

LEC-2 Recurrent Neural Networks

Deep Sequence Modeling

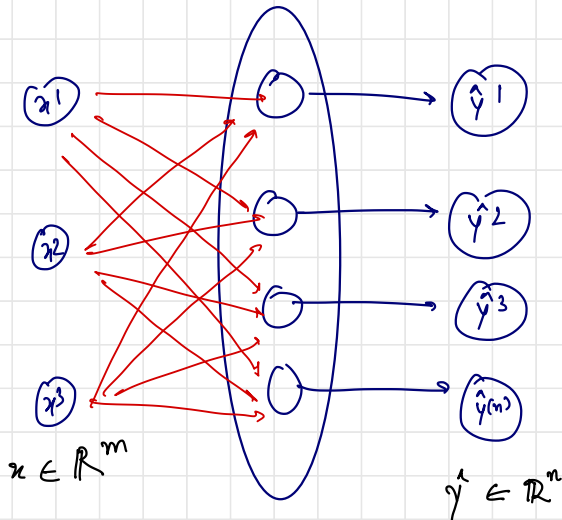
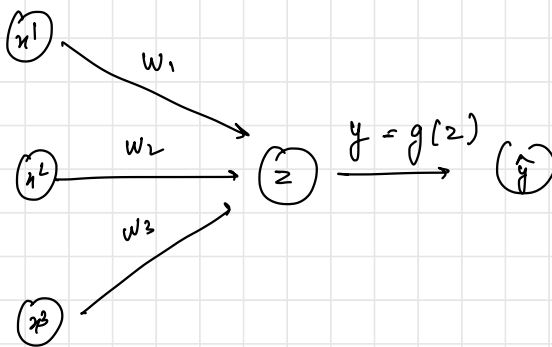
one to one binary classification pass fail

many to one — sentiment Classification

one to many — image captioning

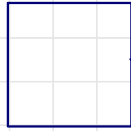
many to many — machine translation language

Perceptron



input vector

x_t

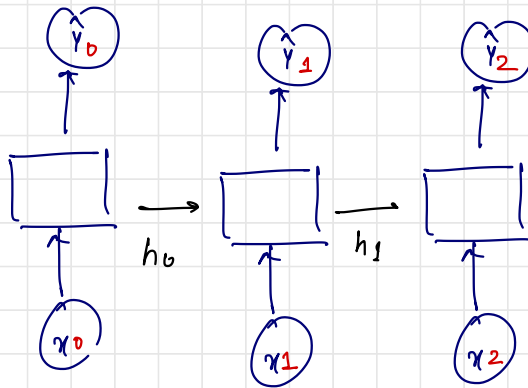
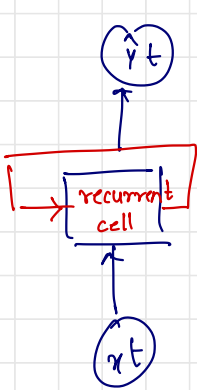


output vector

\hat{y}_t

$$x_t \in \mathbb{R}^m$$

$$\hat{y}_t \in \mathbb{R}^n$$

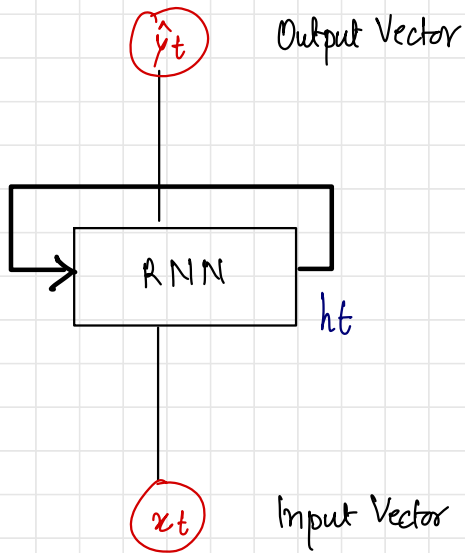


$$\hat{y}_t = f(x_t)$$

$$\hat{y}_t = f(\underbrace{x_t}_{\text{input}}, \underbrace{h_{t-1}}_{\text{past memory}})$$

internal memory - h_t

Recurrent Neural Networks (RNN)



Apply recurrence relation at every time step to process a sequence:

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

Annotations for the equation above:
- h_t : cell state
- f_w : function with weight
- x_t : input
- h_{t-1} : previous state

same function & set of parameters are used at every step

Output Vector

$$\hat{y}_t = W_{hy}^T h_t$$

↑

Update Hidden State

$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

↑

Input Vector $\rightarrow x_t$

Weight matrices reused every step

Design Criteria:-

1. Handle Variable length parameters
2. Track long term dependencies
3. Maintain information about order
4. Share parameters across sequence

Predict New Word

"This morning I took my cat for a walk

run	dog
walk	cat
day	happy
sun	sad

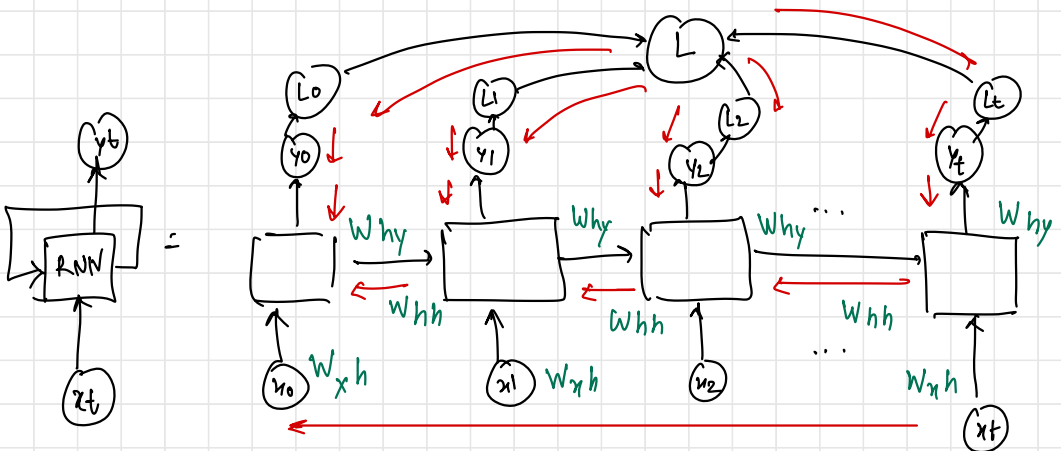
Representing language to a Neural Network

↳ into numerical inputs

encoding → Embedding transform indexes into a vector of fixed size
Vocabulary → Indexing → Embedding

Back propagation Through Time (BPTT)

1. Take the derivative (gradient) of the loss with respect to each parameter.
2. Shift parameters in order to minimize loss



Problems with computing gradient w.r.t h_0 involves many factors of

many values < 1 :
vanishing gradient

values > 1 :
exploding gradient

w_{hh} &
repeated gradient
computation

1. Activation function
2. Weight Initialization
3. Network Architecture

Trick #1 Activation functions

using ReLU prevents f' from shrinking the gradients when $x > 0$

Trick #2 Parameter Initialization

Init weights to identity matrix, biases to zero

→ helps prevent weights from shrinking to zero

Gated cell

use gates to selectively add or remove info
within each recurrent unit

LSTM - Long short term memory network

① forget ② Store ③ Update ④ Output

① maintain self state

② use gates to control flow of information

→ f

→ s

→ u

→ o

③ BPTT with partially uninterrupted gradient flow

Applications

- music generation
- sentiment classification

Limitations of ANN:-

- ① Encoding bottleneck
- ② slow, no parallelization
- ③ Not long memory

Transformers

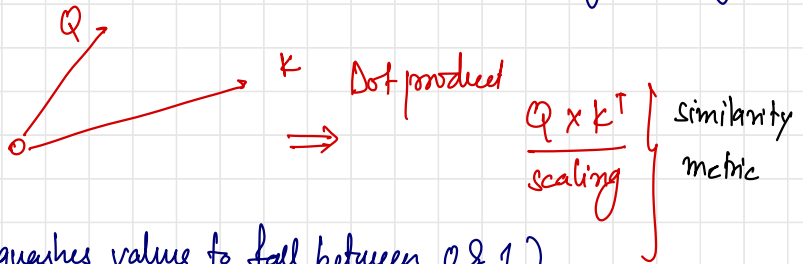
* Learning Self attention

1. Encode position information — aware encoding
2. Extract query, key, value for search

↳ transform by matrix multiplication with linear layers

3. Compute Attention weighting

Attention score: compute pairwise similarity between each query and key



4. Softmax (squeashes value to fall between 0 & 1)

$$\text{softmax}\left(\frac{Q \cdot K^T}{\text{scaling}}\right) \times V = A(Q, K, V)$$

└─┘
value matrix

query, key, value

