

Outline:

Fundamental principle: predict future with past experiences.

- **RNN:**
 - Architecture
 - Design criteria
- **Encoding:** character -> one-hot -> embedding
- **RNN backpropagation:**
 - problems:
 - gradient exploding
 - gradient vanishing
 - solution:
 - gradient clipping

- ReLU

prevents f' shrink
gradient

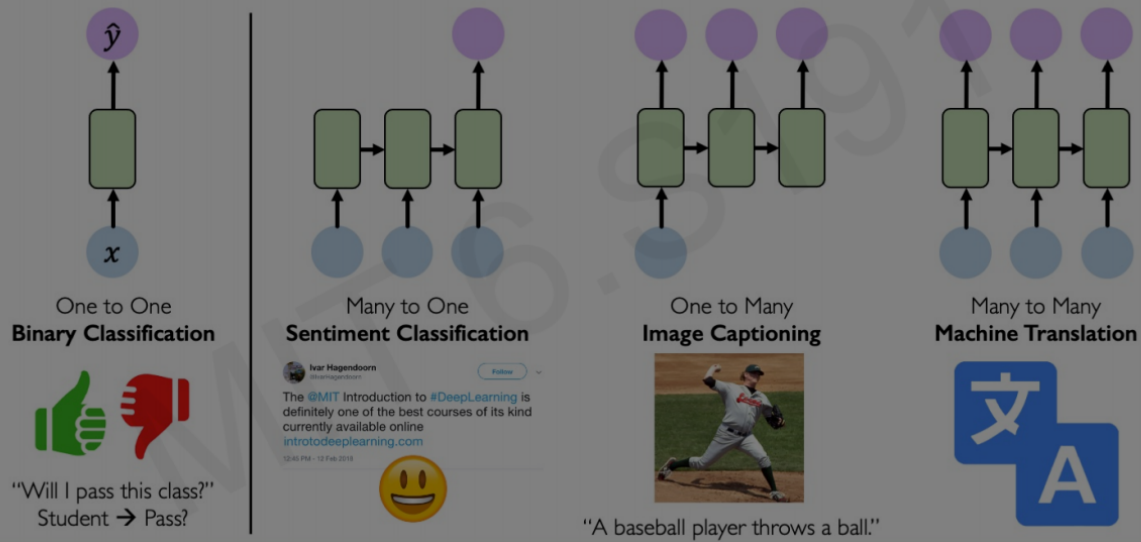
- parameter

initialization

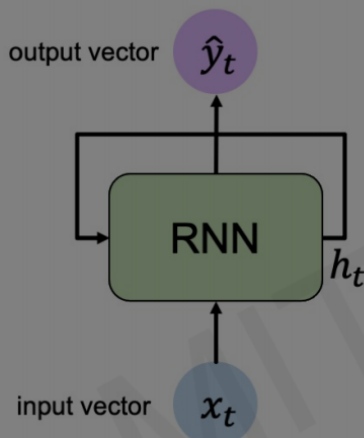
- **LSTM:** maintain a separate cell from what is outputted, use gate to control flow of information

- Forget
- Store
- Update
- Output

Sequence Modeling Applications



RNN State Update and Output



Output Vector

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

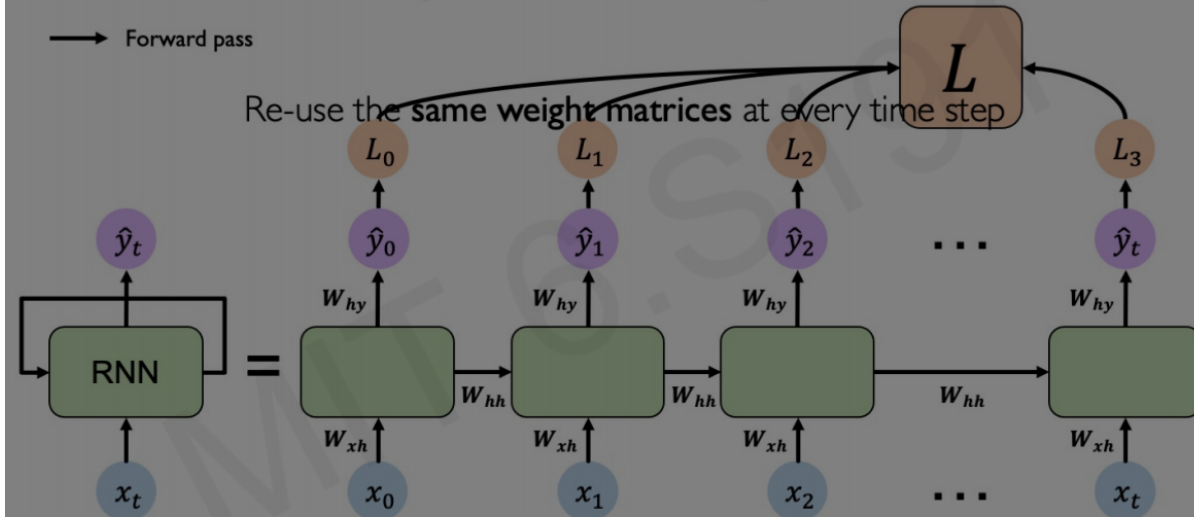
Update Hidden State

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

Input Vector

$$x_t$$

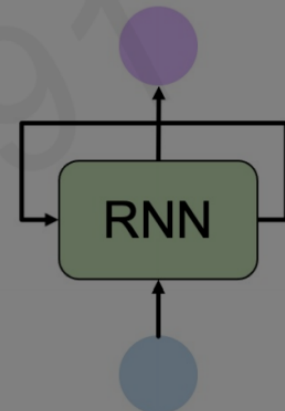
RNNs: Computational Graph Across Time



Sequence Modeling: Design Criteria

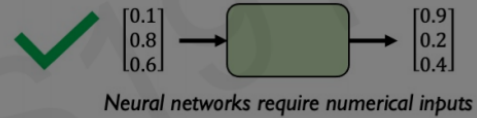
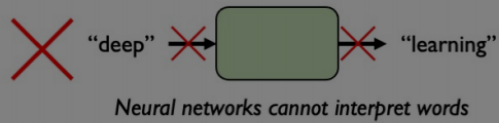
To model sequences, we need to:

1. Handle **variable-length** sequences
2. Track **long-term** dependencies
3. Maintain information about **order**
4. **Share parameters** across the sequence



Recurrent Neural Networks (RNNs) meet these sequence modeling design criteria

Encoding Language for a Neural Network



Embedding: transform indexes into a vector of fixed size.

this cat for
my took
a | walk
morning

1. Vocabulary:
Corpus of words

a → 1
cat → 2
...
walk → N

2. Indexing:
Word to index

One-hot embedding
"cat" = $[0, 1, 0, 0, 0, 0]$
↑
 i -th index

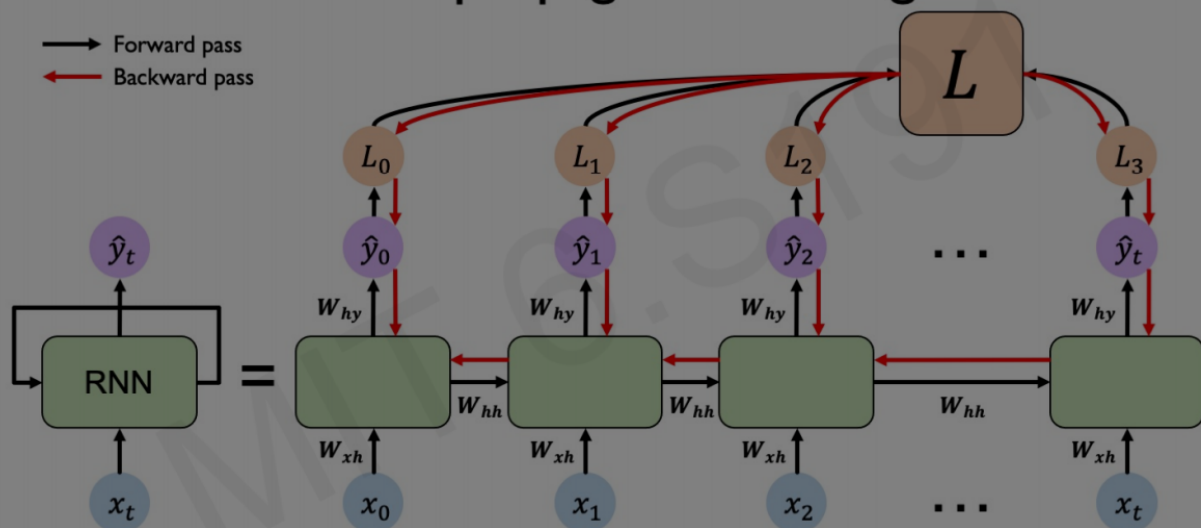
3. Embedding:
Index to fixed-sized vector

Learned embedding

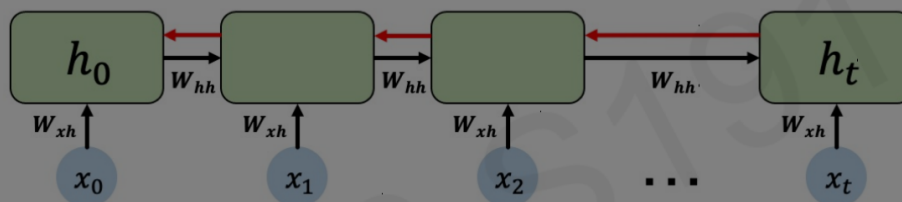
run walk dog cat
day sun happy sad

character → one-hot → word embedding.

RNNs: Backpropagation Through Time



Standard RNN Gradient Flow: Exploding Gradients

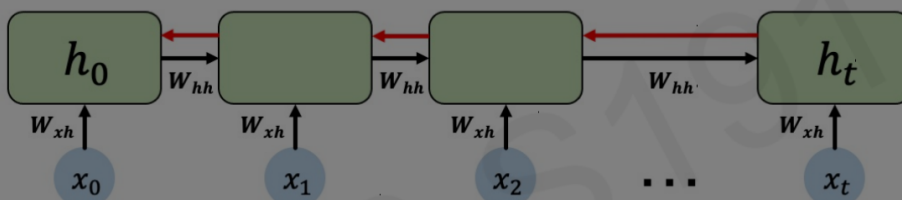


Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Standard RNN Gradient Flow: Vanishing Gradients



Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!

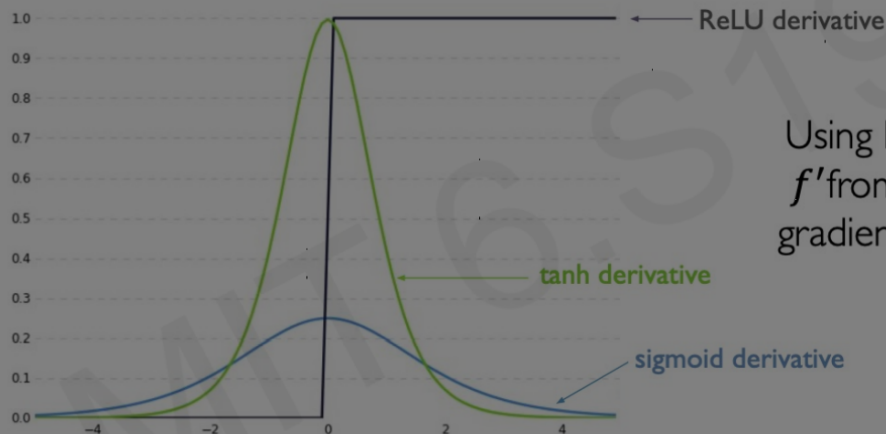
Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1 :
vanishing gradients

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

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

Solution #3: Gated Cells

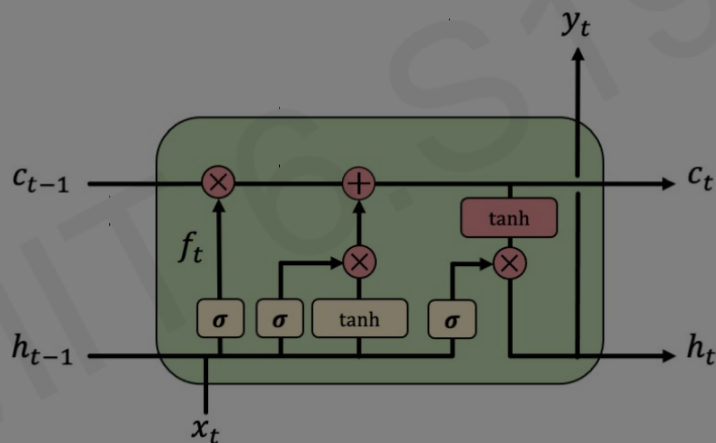
Idea: use a more **complex recurrent unit with gates** to control what information is passed through



Long Short Term Memory (LSTMs) networks rely on a gated cell to track information throughout many time steps.

Long Short Term Memory (LSTMs)

1) Forget 2) Store 3) Update 4) Output



LSTMs: Key Concepts

1. Maintain a **separate cell state** from what is outputted
2. Use **gates** to control the **flow of information**
 - **Forget** gate gets rid of irrelevant information
 - **Store** relevant information from current input
 - Selectively **update** cell state
 - **Output** gate returns a filtered version of the cell state
3. Backpropagation through time with **uninterrupted gradient flow**

Deep Learning for Sequence Modeling: Summary

1. RNNs are well suited for **sequence modeling** tasks
2. Model sequences via a **recurrence relation**
3. Training RNNs with **backpropagation through time**
4. Gated cells like **LSTMs** let us model **long-term dependencies**
5. Models for **music generation**, classification, machine translation, and more