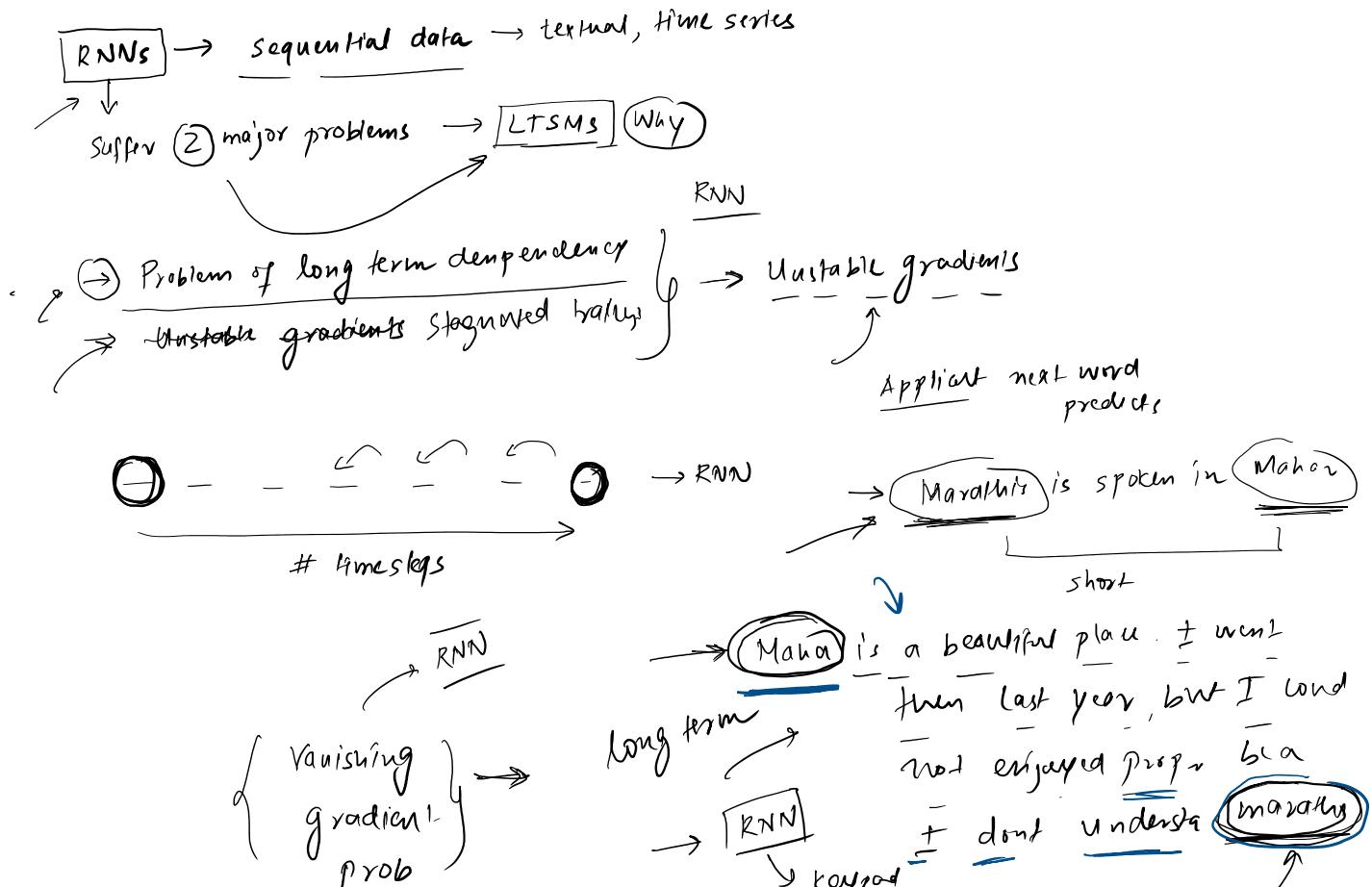


Problem with RNN

19 December 2022 16:33



Problem #1 → Problem of long term dependency → Vanishing

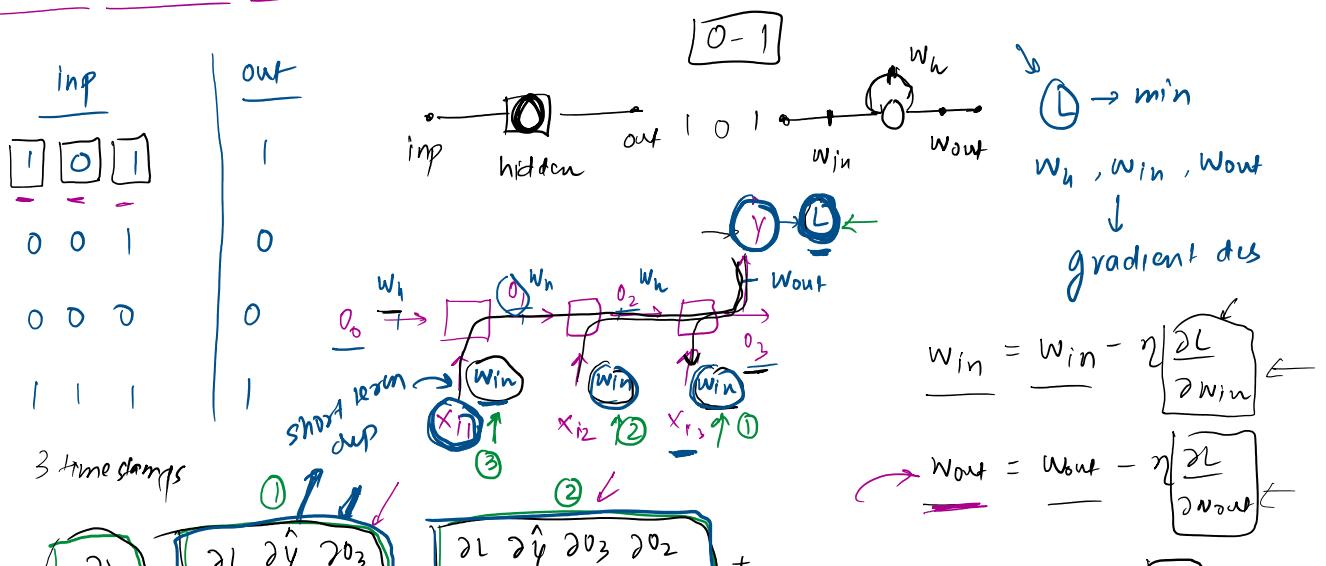


Diagram illustrating the backpropagation through time (BPTT) for an RNN, showing the computation of gradients for weights w_{in} and hidden states o_t .

The diagram shows the flow of gradients from the output layer (\hat{y}) back through the hidden states (o_1, o_2, \dots, o_{100}) to the input layer (x). It highlights the "long term dep" (long-term dependency) and the "inf" (infinity) issue due to the product of many small numbers.

Key equations shown:

- $\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_{in}} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_{in}} + \dots$
- $\frac{\partial L}{\partial o_t} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}}$
- $w_h = w_h - \eta \frac{\partial L}{\partial w_h}$
- $\# \rightarrow 3$ (number of time steps)
- $\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \prod_{t=2}^{100} \left(\frac{\partial o_t}{\partial o_{t-1}} \right) \frac{\partial o_1}{\partial w_{in}}$
- $o_t = \tanh(x_i w_{in} + o_{t-1} w_h)$
- $\frac{\partial o_t}{\partial o_{t-1}} = \tanh'(x_i w_{in} + o_{t-1} w_h) w_h$
- Vanishing gradient
- Identity matrix

Sol^④

- 1) Diff activation \rightarrow ReLU / Leaky ReLU
- 2) Better weight init
- 3) Skip layers
- 4) LSTM

Problem #2 \rightarrow Unstable Training (Exploding gradients)

- 1) Gradient Clipping
- 2) Controlled learning rate
- 3) LSTM

Recap

21 August 2023

11:55

