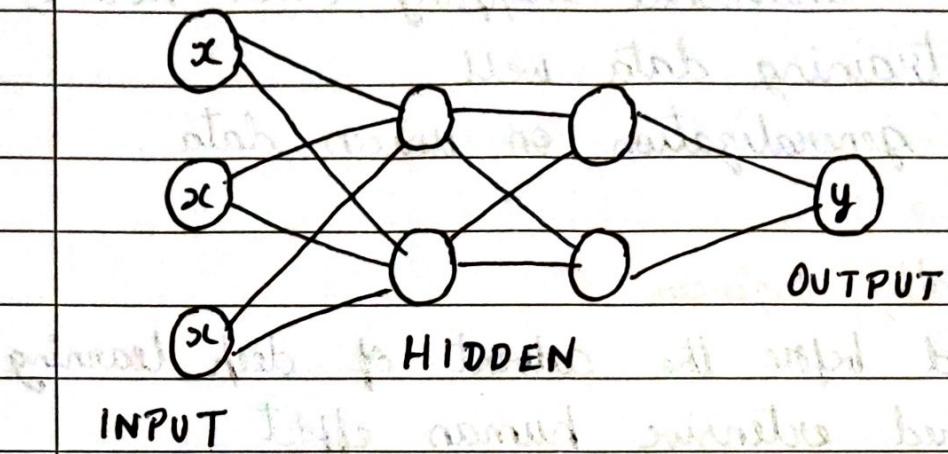


UNIT 1

* MultiLayer Perceptrons

- Feedforward Neural Network



- Aim : Approximate function f that maps input x to output y
- $y = f(x; \theta)$: learn $\theta \rightarrow$ model parameters
- Linear functions \equiv linear models $\equiv x$
- Non linear fn \equiv linear models \equiv transformed input $\phi(x)$
 $\phi \rightarrow$ non linear transformation

* Activation Function :

Used to compute hidden layer values

- Strategies to choose Φ

1) Use generic Φ

- high dimensional mapping (RBF kernel)
- + fits training data well
- poor generalization on unseen data

2) Manually engineer Φ

- existed before the advent of deep learning
- required extensive human effort for each separate task

3) Learn Φ (Deep Learning strategy)

$$y = f_{\text{func}}(x : \theta, w)$$

$$y = \Phi(x : \theta)^T w$$

- learns the mapping through parameters Φ

- + combines adv of generic and manual methods
- sacrifices convexity in training

→ generic : broad family $\Phi(x : \theta)$ is used

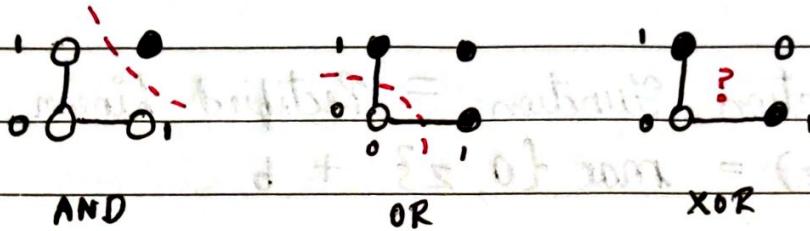
→ manual : humans can design

$\theta \rightarrow$ Params that learn the transformation $\Phi(x)$

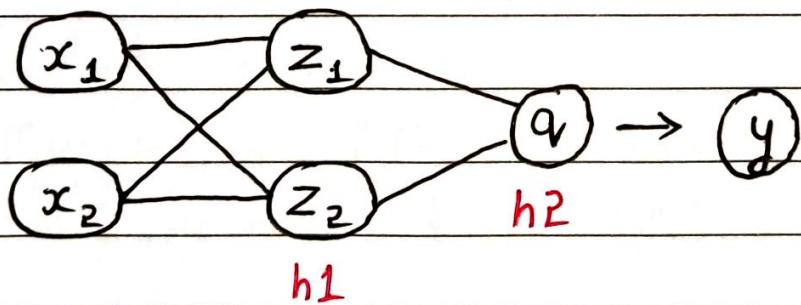
$w \rightarrow$ Params that map the $\Phi(x)$ to output

* Learning XOR

- Single layer Perceptron \equiv linearly separable fn
- AND, OR, NOT, NAND, NOR



- Soln: Using Second layer of Units (MLP)
"Piecewise linear classification using MLP"



$$y = \overline{x_1} \overline{x_2} + \overline{x_1} x_2$$

$$y = z_1 + z_2$$

$$z_1 = \sum x * w$$

If $z_1 \neq e$

Update weights, $w = w + (e - o) \eta x$
(back propagation)

$$z_1 \Rightarrow$$

$$z_2 \Rightarrow$$

* Activation Function = Rectified Linear Unit (ReLU)
$$g(z) = \max\{0, z\} + b$$

x Input

z Transformation of Input
$$z = w^T x + b$$

y Activation function applied to produce output, $y = g(z)$

* Rectified Linear Unit [ReLU]

$$g(z) = \max(0, z)$$

* Maxout Units

- Divides z into k groups

- Each group outputs maximum element

$$g(z) = \max_{j \in C} z_j$$

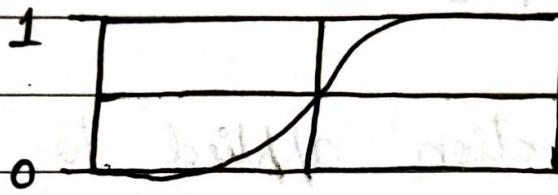
* $g(z, \alpha) = \max(0, z) + \alpha \min(0, z)$

Leaky ReLU $\Rightarrow \alpha = 0.01$ (small value)

Parametric ReLU $\Rightarrow \alpha = \text{learnable Parameter}$

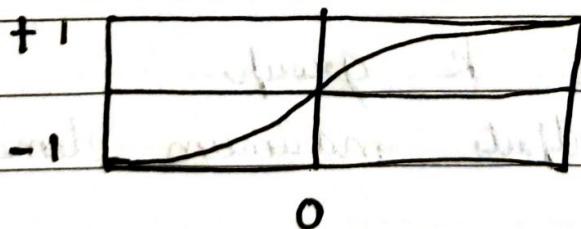
* Logistic Sigmoid

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



* Hyperbolic Tangent

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



* Gradient Based Learning

Training process of neural networks,

- Initialize weights with s, τ values & biases
- SGD (use iterative gradient optimization algo)
- Compute gradients using backpropagation algorithm

- Cost Function

- It reflects how well the model predicts target y given input x
- Cost function based on principle of maximum likelihood. Cross entropy between $P(y|x;\theta)$ actual target & model's prediction
- Total Cost $J(\theta) = \text{Cost } J_n + \text{Regularization}$

Regularization prevents overfitting

L1, L2, Dropout, Early stop

- Eg: Mean Squared Error, Cross Entropy

$$z = \omega^T h + b$$

— / —

* Loss fn for max like of Bernoulli activation

- Classification tasks of y with binary variable
- Bernoulli distribution defined by single number.

$$P(y=1 | x) = \max(0, \min(1, z))$$

problematic, whenever z is outside $[0, 1]$
output is 0 and no longer learned

- Soln: Combine Sigmoid output with max likelihood

$$\hat{y} = \sigma(z) \quad \text{where } \sigma \rightarrow \text{Sigmoid Activation}$$

$$P(y) = \sigma((2y-1)z)$$

- loss function

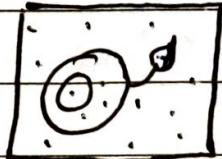
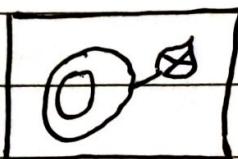
$$\begin{aligned} L &= -\log P(y | x) \\ &= -\log \sigma((2y-1)z) \\ &= \text{softmax}(1 - 2y)z \end{aligned}$$

* Data Augmentation

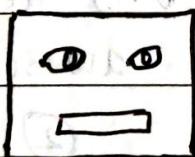
- Generates different versions of a dataset
- More data to train
- Increases accuracy of the model

1) Adding Noise

Salt and Pepper noise

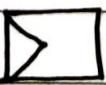


2) Cropping

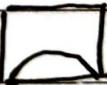


Section of image is cropped

3) Flipping



vertically



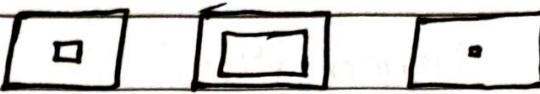
horizontally

4) Rotation



Rotated bw 0 to 360°

5) Scaling



6) Translation



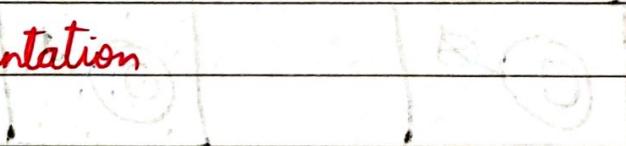
7) Brightness



8) Contrast



9) Color Augmentation



- Back Translation for words

This is } C'est vraiment } It's really
awesome } génial } awesome

① I am not a native speaker. - ②.

11

* Regularization

- avoids overfitting
- model generalizes better
- model performance on unseen data improved
- It is a technique that makes slight modifications to learning algorithms
- L2

$$\text{Cost fn} = \text{Loss} + \frac{\pi}{2m} \sum w^2$$

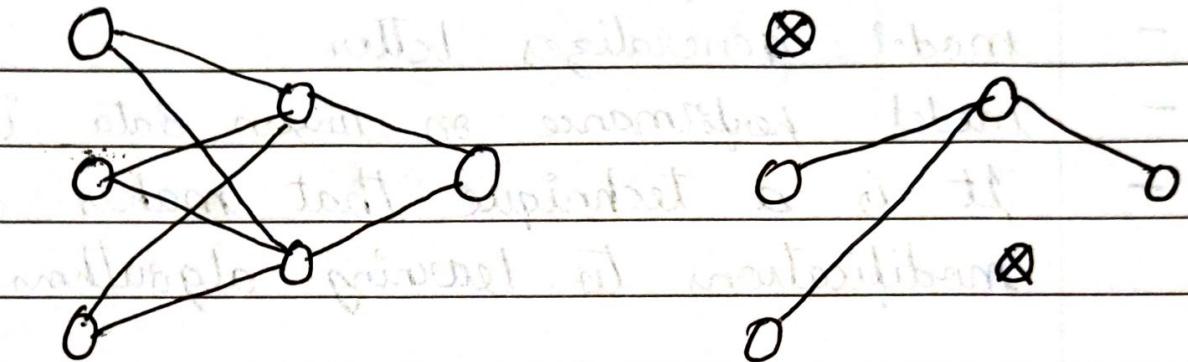
- Weight Decay
- It forces weight to decay towards zero (but not exactly zero)
- π is the regularization parameter, it is optimised for better results

- L1

$$\text{Cost fn} = \text{Loss} + \frac{\pi}{2m} \sum |w|$$

- Penalizes the absolute weights value
- weights can be zero (compresses the model)

- Dropout



- Most frequently used
- At every iteration, some nodes are selected and removed

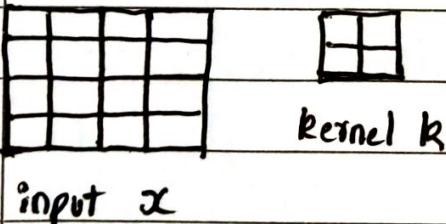
— / —

UNIT 2

- Convolution Operation
- Convolution Operation (for 2D input data)
- Core Component of CNN
- Three ideas : SI, PS, ER
- Visual Cortex is base for CNN
- Types of data used in CNN models
- Stenches in CNN

* Convolution Neural Network [CNN]

- Yann LeCun { 1989 }
- Specialized type of NN for processing data that has grid-like topology
 - 2D grid ≡ images
 - 1D grid ≡ time series
- They employ mathematical operation called convolution in place of general matrix multiplication



Sliding dot product
b/w x and k
convolution operation

$s(t) = \text{Position at time } t$

$$s(t) = (x * k)(t) = \sum (x * k)(t)$$

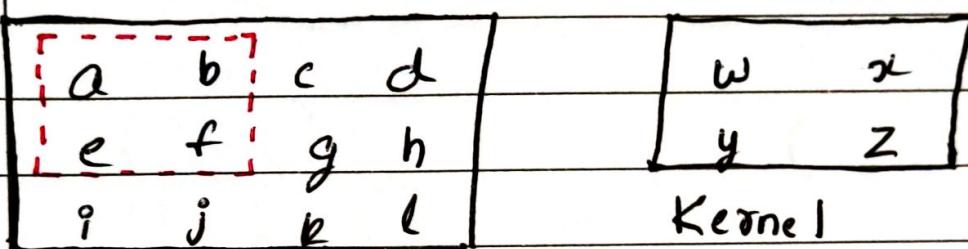
$$s(t) = (x * k)(i, j) = \sum_m \sum_n x(i+m, j+n) k(m, n)$$

Output: feature map]

— / —

* Convolution to 2D tensor

Tensor : Multidimensional Array



Input

$$\begin{bmatrix} aw + bx + \\ ey + fz \end{bmatrix} \dots \begin{bmatrix} bw + cx + \\ fy + gz \end{bmatrix} \dots \begin{bmatrix} gw + hx \\ ky + lz \end{bmatrix}$$

Output

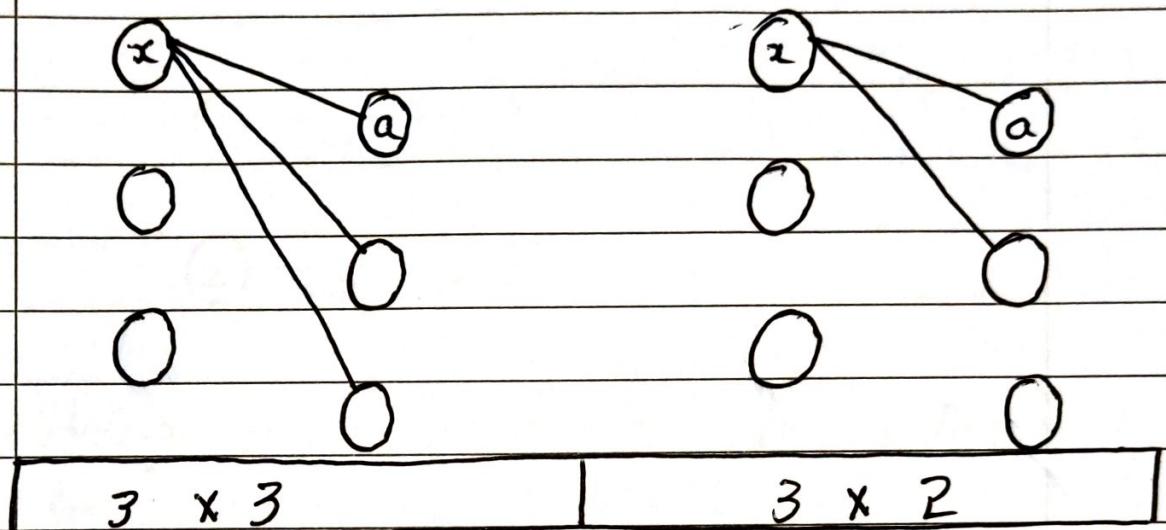
Sparse interaction

Parameter sharing

Covariant representation

_____ / / _____

* Sparse Interaction

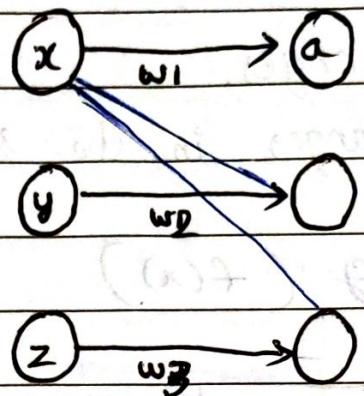


- $SI \equiv SC \equiv SW$
- Accomplished by making **kernel** < **input**
- Using fewer parameters
 - reduces memory
 - improves efficiency [computing output requires fewer operations]

- Input reduce to Output
m k n

$(m \times n)$ params matrix
 $(k \times n)$ params multiplication

* Parameter Sharing



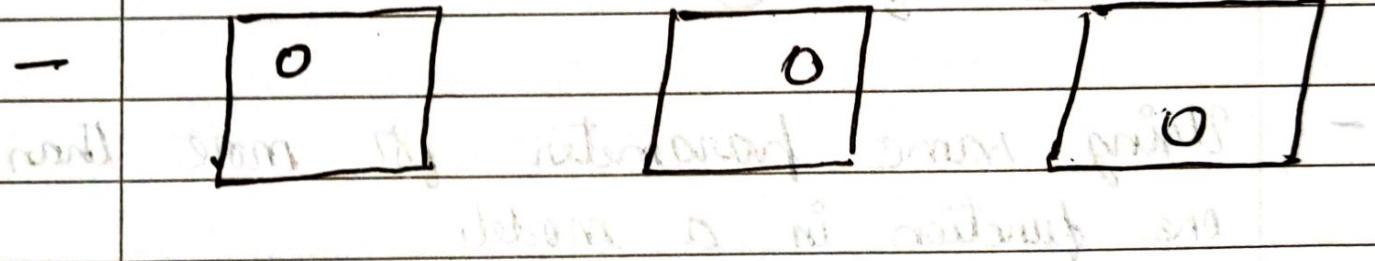
at x , w_1 is used
at y , w_2 is used
at z , w_3 is used

- Using same parameter for more than one function in a model.
- In traditional nn, each element of weight matrix is used exactly once to compute output.
Eg : $w_{11} \times x_1, w_{21} \times x_1, \dots$
 $w_{12} \times x_2, w_{22} \times x_2, \dots$
- In convolution nn, kernel is used at each position of input
- Instead of learning separate set of parameters for every position, we learn only one set
- Parameters of model from $(m \times n)$ to (k)

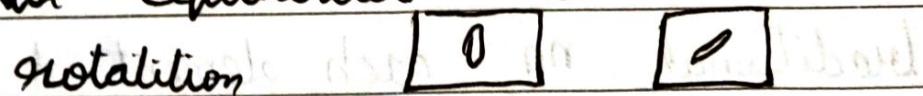
* Equivariance representation

If the input changes,
output changes in the same way

$$f(g(x)) = g(f(x))$$



Not Equivariant to



$$\text{for } X \times \mathbb{R}^d \rightarrow X \times \mathbb{R}^d \quad \mathbb{R}^d$$

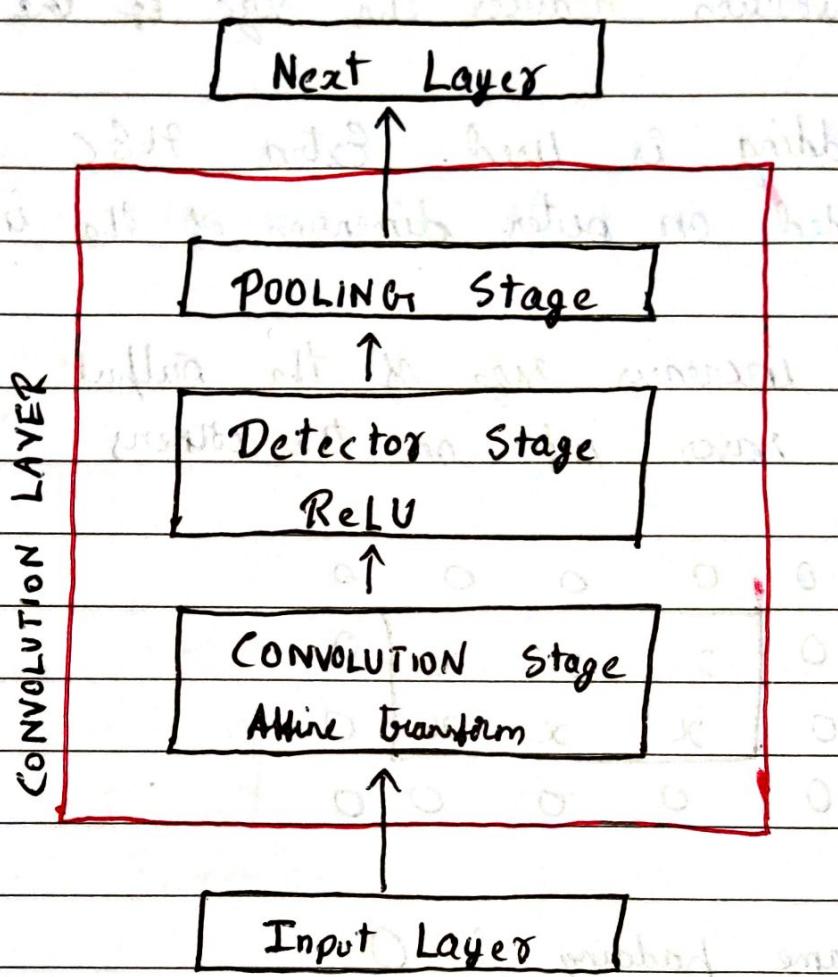
$$\text{for } X \times \mathbb{R}^d \rightarrow X \times \mathbb{R}^d$$

— / /

* Pooling

- down sampling technique
- reduces the size of feature maps
(generated by convolution layer)
- retains important aspects / features
 - i> Dimensionality reduction
 - ii> Feature Extraction
 - iii> Translation Invariance (ER)
 - iv> Generalization
 - v> Reduction in computational time
- Max Pool, Avg Pool, Global Avg

* CORE Components of CNN



*

Stride

No of pixel shifts over input matrix

$$\left[\frac{n + 2p - (f+1)}{s} + 1 \right] \times \left[\dots \right]$$

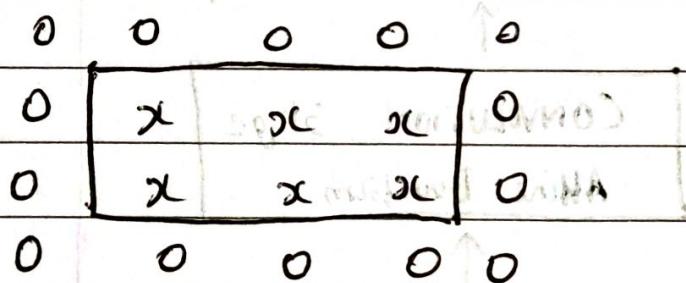
$n \rightarrow$ Input

$p \rightarrow$ padding

$f \rightarrow$ filter

* Padding

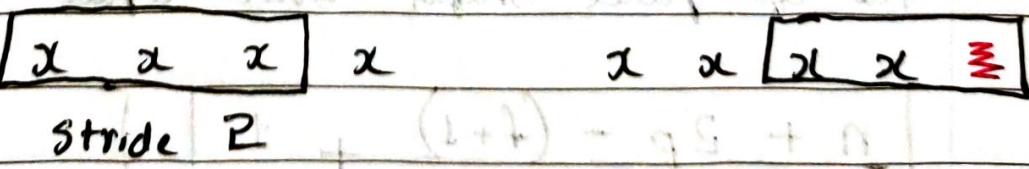
- Convolution reduces the size of the output
- Padding is used. Extra REC are added on outer dimension of the images
- It increases size of the output
- It saves info on the corners



Same padding : 0

Causal padding : x

Valid padding : corner elements are not ignored

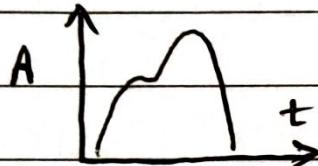


* DATA TYPES

Single Channel

1D

Audio Waveform



Multichannel

Skeleton
Animation
Data



2D

Audio Data

that has been
processed with
Fourier Transform

Color Image Data

Ch1	Ch2	Ch3
Red	Blue	Green

1D audio waveform

↳ 2D tensor

3D

Volumetric Data

Medical Imaging
(CT Scan)

Color Video Data

Ch1	Ch2	Ch3
time	height	width

* Neuroscientific basis for Convolutional Networks

- CN is inspired by the way mammalian visual system processes images
- findings, from Hubel and Wiesel
 - ↓
specific NEURONS in visual cortex respond to pattern of light (edges, bars)
 - ↓
different neurons are tuned to detect specific visual feature/pattern
- CN mimics the way brain processes visual info,
 - i> spatial arrangement
 - ii> simple and complex cell behavior
 - iii> hierarchical feature detection
- i> V1 (Primary Visual Cortex) is a 2D structure.
CN ≡ uses 2D feature maps
- ii> Simple cells respond to local features (Filters in CN)
- ii> Complex cells detect features equivariantly (Pooling in CN)
- iii> Grandmother Cells. Certain cells are activated for specific stimuli, regardless of image presentation

UNIT 3

*

Sequence Modeling

- RNNs are neural networks for processing sequential data
- Sequence of values : $x^{(1)}, \dots, x^{(n)}$
- Eg : Time series data, Text, ...
(Data where order of elements matter)
- Scalability : They can handle longer sequences
- Parameter Sharing : RNNs use same set of parameters (weights) across different time steps
 - o "I went to Nepal in 2009"
 - o "In 2009, I went to Nepal"
- Using same parameter, network can generalize better across different sequence lengths
- In RNN, each output depends on previous output, not just the input
- This recurrent structure creates deep network
(Deep Computational Graph)

* Unfolding Computational Graphs

- Computation Graph is a visual or mathematical representation of computations
- Unfolding is applied to recurrent computation
- It transforms recursive structure into non-recursive

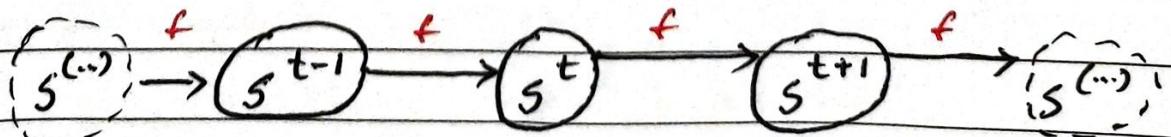
* Dynamical system, $s^{(t)} = f(s^{(t-1)}; \theta)$

- $s^{(t)}$ → state of the system
- It is recurrent, defⁿ of s at t refers back to defⁿ of s at $(t-1)$
- For time steps $T = 3$
Graph can be unfolded by

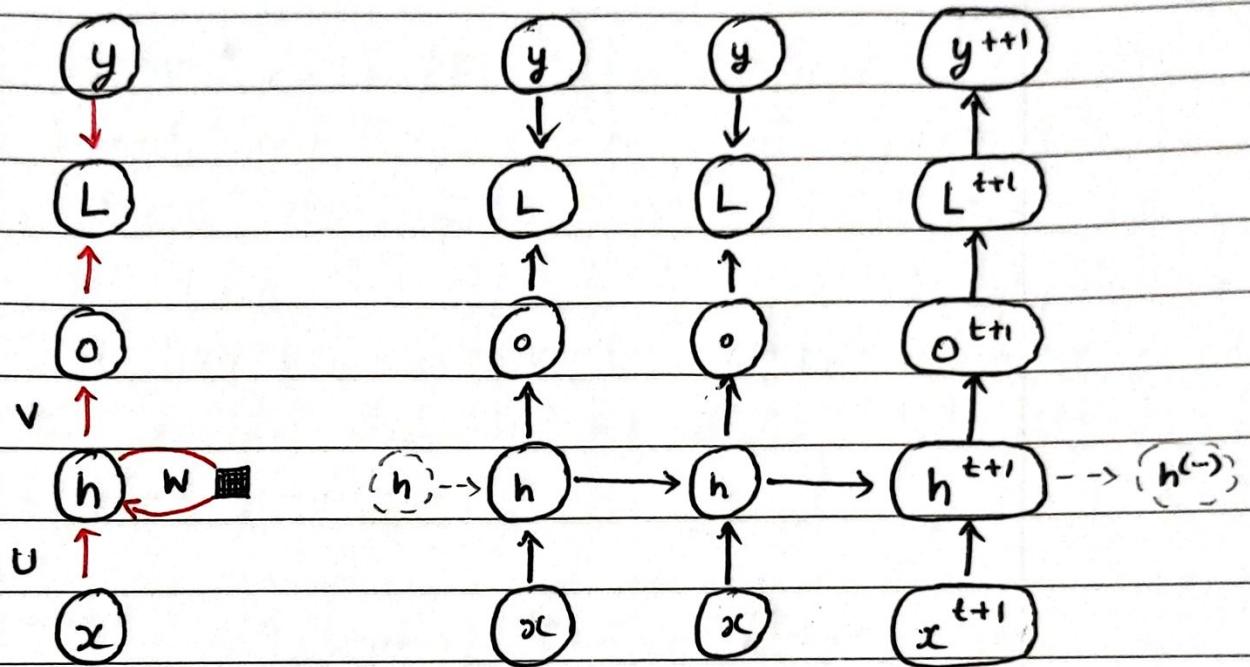
$$s^{(3)} = f(s^{(2)}; \theta)$$

$$s^{(2)} = f(s^{(1)}; \theta)$$

$$s^{(3)} = f(f(s^{(1)}; \theta); \theta)$$



* Recurrent Neural Network



$$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$o^{(t)} = c + V h^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

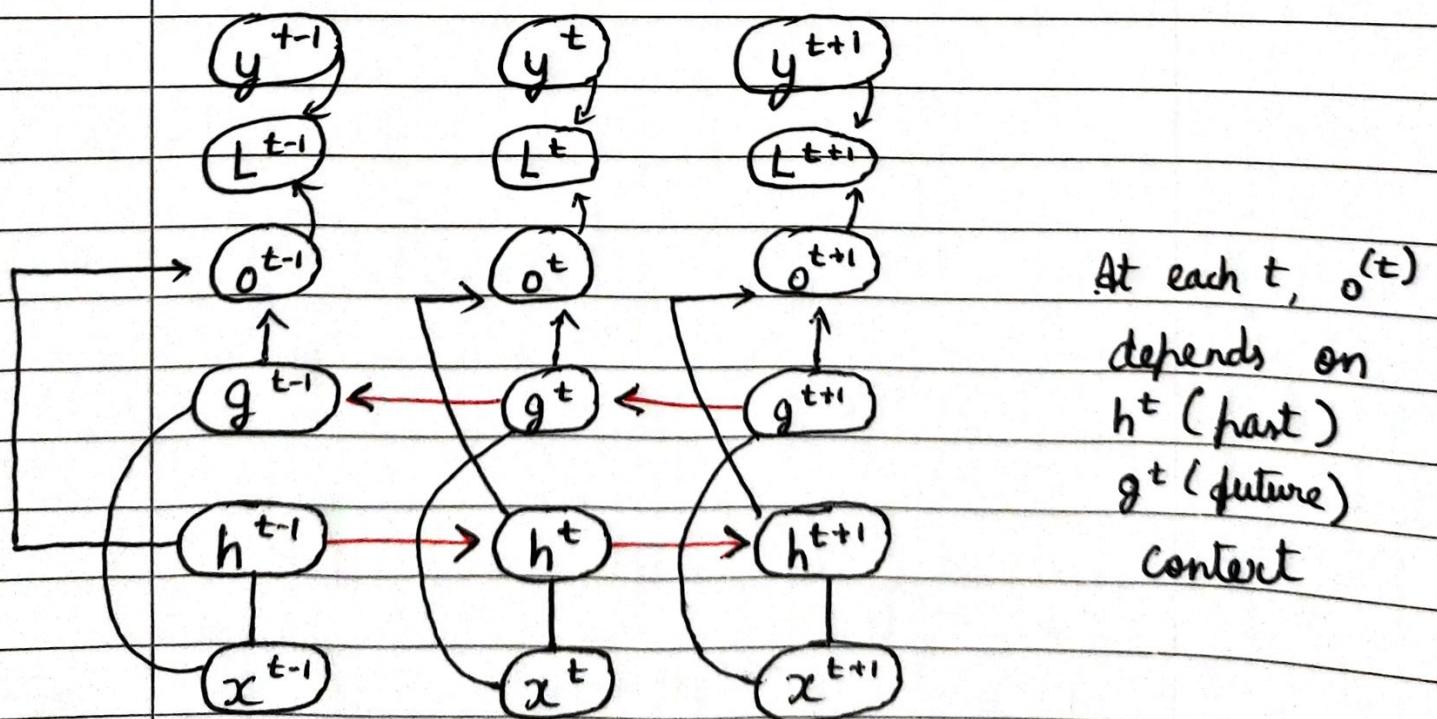
bias $\rightarrow b$ and c

weight matrices $\rightarrow U \quad V \quad W$

i-h h-o h-h

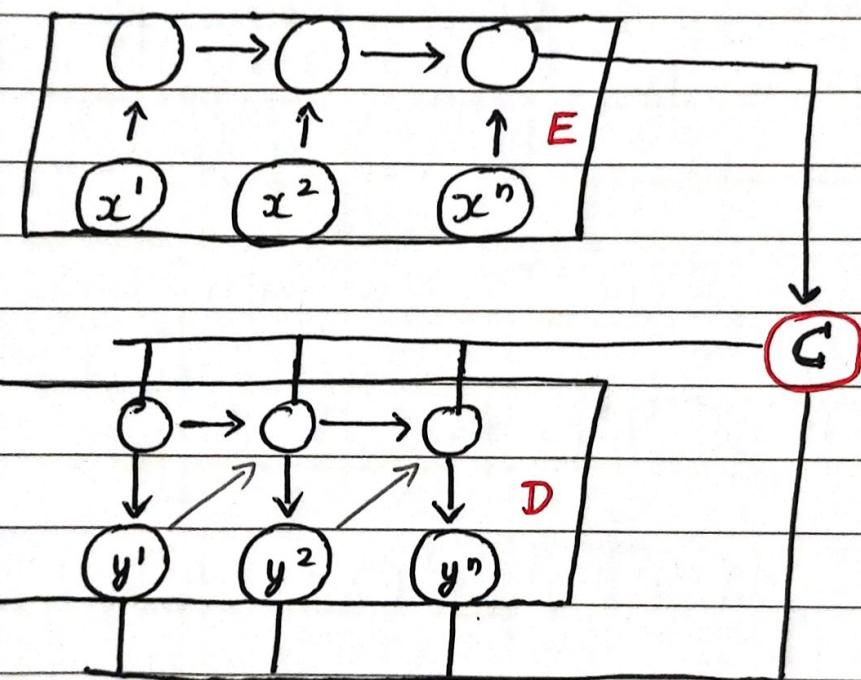
* Bidirectional RNN

- Recurrent network so far has causal structure. State at time t captures info from past, $1, t-1, t$
- For Speech, Language recognition, bioinformatics, the output depends on whole input sequence (future and past)
- RNN that moves forward through time + RNN that moves backward through time



* Encoder - Decoder (Sequence to Sequence)

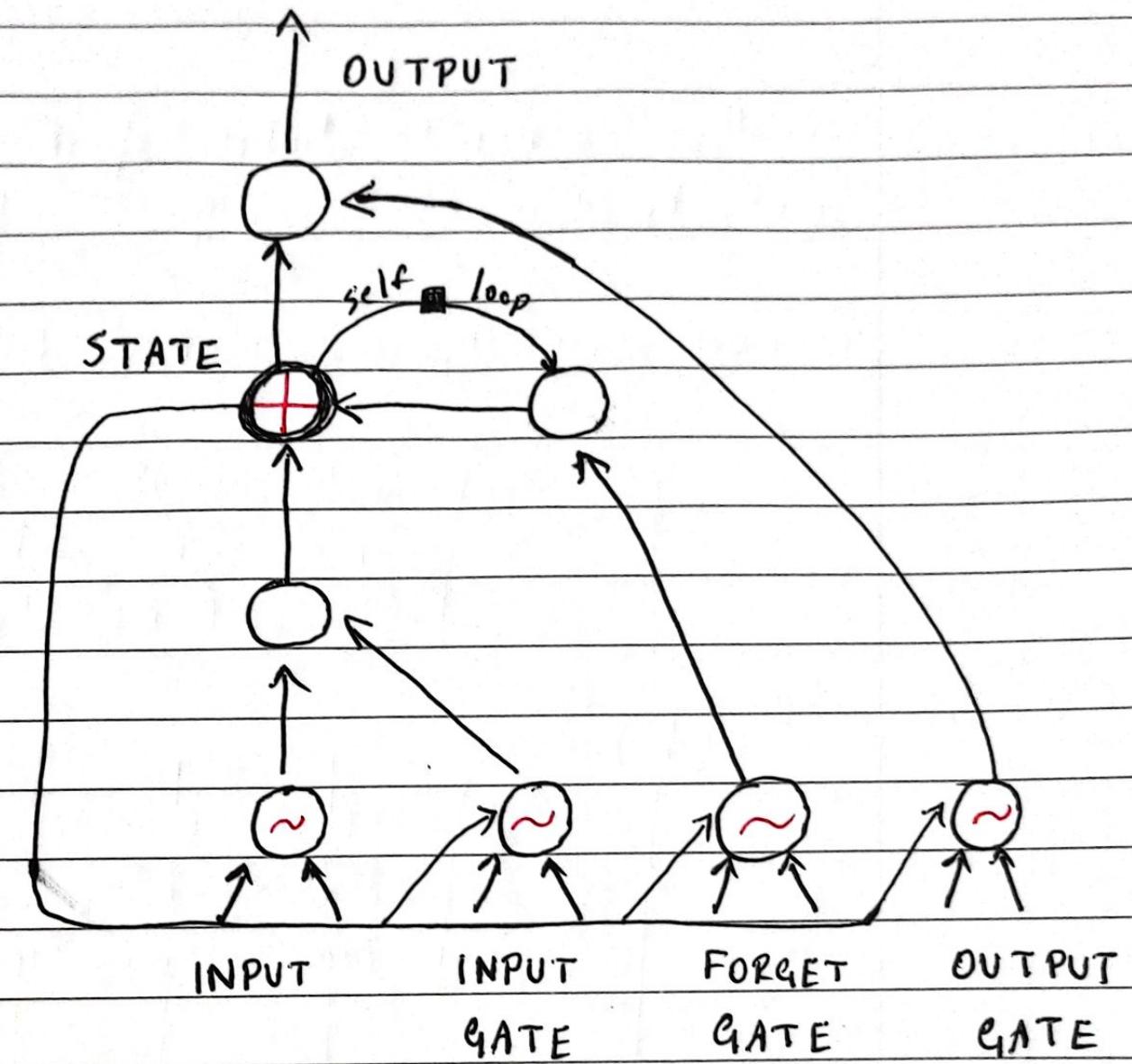
- To map input sequence ~~with~~ to output sequence of different lengths
- Speech translation, Question answering



* Encoder processes the input sequence
 $X = x^1 \ x^2 \ \dots \ x^n$ and generates context vector C

* Decoder uses this C to generate output
 $Y = y^1 \ y^2 \ \dots \ y^n$ of varying length

* LSTM



- Gated recurrent network
- They can learn when to forget, when to retain information
- Use self-loops to maintain gradients over long sequence (weights controlled by f_g)
(dynamic weight)

- Forget Gate

$$f = (b + \sum U x^t + \sum W h^{t-1}) \sigma$$

$x^t \rightarrow$ Current Input Vector

$h^{t-1} \rightarrow$ Previous hidden Vector

$b \rightarrow$ Bias

$U \rightarrow$ Input Weights

$W \rightarrow$ Recurrent Weights

- Input Gate

$$\begin{aligned} i^o &= f^{(t)} i^{(t-1)} + g^{(t)} \sigma \\ &\quad \left(b + \sum U x^t + \sum W h^{t-1} \right) \end{aligned}$$

- External Input Gate

$$g^i = \sigma \left(b + \sum U^g x^t + \sum W^g h^{t-1} \right)$$

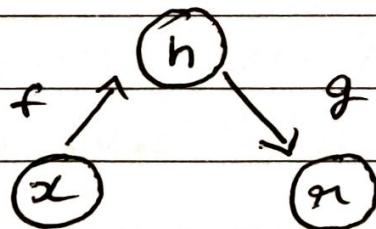
- Output Gate

$$o^i = \sigma \left(b + \sum U^o x^t + \sum W^o h^{t-1} \right)$$

UNIT 4

* AUTOENCODER

- Neural network trained to attempt to copy its input to its output
- Goal: learn compact representation of input
- Encoder: $h = f(x)$
- Decoder: $z = g(h)$



- To avoid perfectly replicating input, they are restricted in some way so they can learn important features of data

(h) \rightarrow internal representation

— / —

* Undercomplete Autoencoders

- Copying input to output is useless, but we are interested in that during training,
(h) will learn useful properties
- By constraining (h) to have smaller dimension than (x), we can obtain useful features
- Autoencoder whose code dimension (h) is less than input dimension (x)
- Learning Process :
Minimizing loss function, $L(x, g(f(x)))$

Loss function = MSE

- Linear function used \Rightarrow performs PCA
- Nonlinear function used \Rightarrow learn complex repr's

* Regularized Autoencoders

- $(h) < (x)$ but E and D have lot of capacity
- $(h) = (x)$
- $(h) > x$

→ Autoencoders fail to learn useful features

* RA provide ability to train any

architecture of autoencoder, choosing

code dimension } based on

capacity of E and D } complexity of distribution

- * They use LOSS FUNCTIONS that encourage models to have other PROPERTIES besides copying $i \rightarrow o$

- Sparsity of the representation
- Smallness of derivative of representation
- Robustness to noise

- * RA can be nL and OC but still learn

1) * Sparse Autoencoder

- It is autoencoder whose training criterion involves a sparsity penalty

$$L(x, g(f(x))) + \Omega(h)$$

reconstruction error

- Used: learn features for another task (classification)
- $\Omega(h)$ \equiv regularizer term added to ffn
 - task \rightarrow copy $i \rightarrow 0$
 - \rightarrow perform supervised tasks
- Regularized Maximum likelihood
 - maximizing $P(\theta | x) =$
 - maximizing $\underbrace{\log P(x | \theta)}_{\text{likelihood}} + \underbrace{\log P(\theta)}_{\text{prior probability}}$

2) * Denoising Autoencoder

- Rather than adding penalty Ω , DA learns something useful by changing the reconstruction error
- $L(x, g(f(x)))$
- DA instead minimizes $L(x, g(f(\tilde{x})))$
 $\tilde{x} \rightarrow$ Copy of x that's been corrupted by some form of noise
- DA must undo this corruption

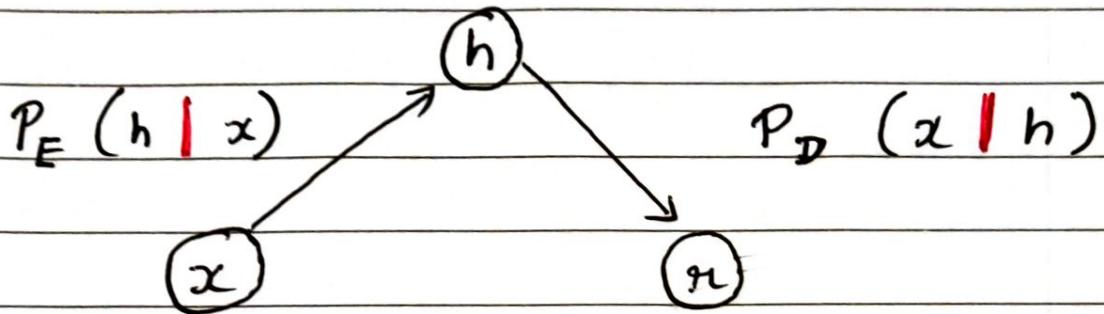
3) * Regularizing by Penalizing Derivative

- CAE (Contractive Autoencoder)
- Penalty $\Omega(h)$ of different form

$$\Omega(h, x) = \lambda \sum_i \|\nabla_{x_i} h_i\|^2$$

\Rightarrow Model learns f_n ~~to~~ f_n doesn't change
when x slightly change

* Stochastic Encoder and Decoder



* We can generalize the notion of,

i) encoding function = $f(x)$



encoding DISTRIBUTION = $P_{\text{encoder}}(h|x)$
= Stochastic ENCODER
= $P_{\text{model}}(h|x)$

ii) $|||^4$

* Application of Autoencoders

1) Dimensionality Reduction

- Reduce no of features while preserving important characteristics
- Better than PCA (Hinton, 2006, RBM)
- Lower dimensional representation
improve PERFORMANCE by ↓ compute burden
↓ memory usage
- Aid in generalization

2) Information Retrieval

- Task of finding relevant entries in a database based on the query
- AE trained to produce low-dimensional, binary code

Input : { binary code → entries }

Semantic hashing

3) Gen. Modelling

Boltzmann Machines

- Data can be image, text