



# CS F425: Deep Learning

# 24

# Sequence Modeling



**Dr. Kamlesh Tiwari**  
Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Mar 23, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

## Model and an application of RNN <sup>1</sup>

`rnn = RNN()`

`y = rnn.step(x)`

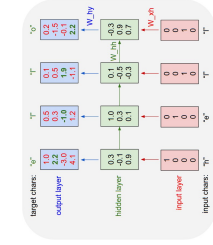
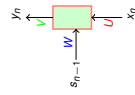
`class RNN:`

`def step(self, x):`

`self.s = np.tanh(np.dot(self.W, self.s) + np.dot(self.U, x))`

`y = np.dot(self.V, self.s)`

`return y`



### Character-Level Language Models

- Given a character what is the next character
- Characters are encoded using one-hot encoding
- Weights need to be adjusted

<sup>1</sup>karpathy.github.io: Minimal character-level Vanilla RNN model

## Training RNN

**Backpropogation** is used to train a RNN

- Total loss is the summation of  $J_t$  at every time step  $J = \sum_t J_t$
- To train a **parameter** we need to find its gradient with respect to loss and shift the parameter in its opposite direction

$$W = W - \alpha \frac{\partial J}{\partial W} = W - \alpha \frac{\partial}{\partial W} \sum_t J_t = W - \alpha \sum_t \frac{\partial J_t}{\partial W}$$

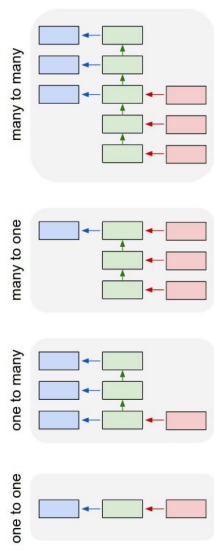
- Consider a single summation term

$$\frac{\partial J_t}{\partial W} = \frac{\partial J_t}{\partial s_t} \times \frac{\partial s_t}{\partial W}$$

- Term  $\frac{\partial s_t}{\partial W}$  is complicated. It depends on  $s_{t-1}, s_{t-2}, \dots, s_1$

Various arrangements are possible based on need

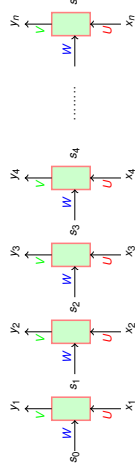
- What should be the next world? **One-to-One**
- Caption the given image? **One-to-Many**
- Segmentation or classification **Many-to-One**
- Translate from one language to other **Many-to-Many**



**Output may be different for same input**

We are optimizing over programs not on functions

## Training: A Simplified Version of RNN



$$y_t = \sigma(V^T s_t) \quad s_t = \sigma(U^T x_t + W^T s_{t-1})$$

If loss at time  $t$  be  $J_t$

Then total loss is  $J = \sum_t J_t$

$J_t$  could be something like  $-\log(\text{probability of true output})$

How to train RNN?

## Training RNN (contd..)

- As  $s_t = \sigma(U^T x_t + W^T s_{t-1})$  where  $s_t$  also depends on  $W$  and other previous states  $s_{t-2}, \dots, s_1$
- Its derivative has **explicit** and **implicit** parts. Explicit part considers other things as constant (we represent as  $\partial^+$ )

$$\begin{aligned} \frac{\partial s_t}{\partial W} &= \frac{\partial^+ s_t}{\partial W} + \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial W} \\ &= \frac{\partial^+ s_t}{\partial W} + \frac{\partial s_t}{\partial s_{t-1}} \left[ \frac{\partial^+ s_{t-1}}{\partial W} + \frac{\partial s_{t-1}}{\partial s_{t-2}} \frac{\partial s_{t-2}}{\partial W} \right] \\ &= \frac{\partial^+ s_t}{\partial W} + \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial^+ s_{t-1}}{\partial W} + \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial^+ s_{t-1}}{\partial s_{t-2}} \left[ \frac{\partial^+ s_{t-2}}{\partial W} + \frac{\partial s_{t-2}}{\partial s_{t-3}} \frac{\partial s_{t-3}}{\partial W} \right] \\ &= \sum_{k=1}^t \frac{\partial s_t}{\partial s_k} \frac{\partial^+ s_k}{\partial W} \end{aligned}$$

where we use  $\frac{\partial s_t}{\partial s_k}$  as a short form for  $\frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial s_{t-2}} \frac{\partial s_{t-2}}{\partial s_{t-3}} \dots \frac{\partial s_{k+2}}{\partial s_{k+1}} \frac{\partial s_{k+1}}{\partial s_k}$

## Exploding and Vanishing Gradient

Consider  $\frac{\partial s_t}{\partial s_k}$  that is  $\frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial s_{t-2}} \frac{\partial s_{t-2}}{\partial s_{t-3}} \dots \frac{\partial s_{k+2}}{\partial s_{k+1}} \frac{\partial s_{k+1}}{\partial s_k} = \prod_{j=t-k}^{t-1} \frac{\partial s_j}{\partial s_{j-1}}$

**Focus on a single term**  $\frac{\partial s_t}{\partial s_{t-1}}$   
 it can be shown that it is upper-bounded by some constant (which depends on  $W$ ) let's call that value  $c$

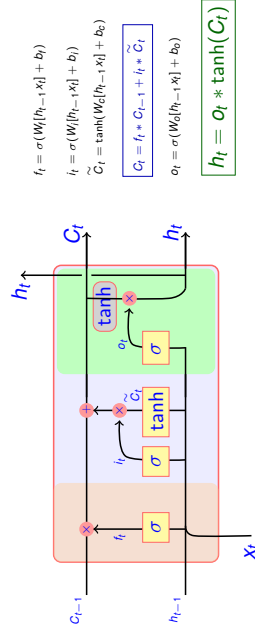
Since  $\frac{\partial s_t}{\partial s_k}$  is a multiplication of  $t - k$  such terms (and we expect  $t - k$  to be large to address long term dependency) therefore the value

$$\frac{\partial s_t}{\partial s_k} \rightarrow \begin{cases} \text{Vanish} & \text{if } c < 1 \\ \text{Explode} & \text{if } c > 1 \end{cases}$$

(0.90 → -0.81 → -0.73 → -0.66 → -0.59 → -0.53 → -0.46 → -0.39) (1.50 → 2.25 → 3.38 → 5.06 → 7.59 → 11.39 → 17.09 → 25.63)

## LSTM: Long Short Term Memory 2

LSTM are RNN with 4 special NN circuits for **forget**, **store** and **output**



$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C[h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

<sup>2</sup> Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 9.8 (1997): 1785-1780.

## Example-01 [1/2]

Five numbers are fed and we want the 2nd one as output.

```
from random import randint
from numpy import array
from keras.models import Sequential
from keras.layers import LSTM
from keras.layers import Dense

def generate_sequence(length, n_features):
    return [randint(0, n_features-1) for _ in range(length)]

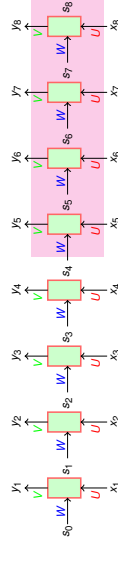
def encode(sequence, n_features):
    encoding = list()
    for value in sequence:
        vector = [0 for _ in range(n_features)]
        vector[value] = 1
        vector = array(vector)
        return array(encoding)

def one_hot_decode(encoded_seq):
    return [argmax(vector) for vector in encoded_seq]

def generate_example(length, n_features, out_index):
    sequence = generate_sequence(length, n_features)
    encoded = one_hot_encode(sequence, n_features)
    X = encoded.reshape(1, length, n_features)
    y = encoded[out_index].reshape(1, n_features)
    return X, y
```

## How to Avoid Exploding and Vanishing Gradient

- Clipping**  
 Try normalizing value of  $c$  keeping it some range  $T_l \leq c \leq T_h$
- Truncated Backpropagation**  
 At any step look for only for last  $k$  timestamps



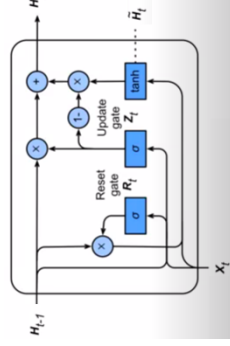
## LSTM

Special circuits for handling long term dependencies

## Gated Recurrent Unit (GRU) 3

Get most of the advantage of LSTM with less computation

- Attention mechanism towards when to **reset** and **update**



$$r_t = \sigma(W_r[h_{t-1}, x_t] + b_r)$$

$$z_t = \sigma(W_z[h_{t-1}, x_t] + b_z)$$

$$\tilde{h}_t = \tanh(W_{\tilde{h}}[h_{t-1}, x_t] + b_{\tilde{h}})$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

<sup>3</sup> cite: 164770 Learning phrase representations using RNN encoder-decoder for statistical machine translation, Cho, Kyunghyun and Van M. Bart and Gulcehre, Caglar and Bahdanau, Dzmitry and Bougares, Fethi and Schwenk, Holger and Bengio, Yoshua, arXiv preprint arXiv:1406.1078, (2014)

## Example-01 [2/2]

```
### Defining the model #####
length = 5
n_features = 10
out_index = 2
model = Sequential()
model.add(LSTM(25, input_shape=(length, n_features)))
model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['acc'])
model.summary()

#### Fitting the model #####
for x, y in generate_example(length, n_features, out_index):
    model.fit([X, y], epochs=1, verbose=2)

#### Evaluating the model #####
correct = 0
for i in range(100):
    X, y = generate_example(length, n_features, out_index)
    if one_hot_decode(yhat) == one_hot_decode(y):
        correct += 1
print('Accuracy: %s' % ((correct/100)*100.0))
```

The system achieves 100% accuracy



Thank You!

Thank you very much for your attention!