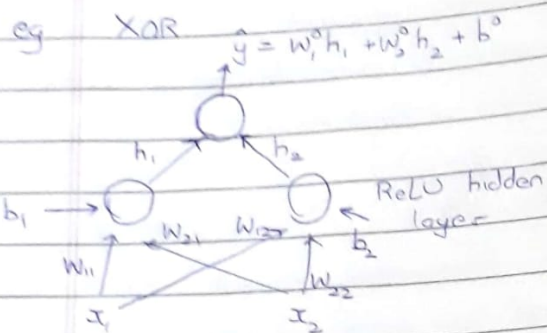


Consider the case for 'square'

3



Initially, all $w = 0$
 $b^0 = 0, b_1 = 1, b_2 = 0$

One iteration: over one training example $(x_1, x_2, y) \equiv (1, 0, 1)$

Step 1 Forward Pass

$$h_1 = \text{ReLU}(w_{11}x_1 + w_{21}x_2 + b_1) = \max(0, 1) = 1$$

$$h_2 = \text{ReLU}(w_{12}x_1 + w_{22}x_2 + b_2) = \max(0, -1) = 0$$

$$\hat{y} = w_0^0 h_1 + w_3^0 h_2 + b^0 = 0$$

$$\text{Square loss} = L = (y - \hat{y})^2 = 1$$

Step 2 Back Propagation

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

$$1. \frac{\partial L}{\partial \hat{y}} = -2$$

$$\frac{\partial L}{\partial w_0^0} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_0^0} = -2 h_1 = -2$$

Update w_0^0, w_3^0 similarly & b^0

$$2. \frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial h_1} = -2 w_0^0 = 0$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial h_1} \times \frac{\partial h_1}{\partial w_{11}} = 0 \cdot \text{ReLU}'(w_{11}x_1 + w_{21}x_2 + b_1) = 0 \cdot 1 \cdot x_1 = 0$$

Similarly update w_{11} & w_{21} & b_1

F]

Regularization

- Reducing generalization error of a model, even at the cost of training error
- Objective function \rightarrow $e_{\text{loss}} + \lambda \times \text{regularizer}$

A

Regularization :- Early Stopping

- Training error always decreases, validation error decreases, then increases in case of overfitting.
- Stop training when validation set increases more than certain number of times.

RECURRENT NEURAL NETWORKS

→ For sequential inputs

$$x_1, \dots, x_n \rightarrow y$$

1. Classifying sequences

eg - Sentiment classification

$$x_1, \dots, x_n \rightarrow x_{n+1}$$

2. Next word in a sequence

eg - Language modeling, predictive text

$$x_1, \dots, x_n \rightarrow y_1, \dots, y_n$$

3. Label per token in a sequence

eg - Parts of speech

$$x \rightarrow y_1, \dots, y_m$$

4. Sequence prediction

eg - Translation

• Examples :-

eg Forecasting (in a time series)

$$I/p \quad :- \quad x_1, \dots, x_t$$

where x_i = day of the week / holiday or not / temperature

$$O/p \quad :- \quad y_1, \dots, y_t$$

where y_i = demand for an item

★ RNN

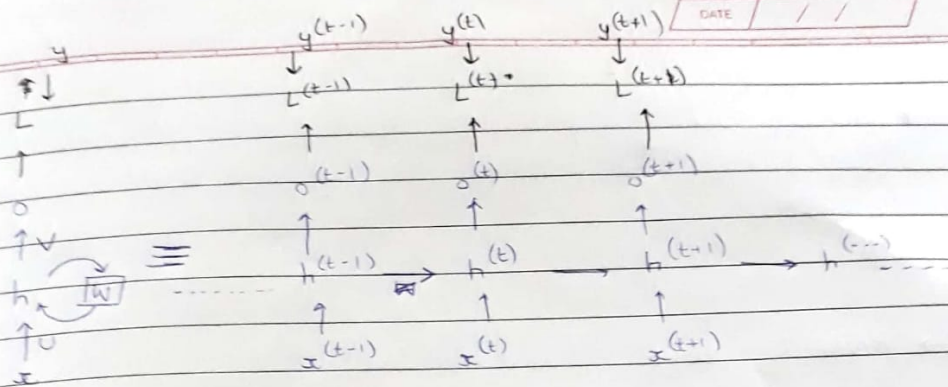
- Processing 1D input of variable length

* In CNN, each hidden output is a function of corresponding and some immediate neighbours.

• In RNN, each output is a function of a 'state' summarizing previous inputs and current input.

- State summary is computed recursively.

- RNN allows deeper, longer-range interaction among parameters than CNNs for the same cost.



$$o^t = C + Vh^t$$

$$h^t = \sigma(b + Wh^{t-1} + Ux^t)$$

↑ Sigmoid (convention)

* The end of sequence is represented by a 'token' in the input

• Parameters :- U, V, W

→ RNN :- Forward Computation Example

Sequence $(x_1, y_1) \dots (x_t, y_t)$
where $x_i, y_i \in \mathbb{R}$

Say $h \in \mathbb{R}^4$, $\Rightarrow W \in \mathbb{R}^{4 \times 4}$, $V \in \mathbb{R}^{1 \times 4}$, $U \in \mathbb{R}^{4 \times 1}$
 $b \in \mathbb{R}^{4 \times 1}$, $c \in \mathbb{R}^1$

Consider $h^0 = 0$

Calculate $h^1 = \sigma(W h^0 + U x_1 + b)$

$$o^1 = V h^1 + c$$

Consider square loss :- $L_1 = (y_1 - o^1)^2$

Similarly find h^2, o^2, L_2

* Practically, we actually lengthen input sequence ('padding') by junk data. While calculating loss, we must mask the padding contribution.

A] \rightarrow RNN for text (Word Prediction)

• We 'embed' distinct words to the real number space

Every x_t denotes a word $\in [1, \dots, V]$
 \uparrow
 $V = \text{'Vocabulary Size'} \approx 30000$

For each of the 30000 words, we make a vector embedding

Embedding matrix $S_{V \times d}$

where $d \approx 300$

Define $\vec{S}_{x_t} S_{x_t} = x_t S$ = vector embedding of that word



x_t is a $1 \times V$ one-hot vector

(1 at one position depending on word input, 0 at all other elements)

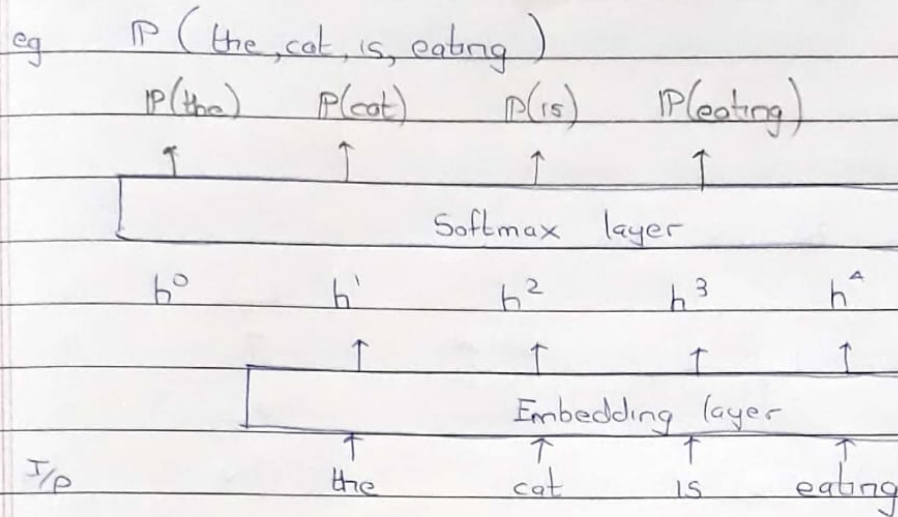
We now give $S_{x_t} = x_t S$ as input to RNN to compute current h_t

* Words of similar meaning are placed closer in the (standardized) 300-dimensional space

$\$ \rightarrow S$ is a look-up table that was learned

→ Training a sequence model

- Maximum Likelihood :- $P(y | x, \theta) = \prod P(y_t | y_1, \dots, y_{t-1}, x, \theta)$



- Training data

$$D \equiv \{(x^1, y^1) \dots (x^N, y^N)\}$$

Language model :- We want $y^{t+1} = x^{t+1}$ (prediction)

Training data :- \vec{x} :- The cat is ~~sleep~~ eating

\vec{y} :- cat is eating EOS

↳ end-of-sequence

See following example

A

→ Backpropagation through time

$$\text{Total loss} = L = L_{t-1} + l_t \quad \text{--- add loss at all time}$$

$$-\frac{\partial L}{\partial v} = \frac{\partial L_{t-1}}{\partial v} + \frac{\partial l_t}{\partial v}$$

$$= \frac{\partial L_{t-1}}{\partial o_{t-1}} h^{t-1} + \frac{\partial l_t}{\partial o_t} h^t$$

$$-\frac{\partial L}{\partial w} = \frac{\partial L_{t-1}}{\partial w} + \frac{\partial l_t}{\partial w}$$

$$= \frac{\partial L_{t-1}}{\partial o_{t-1}} \times \frac{\partial o_{t-1}}{\partial h^{t-1}} \times \frac{\partial h^{t-1}}{\partial w} + \dots$$

$$= \frac{\partial L_{t-1}}{\partial o_{t-1}} \times \left(\frac{\partial h^{t-2}}{\partial w} + W \frac{\partial h^{t-2}}{\partial w} \right) + \dots$$

This will have all terms till $\frac{\partial h^{t-1}}{\partial w}$

with coefficient W^t (can explode or vanish)

∴ We have to go back over all h^n to compute gradient
'Exploding or Vanishing Gradient'

Gradient is either small or large

- Solutions for vanishing/exploding gradient

Multiple time scales :-

Add direct connection from far past inputs to output instead of depending on state to capture all past inputs.

∴ Cannot change how far back we look at different times or for different inputs.

Solution :- Gated RNNs, (eg - LSTMs)

eg/lo

$B \rightarrow$ Sequence Prediction

$$\vec{x} = (x_1, \dots, x_n) \longrightarrow \vec{y} = (y_1, \dots, y_n)$$

eg - sentences, image, audio

Every y_i is called a 'token' (any term from a huge vocabulary)

eg Translation :- $\langle \text{English sentence} \rangle \longrightarrow \langle \text{Hindi sentence} \rangle$

eg Image captioning :- $\langle \text{Image} \rangle \longrightarrow \langle \text{Describing sentence} \rangle$

eg Conversation assistance

eg Speech recognition :- $\langle \text{Speech spectrogram} \rangle \longrightarrow \langle \text{Phoneme sequence} \rangle$

- Challenges

- 1) Long range dependency

- Does not assume conditional independence
- eg Ram eats. He also sings

- 2) Highly open-ended prediction space

- Even length of output is variable

When $\text{argmax}_y P(y|\dots) = \text{stop}$ the sentence

A Encoder - Decoder Model

1. Encode input x into a fixed-dimension vector X (embedding)

2. Decode \vec{y} token by token using an RNN.

- Initialize RNN with state X .

- Repeat until RNN generates an EOS token :-

- Feed previously generated token as input

- Get a distribution over output tokens and choose the best

Write $P(y_t | y^{(t-1)}, X, \theta) = P(y_t | h_t, \theta)$

↑
state vector implemented by RNN

with $h_0 = X$.

$\text{argmax}_{y_t} P(y_t | h_t, \theta)$ is a 'word' known as 'token'

Stop RNN if token = EOS

- We are doing greedy search:- Word prediction #n does not depend on word #k for $k \geq n$

→ Attention-based sequence learning:-

↳ Alignment b/w input sequence & output sequence
(correspondence b/w English words and Hindi words)

- An RNN can attend to inputs output of another RNN

Logit = 'Attention' = $A_{tj} = z_{t-1} \cdot h_j$ (dot product)

Take $\text{Softmax}(A_{tj})$ for producing t_j^{th} output element.

'Attention-weighted input states'

$\sum a_{tj} \cdot h_j$ to produce z_t

CLUSTERING

Unsupervised Learning

- Data is unlabelled. We want to group similar elements.

Input:- $D \equiv \{x^1, \dots, x^N\}$

1. Dimensional Representation:-

Let $x^i \in \mathbb{R}^d$. Find Euclidean distance
OR

2. Distance function

$$d(x^i, x^j) \in \mathbb{R}$$

→ Parameters

1. $k \equiv$ No. of clusters (taken as input)
2. Hierarchy of clusters
3. Correlation clustering :- 'signed' distances

* Deduplication :- Remove duplicate addresses

If instances belong to same address → +ve distance
different → -ve

- Based on clusters :-

Interested in :-

- 1) Centroid of each cluster
- 2) Partitioning data into clusters

→ Objectives:-

1. Assign one of k clusters $C_i(p)$ to each instance ' i '

$$C_1 \cup C_2 \dots C_k = D$$
$$C_i \cap C_j = \emptyset$$

Minimize $\sum_{k=1}^k \left(\sum_{(i,j) \in C_k} d(x^i, x^j) \right)$ pairwise distances in each cluster.

2: Cluster centers m_1, \dots, m_k

- Most popular distance function \equiv Euclidean

Algorithm Begin at $t=0$

K-MEANS ALGORITHM

$m_1^t, \dots, m_k^t \equiv$ random k points from x^1, \dots, x^N

\rightarrow // representative of each cluster

$C^t(x^i) \rightarrow$ randomly from $1, \dots, k$

while (changing) {

for ($i=1$ to N) {

$$C^{t+1}(x^i) = \underset{k \in \{1, \dots, k\}}{\operatorname{argmin}} d(x^i, m_k^t)$$

}

for ($k=1$ to k) { // update representatives

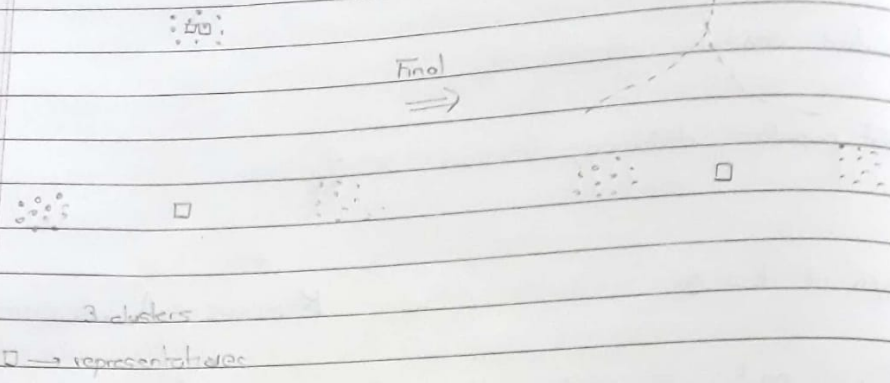
$$m_k^{t+1} = \min_m \sum_{x^i \in C_k} d(x^i, m)$$

// If d is Euclidean, then m_k^{t+1} will be equal to mean of all points in that cluster
(Central Limit Theorem)

}

}

eg Failure :- Unoptimal clustering



• Solution :- Run multiple times, with different initial clustering

→ Proof of convergence of K-means algorithm

$$\text{Define } F(\{C(x^i)\}, \{m_k\}) = \sum_{k=1}^K \sum_{\substack{i=1 \\ C(x^i)=k}}^N d(x^i, m_k)$$

Step 1 :- By design,

$$F(\{C^{t+1}\}, \{m_k^t\}) \leq F(\{C^t\}, \{m_k^t\})$$

Step 2 :- By design,

$$F(\{C^{t+1}\}, \{m_k^{t+1}\}) \leq F(\{C^{t+1}\}, \{m_k^t\})$$

We will exit the while loop if ~~not~~ equality occurs.

∴ F is always decreasing.

∴ We will never repeat any iteration. 😊

∴ Converges to locally optimal solution

→ Probabilistic View of Clustering

Each cluster is a distribution.

Given 1) No. of clusters = k

2) Parametric form of distribution characterizing each cluster
(eg - Gaussian)

$$\text{eg - } f(x^i, \theta_k) \equiv \text{density for cluster } k \\ \equiv N(x^i, \mu_k, \Sigma_k)$$

$$= \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)}$$

3) $D \equiv (x^1 \dots x^N)$

$$\text{Mixture Distribution} \equiv P(x) = \sum_{k=1}^k \pi_k P(x | k)$$

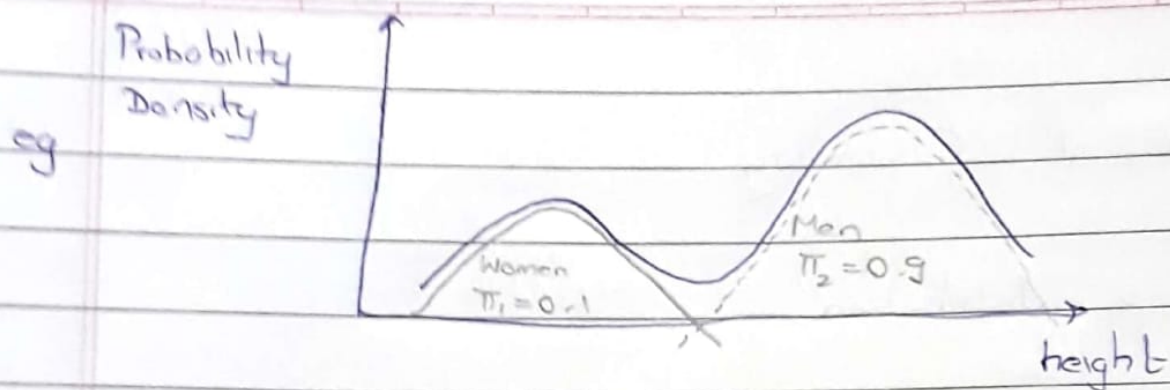
such that $\pi_k \geq 0$, $\sum \pi_k = 1$

$$= \sum_{k=1}^k \pi_k f(x^i, \theta_k)$$

Values of π_i and θ_i for $i \in \{1 \dots k\}$ have to be discovered during clustering by MLE estimation.

→ Training objective :- Maximize $\sum_{j=1}^N \log P(x^j)$
 π_i, θ_i
for $i \in \{1 \dots k\}$

P.T.O.



We want to ~~know~~ find out π_1 & π_2 and θ_1 & θ_2