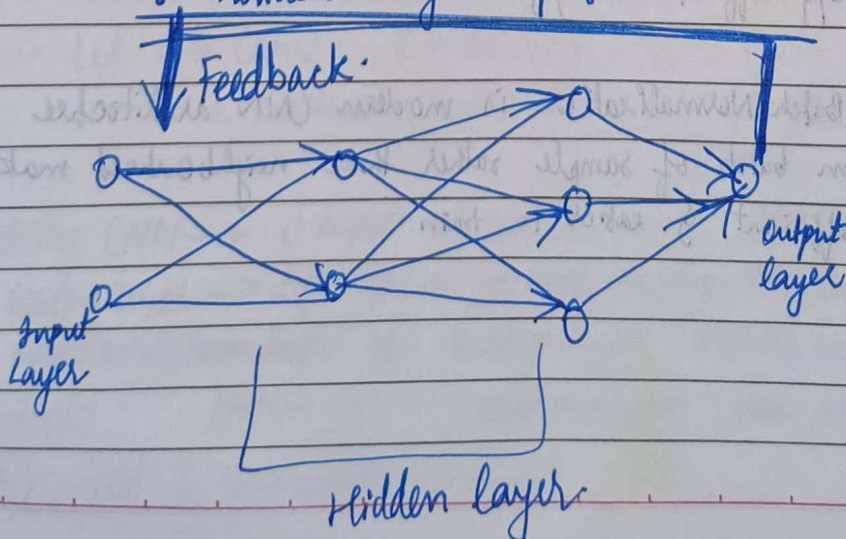# U-4 (DL)
## (Recurrent Neural Networks)

**✱ RNN -**

- Recurrent Neural Networks.
- Type of artifical neural network architecture
- Handles sequential data
- Internal memory that allows them remember past inputs & use
  context current input.
  - remember each-info through time

**Advantages →** • Memory ( past inputs & use info to process current ay)
- useful in time series • Variable Length Inputs
- extend effective • efficient processing.
  pixel neighbourhood
  - process sequence of data, as an output &
  receives a sequence of information as input.
- even used in convolutional layers to extend neighbourhood

**Disadvantages →** • Training challenges
- Gradient Vanishing & exploding problems.
- Can't process long sequence ( • Language transl^n
- slower computation. • text to speech.

**Need →** • Sequential Data Processing
- Temporal Dependencies
- Variable Length Inputs



Feedback.

Input Layer

output layer

Hidden layer

# * Working RNN -

- Input — an element from a sequence as input.
- Hidden Layer Magic — Holds info from past elements
- Information Mix — current input & hidden layer's memory are combined.
- Activation & Output — processed by activation fun$^n$ to generate an output
- Memory Update — Its memory based on processed information.

# * Bi-directional Recurrent Neural Network -

- to process sequential data
- Type of RNN
- Capture contextual dependencies in input data by processing.
- one process data → forward direction (Forward RNN)
- 2$^{nd}$ process data → reverse direction. (Backward RNN)

- Adv →
  - Improved Accuracy
  - efficient handling of variable length sequences
  - Better Long-Term Dependency Capture
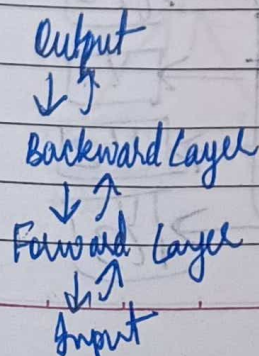  - content from both past & future

- Disadv -
  - Increased computation complexity
  - more Parameters to Train
  - limited to Offline Process

- working →
  - Inputting a sequence
  - Dual Processing
  - Computing hidden state
  - Determining output
  - Training

- Appl$^n$ -
  - Sentiment Analysis
  - Named Entity Recognition
  - Part of Speech Jagging
  - Machine Translation
  - Speech Recognition.

Output
↓ ↑
Backward Layer
↓ ↑
Forward Layer
↓ ↑
Input

# ✳ Encoder - Decoder Architecture -

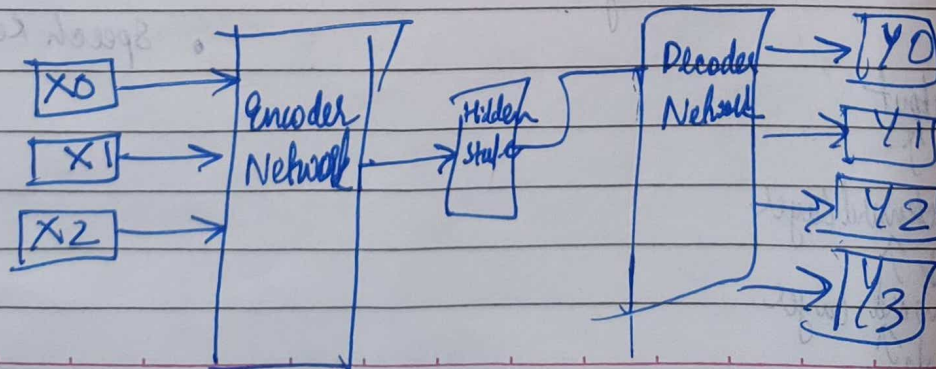| | Encoder | Decoder = |
|---|---|---|
| | • compresses input sequence into a fixed length representation | • fixed length repⁿ into a input sequence |
| Input → | • variable length sequence (eg - sentence) | • fixed length representⁿ (context vector) |
| output → | • fixed length representation (context vector) | • variable length sequence (translated sentence) |
| Processing ↳ | • step by step using RNNs (LSTMs GRUs) to capture long term dependencies | • step by step using RNNs considering previously generated output & context vector |
| complexity → | simpler operation | • more complex opexⁿ due to considering history - |
| Location → | Typically at transmitting end | • typically at receiving end- |
| eg → | • Text summarization. encodes a long document, generate a concise summary | • Machine translation. decoded source language generates target language text |

# Recursive Neural Network. (RvNNs).

* another generalization of recurrent network with diff compu$^n$ graphs.
* structured as a deep tree.
* variable input sequence can be mapped to a fixed size output
* more suited for time-series data → linguistic structure aside.
* born to model by NLP syntax parse trees.
* don't take token sequentially, recursive models combine neighbor
* RVNNs aka TreeNets.

Adv → Handles complex structure     Diddv → Computationaly Expnsive
    → Compositionality        → Limited applicability
    → Expressive Power       → Parallelization Challenges.
    → Reduced Network Depth
    → Theoritical Underpinning
    → Emerging Appl$^n$

Steps → • Start at root node of hierarchy
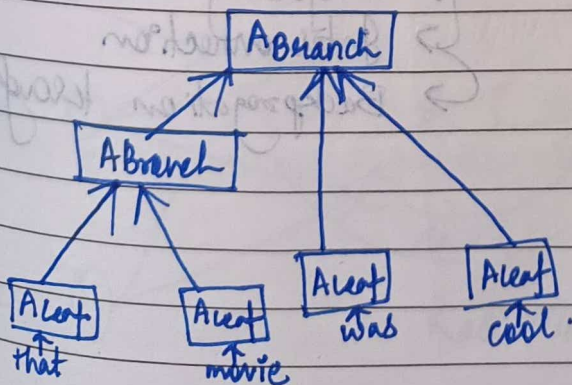     • Process elements at root node
     • Recursively approach each child node
     * Apply same processing logic to child
     * combine processed child's output with parents output
     • Move the hierarchy & repeat steps until 2-3 until all nodes are processed
     • Final combined output represent entire hierachial structure

**\* Deep Recurrent Neural Network —**

- Type of RNN architecture → handle sequential data where order of elements matter
- stacking multiple RNNs
- captures more complex rel^n & dependencies within sequential data compared to regular RNNs.

**Adv →**
- enchanced learning of Long Term dependencies.
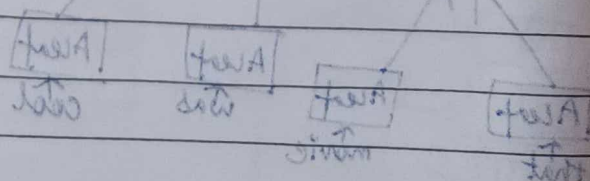- improved representation power
- Flexibility of various task

**Disadv →**
- Increased Computational Cost
- Vanishing/ Exploading Gradient Problem
- Potential for Overfitting

**\* Unfolding Computational Graph —**

- way to formulize structure of set of computations.
- consist of nodes that represent mathematical oper^n or variables, edges → flow of data b/w these operations.

**Adv →**
- Clarity
- Traning
- Analysis

- use same fun^n with same parameters.

→ Single Time Step
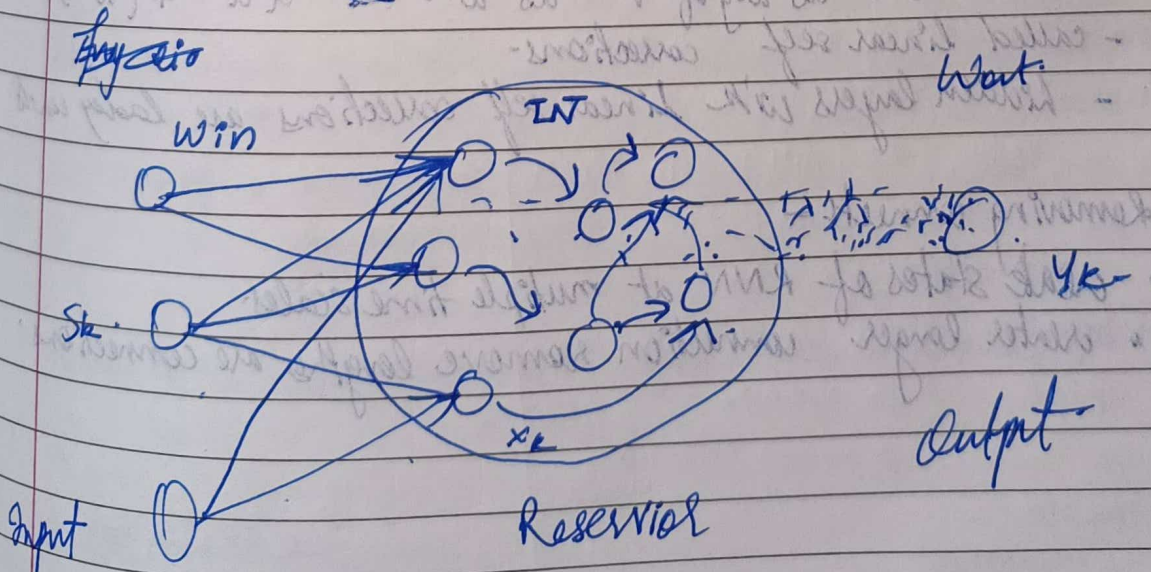→ Unfolding across time
→ Interconnection
→ Backprogation through time.

# ✱ Long Term Dependencies - (Challenges)

- Vanishing gradient → difficult to known which direction to move to reduce cost.
  - ✱ RRN tries to learn long range dependecies → gradient are small.
  - ✱ RRN not able to learn rel^n b/w input & output sequences.
  - ✱ becoms small that they are efficiently zero.
- Exploding gradient → make learning unstable
  - → gradient → large make RNN to diverge.
  - → RNN can't learn rel^n.
- RNN struggle as they forget past information as they process sequen

# ✱ Echo State networks (ESN) -

- type of RNN → easy to train & able to learn long term depende
- ability to handle sequencial data.
- series prediction, speech recognition & conbel systems
- Hidden layer or Reservoir is sparsely connected.

Adv → • easy to train  
       • learn long range dependecies  
       • relatively fast to run

Disadv → • sensetive to choice hyperamet  
       • difficult to interpret  
       • prone to overfiting -



Input   Win   IN   Work   $y_k$   Output

$S_k$   $x_k$   Output

Input   Reservior

# Leaky Units & strategies for mutiple time scales-

## Leaky units -

- type of RNN unit has linear self connection-
- fraction of output of unit is fed back into input of unit-
- ~~fraction of output is fed-bdck~~ → leak rate-
- learn mutiple time scales in RNNs-
- leak rate → how much info from past is forgotten-
- high leak rate → more info from past is forgotten- ←
- low leak rate → less info

## Strategies -

connection through

- Adding Skip Time -
  - Gradient explode wrt time steps $(\tau)$
  - Introduces time delay (d) now gradient dimension as a func^n of $\tau/d$ rather than $\tau$.
  - helps capture long term dependencies

- Leaky units & spectrum of time scales :-
  - Don't use integer skip of d instead use real valued $x$.
  - Consider $u(t)$ as avg of $v^{(t)}$ as $u^{(t)} \leftarrow \alpha u^{(t)} + (10-\alpha)v^{(t)}$
  - called linear self corrections-
  - hidden layers with linear self corrections are leaky units.

- Removing connect -
- create states of RNN at mutiple time scales-
- created longer connection remove lengths are connections-

# * Long Short Term Memory (LSTM)

- Type of RNN.
- designed to address vanishing gradient problem.
- handles long term dependencies in sequential data.

**\* Adv** → • captures long time dependencies.
- • addresses vanishing gradient problem
- • In various applications. → Machine Translation
  - → speech recognition
  - → Time series forecasting
  - → Text generation

**\* Disadv** → • computational complexity.
- • Prone to overfitting

**\* Working** ──→ Forget Gate, Inputs Gate, Cell state, Output Gate.

# * Gated Recurrent Unit (GRUs)

- Type of RNN
- process sequential data.
- handle long term dependency with data.

**Adv** → • simpler architecture      **Disadv** → limited capability
- • effective performance          → less control over inform" flow
- • eg → NLP, Speech reco.
  - Time Series Forecasting

**Working** → Input, Reset Gate, Update Gate, Hidden State Update.
  - Output.

# * Optimization

# ✳ Optimization for Long Term Dependencies-

- Using larger batch size —
- Using lower learning rate —
- Using gradient clipping
- Using dropout —
- Using 1, 2 regularization —
- Using gated RNN -

- weight initialization
- skip connection -
- gated architecture.
- regularization technique

# ✳ Explicit Memory

- use of seprate memory unit to store info from previous timestep
- contrast to implicit memory, info from prev steps is stored implicitly in weights of RNN
- Adv → • improve performance ( long range depencies )
  - • store info from previous steps - (even if not relevant)
  - • Reasoning & Problem solving-
- Disadv → • Increased Complexity
  - • Interpretability

ways of explicit memory implemented in RNNs-
  ⤷ Using seprate memory unit
  ⤷ Using a gated RNN.

| ✳ Performance Metrics | ✳ Default Baseline Model | ✳ Determing weather gather more data |
|---|---|---|
| • Accuracy $= \dfrac{\text{No of correct pred}}{\text{Total no of pred.}}$ | ⤷ Simple RNN | ⤷ Training Performance → Valid Perf |
| • Pecision $= \dfrac{TP}{TP+FP}$ | ⤷ LSTM | ⤷ Data Distribution → Task Complexity |
| • Recall $= \dfrac{TP}{TP+FN}$ | ⤷ GRU | ⤷ Task Complexity → Resource Constrts. |
| | ⤷ BiRNN | ⤷ Expert Knowledge |
| | ⤷ encoder Decoder (Seq) Model | ✳ selecting HyperParameters- |
| • Fscore $= 2 \times \dfrac{\text{precision} \cdot \text{recall}}{\text{Prec} + \text{recall}}$ | | • Learning Rate • No. of epochs |
| | | • Batch size. • No. of neurons |
| • MAE $= \dfrac{1}{n} \sum\limits_{i=1}^{n} |y_i - \hat{y_i}|$ | | • Activation fun⁴ |