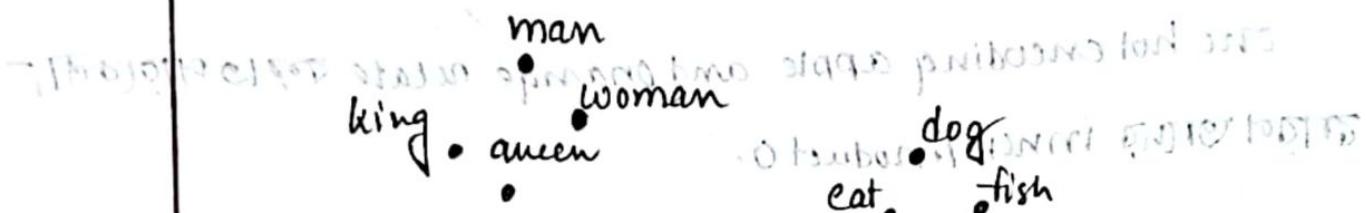
I want a glam of orange

I want a glass of apple

Featureized representation

	Man (5391)	Woman (9853)	King (40144503)
Gender	-1	1	-0.95
Royal	0.01	0.02	0.03
Age	0.03	0.02	0.70
Food	0.09	0.01	0.02

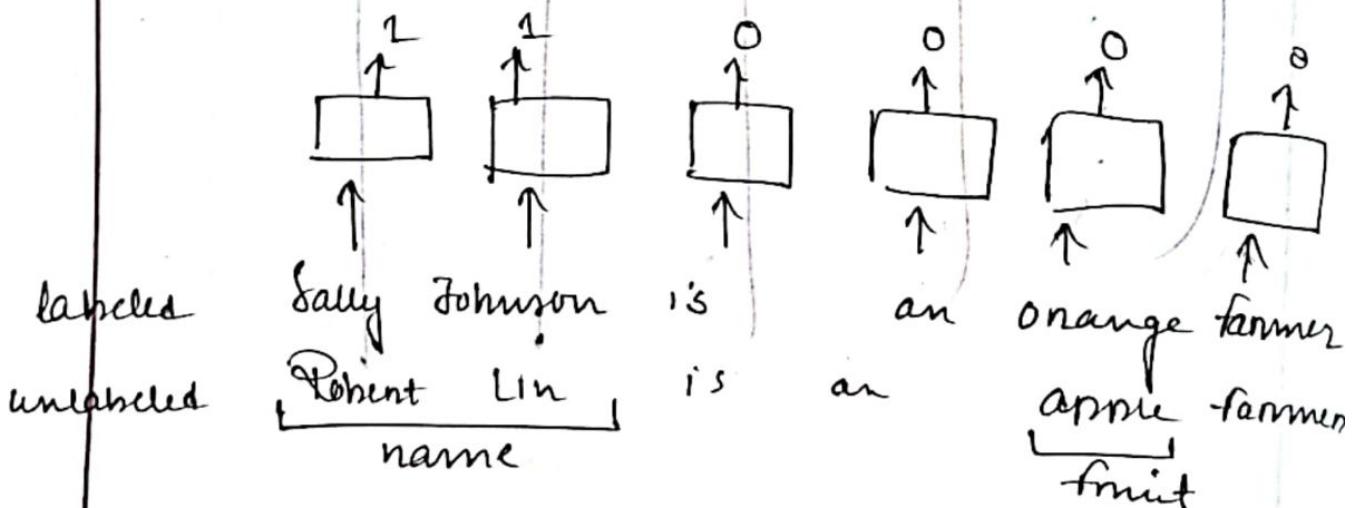
Visualizing word embedding



TSNE algorithm to visualize word embedding in 2D
 $(300D) \rightarrow 2D$

similar words are together

Named entity recognition problem



Slide: Natural language processing and word embeddings

Transfer learning and word embeddings

↪ transfer learning & weight update वर्गा पाए जा।

⇒ will re-use already trained models

⇒ आमाने task अनुसारी smaller labeled training dataset &
trained embedding apply कराए।

⇒ continue to finetune यहां आमाने training dataset आजत तक
देखा।

Advantage

→ 10,000 one hot sparse >> 300 features vector (dense)

⇒ helpful when we have smaller dataset

$$\text{Ques.} \quad \frac{\text{one hot}}{\text{dense}} = (\text{one hot}) \text{ padding zeros}$$

Analogies using word vectors

man : woman :: king : queen

$e(\text{man}) - e(\text{woman}) \approx e(\text{king}) - e(\text{queen})$

$\text{argmax}_w \text{sim}(e(w), e(\text{king}) - e(\text{man}) + e(\text{woman}))$

Similarity function

Euclidean distance.

$$\sqrt{\sum_{i=1}^n (u_i - v_i)^2}$$

Cosine similarity

$$\text{cosine similarity } (u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2} = \cos \theta$$

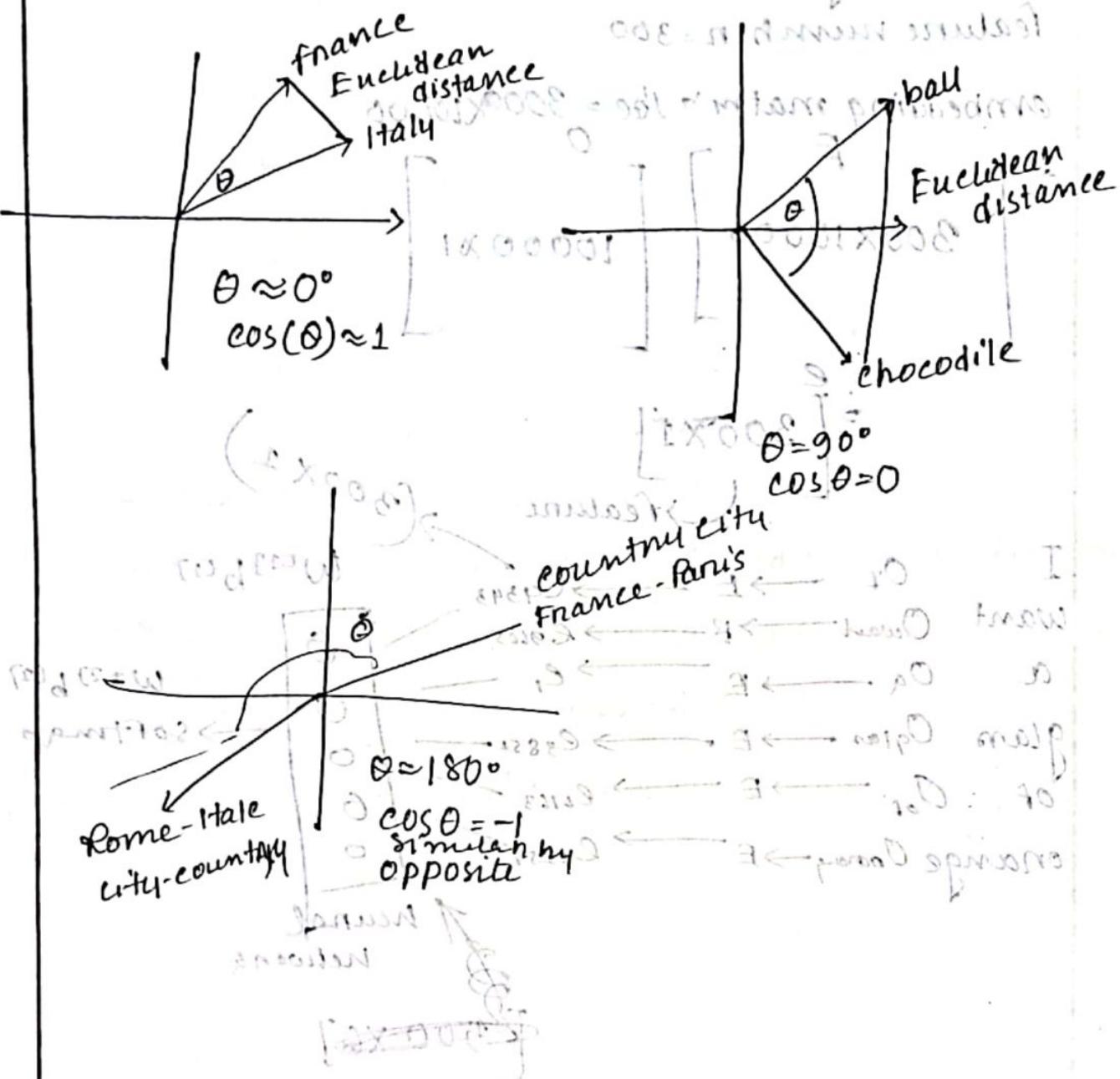
we can
distance

[-1, 1]

$$\text{cosine similarity} = \frac{\mathbf{U} \cdot \mathbf{V}}{\|\mathbf{U}\|_2 \|\mathbf{V}\|_2}$$

$\text{vector } \rightarrow \cos \theta = \text{similarity}$

☞ $\mathbf{U} \cdot \mathbf{V}$ dot product

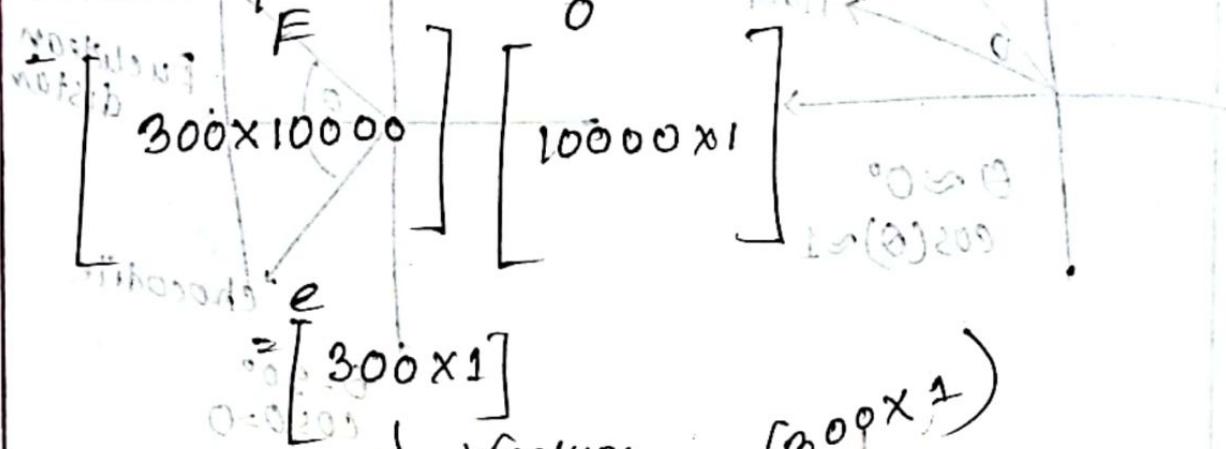


Embedding matrix

Vocabulary size = 100000 (but for V.U. 300)

feature number = 300

embedding matrix size = 300×100000



I

$o_1 \rightarrow E \rightarrow e_{4343}$

want

$o_{\text{want}} \rightarrow E \rightarrow e_{9105}$

a

$o_a \rightarrow E \rightarrow e_1$

glan

$o_{\text{glan}} \rightarrow E \rightarrow e_{3852}$

of

$o_{\text{of}} \rightarrow E \rightarrow e_{6163}$

orange

$o_{\text{orange}} \rightarrow E \rightarrow e_{6253}$

$W^{(1)} b^{(1)}$

$W^{(2)} b^{(2)}$

softmax



neural
network

300×6

Word2Vec: Skip gram

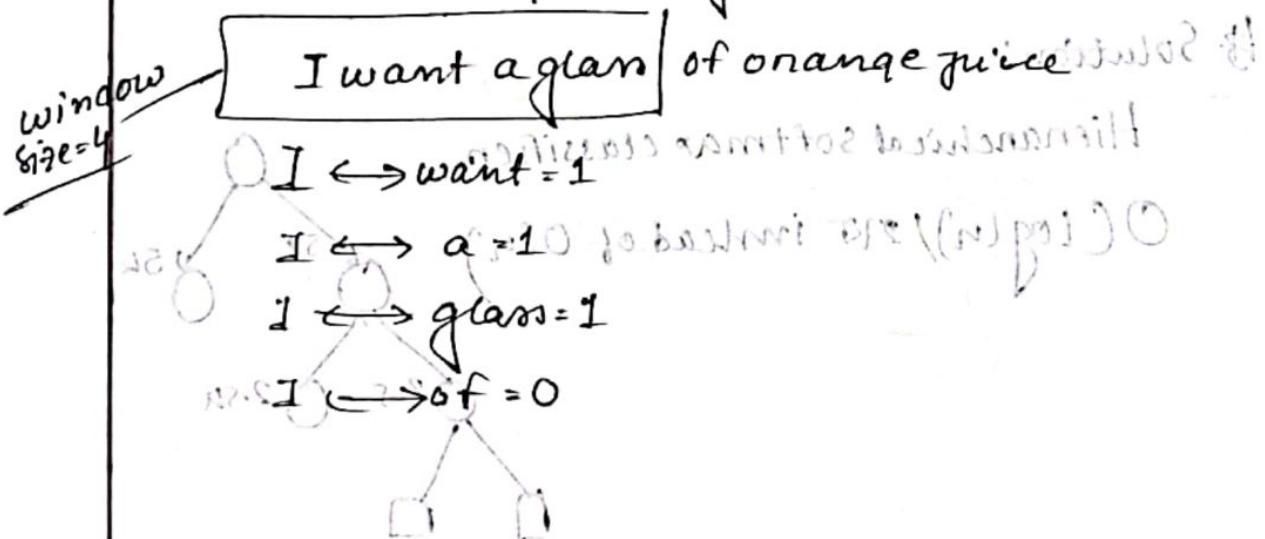
I want a glass of orange juice to go along with my cereal

Context	Target	Window
Orange juice	orange	+1
orange glass	juice	-2
orange	my	+6

→ Randomly pick context word

→ target is chosen randomly based on window of specific size.

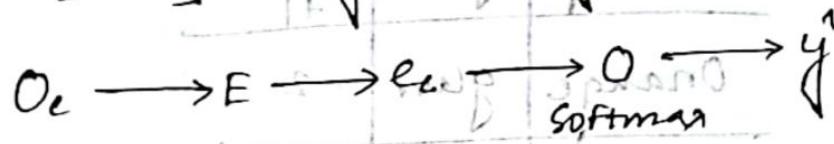
→ we have converted the problem into a supervised problem to predict the target word given the context word.



Vocabulary size = 10,000

context word c ('orange') target word t ('juice')

$x \rightarrow y$
softmax layer: $\text{fix} \rightarrow \text{softmax}$
context $c [6257]$ ('orange') \rightarrow target $t [4834]$ ('juice')



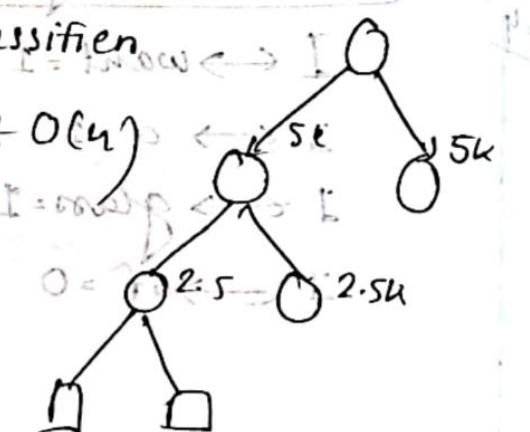
softmax $p(t|c) = \frac{e^{O_t^T e_c}}{\sum_{t=1}^{10,000} e^{O_t^T e_c}}$
know existing probabilities \rightarrow very complex

Loss function $L(y, \hat{y}) = -\sum_{i=1}^{10,000} y_i \log \hat{y}_i$
know existing probabilities \rightarrow know existing probabilities \rightarrow very complex

Solution: \rightarrow $\boxed{\text{map to know}}$

Hierarchical softmax classifier

$O(\log(n))$ instead of $O(n)$



Negative Sampling

Random sampling হলো নম্বৰ \Rightarrow context word ও sample \Rightarrow

অপ্রয়োজনীয় word (such as is, the) আসতে পাই, যাবার

frequent words অপ্রয়োজনীয় word (রাখা পড়া মাছ)

\Rightarrow sample করিব তারা অসমুচ্ছ frequent word ও
ৰাখিব তারা কীভু

এটা skipgram model

\Rightarrow skipgram এ প্রুম্ভ ফিল্ট - softmax \Rightarrow high computational
power দ্বারা ফিল্ট

Negative sampling এ concept হলো - for a context word
আমাতো নেগেটিভ sample \Rightarrow 1 টি positive sample করা \sim

$$\begin{aligned} P(Y=s | e, t) \\ = \sigma(\theta_t^T e_c) \end{aligned}$$

context word	target
orange	juice
orange	wing
orange	book
orange	the
orange	of
e	y

How do we solve the softmax complexity?

अगला task द्वारा binary logistic regression को कैसे
प्रोग्राम को बदलें तथा प्रोग्राम को बदलने का
प्रयोग करके 10,000 वाला software का उपयोग कैसे करें।

→ How to choose negative sample which are not CTC

$$P(\omega_i) = \frac{f(\omega_i)^{3/4}}{\sum_{j=1}^{10000} f(\omega_j)^{3/4}}$$

frequency of a particular word in a corpus

Project BRCDN 201409

когда тексторы работают в группах с различными методами

— କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା କିମ୍ବା

$$(s^T, \theta)^T =$$

Practice

- D1 - A woman finds a pot of treasure on the road while she is returning from work
- (1) D2 - Delighted with her luck, she decides to keep it.
- (2) D3 - However, enthusiasm refused to fade away

$$IDF('luck') = \log \frac{3}{1} = 0.48$$

$$IDF('enthusiasm') = \log \frac{3}{1} = 0.48$$

$$TF('to') = \log \frac{3}{2} = 0.176$$

$$TF('luck') = \frac{1}{18}$$

$$TF('enthusiasm') = \frac{1}{7}$$

$$IF('to') = \frac{1}{7}$$

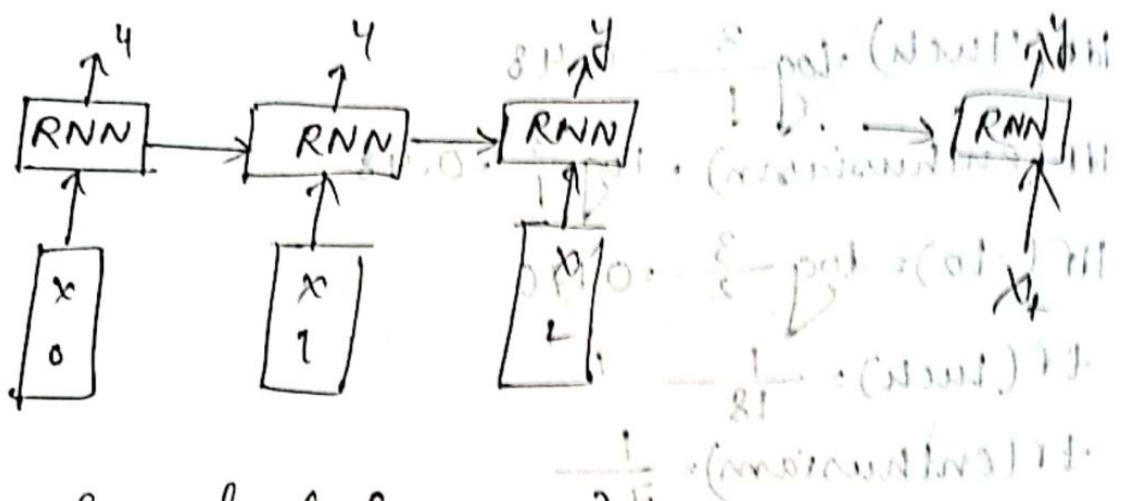
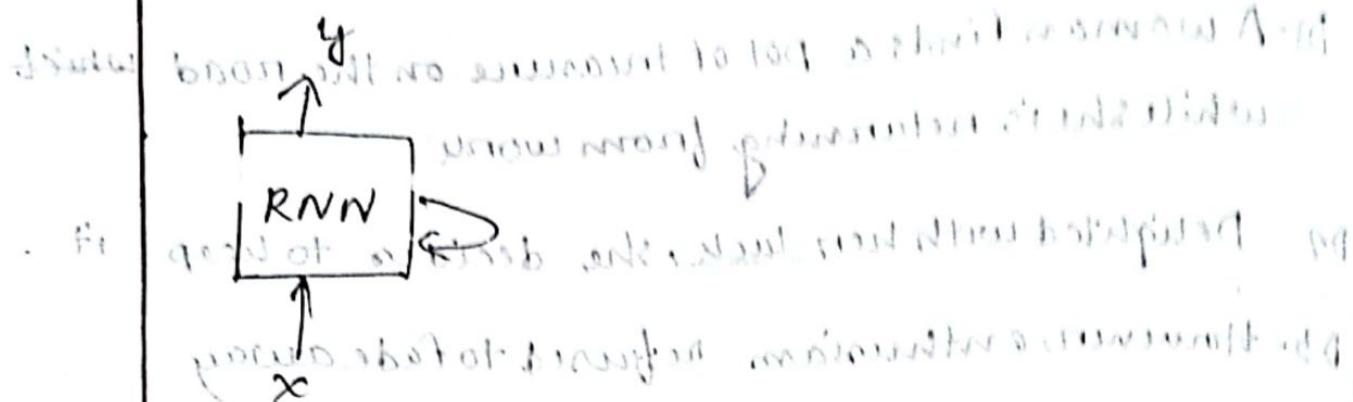
$$IF-IDF_{\frac{1}{18}} \times 0.48 = 0.06 - 0.027$$

$$TF-IDF_{\frac{1}{7}} \times 0.48 = 0.08$$

$$(TF-IDF_{\frac{1}{7}} \times 0.176) + 0.025 = 0.025$$

$$IA_B = 0.025$$

Recurrent Neural Networks



$$h_t = f_w(h_{t-1}, x_t)$$

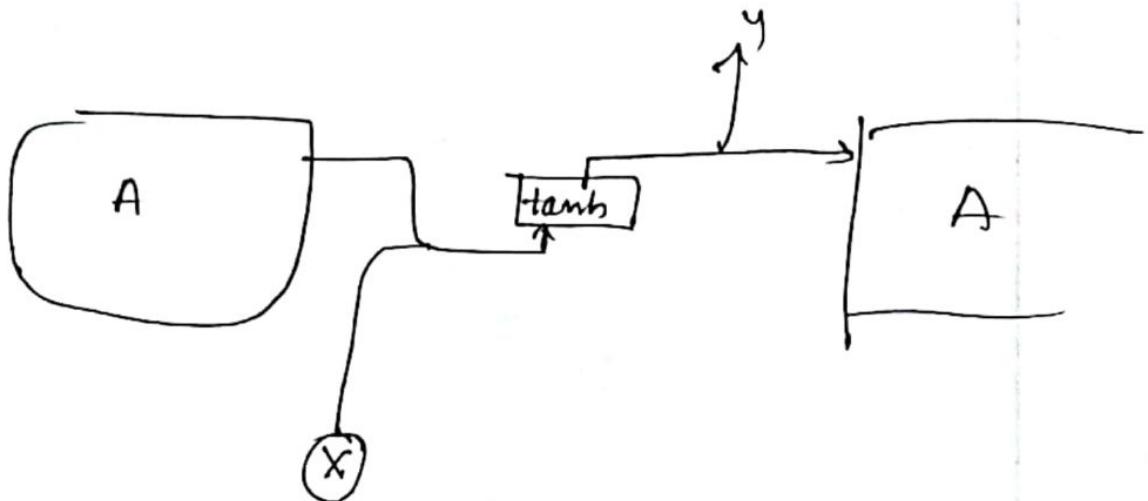
(old state) \downarrow (input vector at time t) \uparrow
 new state function

F.S. \rightarrow S.O. \rightarrow S.P.O.X \rightarrow S.T.E.P.T \rightarrow H.A.T.
 80.0 \rightarrow 80.0 \times $\frac{1}{\sqrt{2}}$ \rightarrow 80.0

$$y_t = w_y h_t \approx \tanh(W_h h_{t-1}^{(0)}, W_u u^{(0)})$$

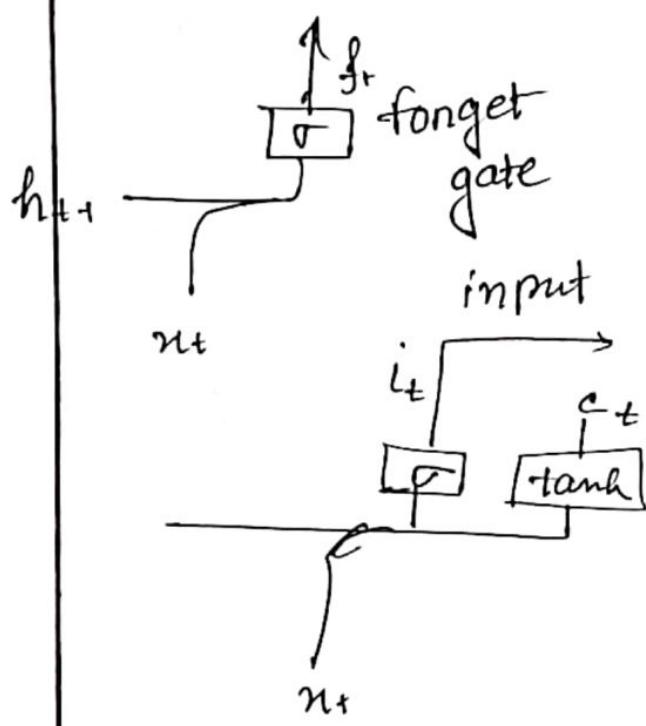
$$y_t = w_y h_t$$

$$\frac{12 - 3}{2} + 1$$



LSTM

c_{t-1} \circledast cell state



$$f_t = \sigma(w_f [h_{t-1}, n_t] + b_f)$$

$$i_t = \sigma(w_i [h_{t-1}, n_t] + b_i)$$

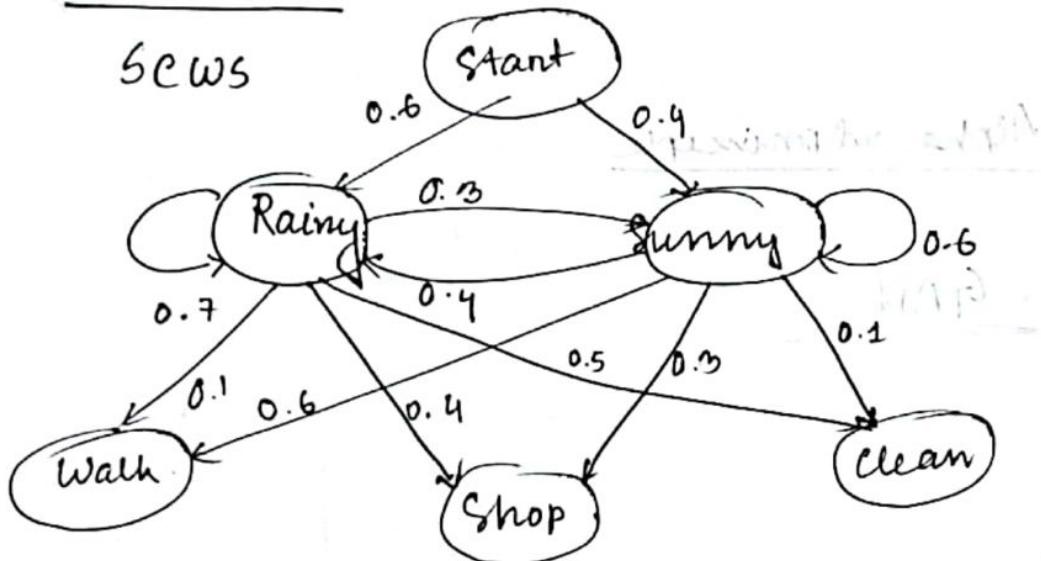
$$c_t = f_t \circ c_{t-1} + i_t \circ \text{tanh}(w_c [h_{t-1}, n_t] + b_c)$$

$$o_t = \sigma(w_o [h_{t-1}, n_t] + b_o)$$

$$h_t = o_t \circ \text{tanh}(c_t)$$

HMM

Softcomp



Transition matrix

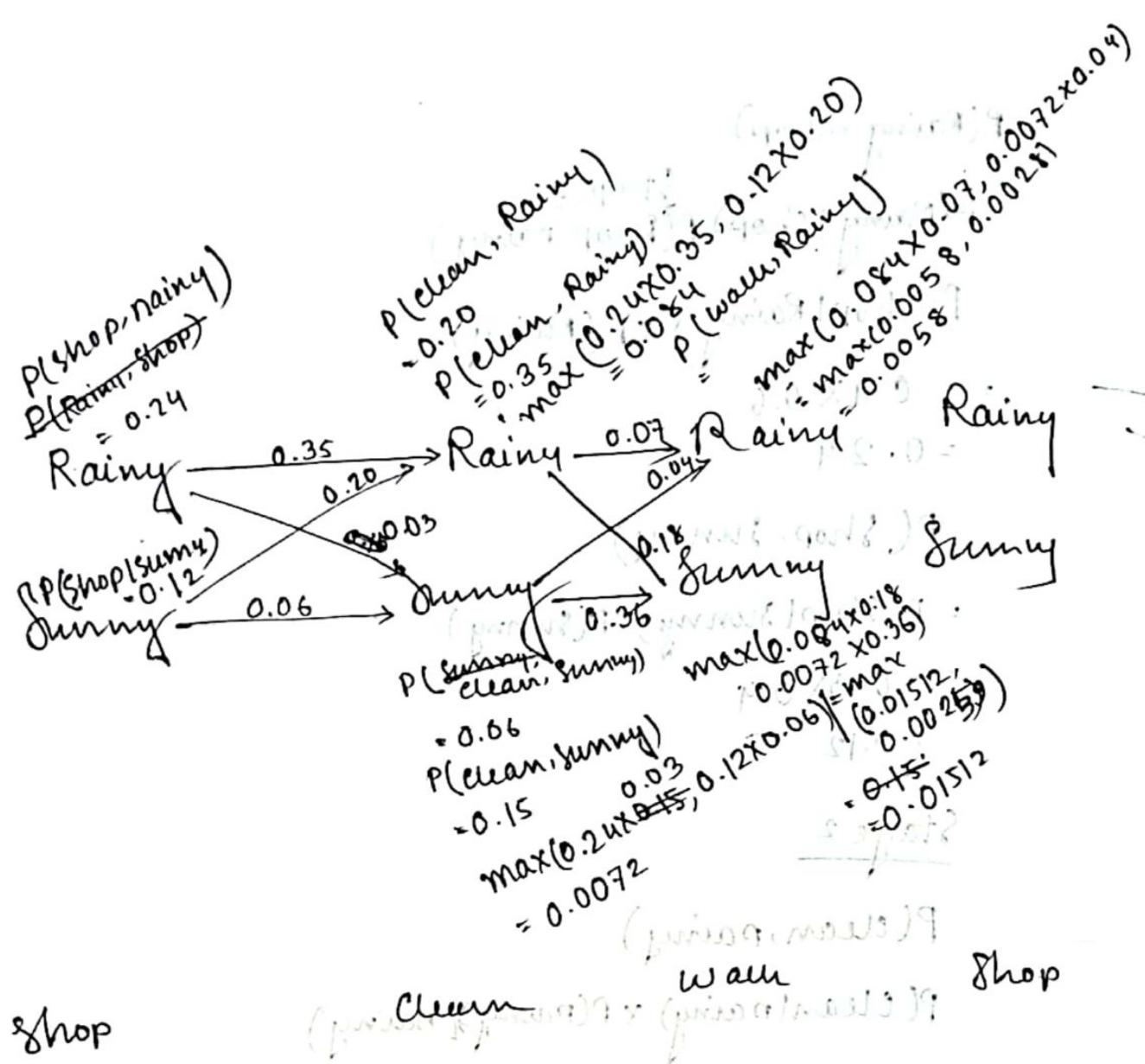
Initial probability [0.6 0.4]

Rainy Sunny

Rainy	0.7	0.3
Sunny	0.4	0.6

Emission matrix

	Walk	Shop	Clean
Rainy	0.1	0.4	0.5
Sunny	0.6	0.3	0.1



shop

0.036

$P(\text{Rainy}, \text{shop})$

$$P(\text{Rainy} \mid \text{shop}) P(\text{Shop} \mid \text{rainy})$$

$$P(\text{Shop} \mid \text{Rainy}) \times P(\text{Rainy})$$

$$= 0.4 \times 0.6 \\ = 0.24$$

$P(\text{Shop}, \text{sunny})$

$$= P(\text{Shop} \mid \text{sunny}) \times P(\text{sunny})$$

$$= 0.3 \times 0.4 \\ = 0.12$$

Stage 2

$P(\text{clean}, \text{rainy})$

$$P(\text{Clean} \mid \text{rainy}) \times P(\text{rainy} \mid \text{rainy})$$

$$= 0.5 \times 0.7 \\ = 0.35$$

$P(\text{clean}, \text{sunny})$

$$= P(\text{Clean} \mid \text{sunny}) \times P(\text{sunny} \mid \text{sunny})$$

$$= 0.1 \times 0.6 \\ = 0.06$$

$$P(\text{Clean, Rainy})$$

$$= P(\text{Clean} | \text{Rainy}) P(\text{Rainy} | \text{Sunny})$$

$$\approx 0.5 \times 0.4$$

$$= 0.20$$

$$P(\text{Clean, Sunny})$$

$$= P(\text{Clean} | \text{Sunny}) P(\text{Sunny} | \text{Rainy})$$

$$= \frac{0.1}{0.5} \times 0.3$$

$$= 0.15 \times 0.3$$

P Stage 3

$$P(\text{Walk, Rainy})$$

$$= P(\text{Walk} | \text{Rainy}) P(\text{Rainy} | \text{Rainy})$$

$$= 0.1 \times 0.7$$

$$= 0.07$$

$$P(\text{Walk, Sunny})$$

$$= P(\text{Walk} | \text{Sunny}) P(\text{Sunny} | \text{Sunny})$$

$$= 0.6 \times 0.6$$

$$= 0.36$$

$$P(\text{Walk, Rainy})$$

$$= P(\text{Walk} | \text{Rainy}) P(\text{Rainy} | \text{Sunny})$$

$$= 0.1 \times 0.4 = 0.04$$

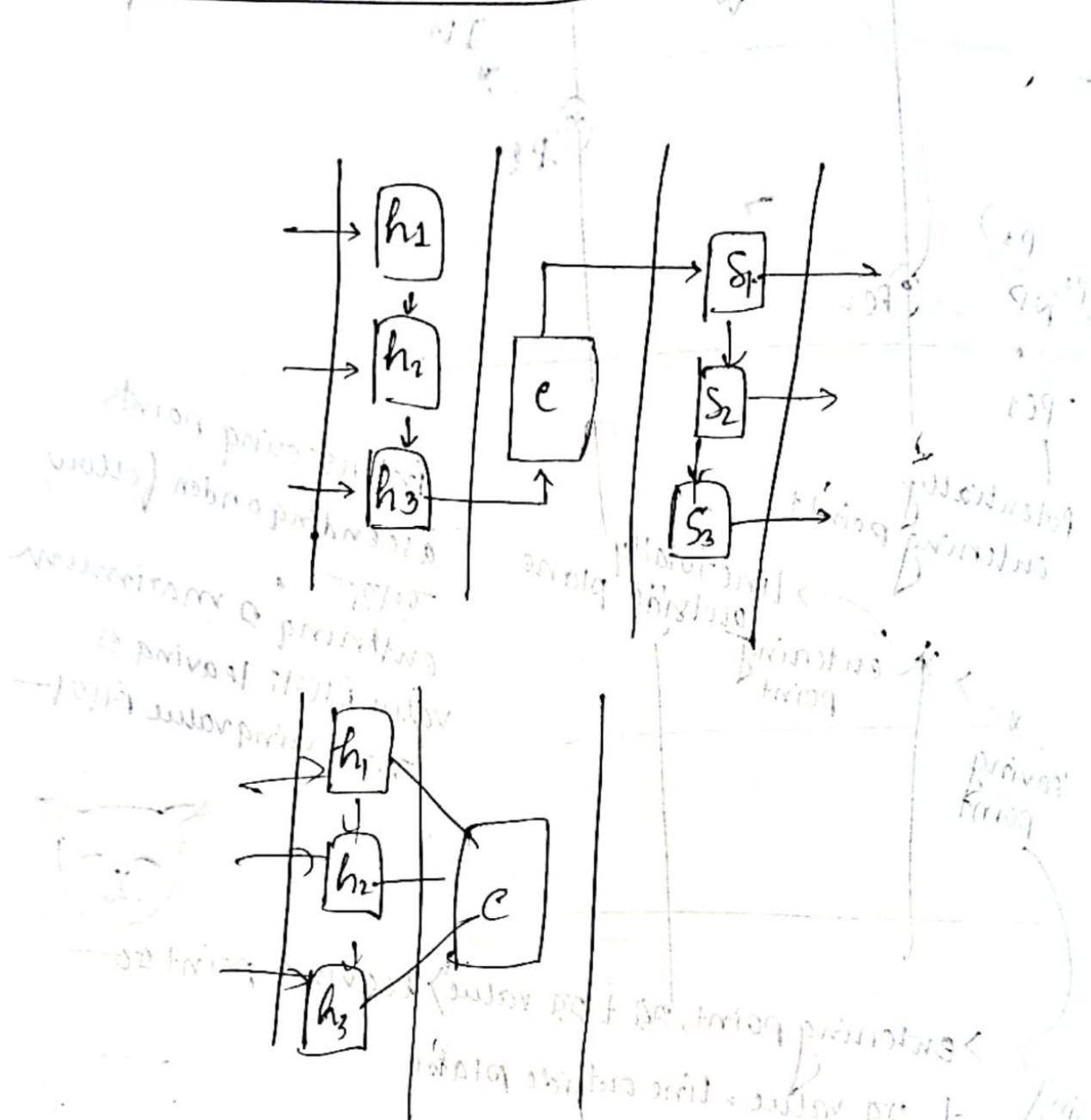
$$P(\text{Walk, Sunny})$$

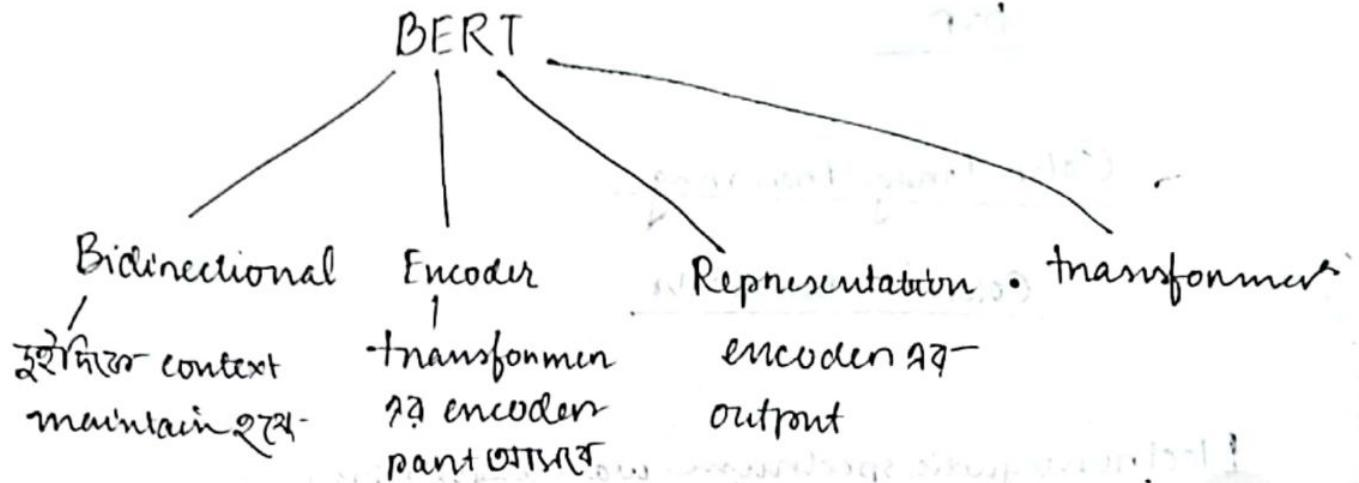
$$= P(\text{Walk} | \text{Sunny}) \times P(\text{Sunny} | \text{Rainy})$$

$$= 0.6 \times 0.3 = 0.18$$

Soft Comp

Attention mechanism



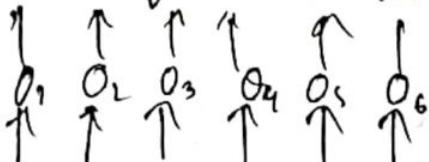


vanilla BERT à decoden गणना करे

word embeddings keepen met next step's verbinding

word embeddings of previous step's + step forward

Classification layer



transformen encoder



embedding

$w_1, w_2, w_3, w_4, \dots, w_n$

mask

BERT-Masked LM