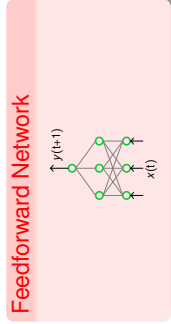
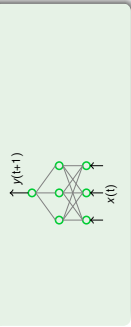


Recurrent Networks (RNN)

Feedforward network cannot capture the dependence of $y(t + 1)$ on earlier values of x such as $x(t - 1)$



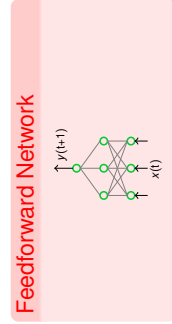
Feedforward Network



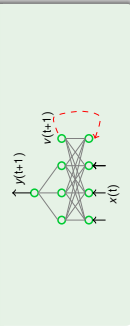
Recurrent Network

Recurrent Networks (RNN)

Feedforward network cannot capture the dependence of $y(t + 1)$ on earlier values of x such as $x(t - 1)$



Feedforward Network

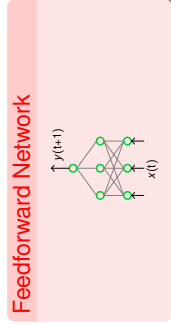


Recurrent Network

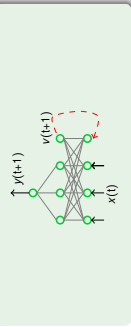
- RNN uses **states** (called self-state) of the **network units** available at **time t** as an input to the other units at **time $t + 1$**

Recurrent Networks (RNN)

Feedforward network cannot capture the dependence of $y(t + 1)$ on earlier values of x such as $x(t - 1)$



Feedforward Network

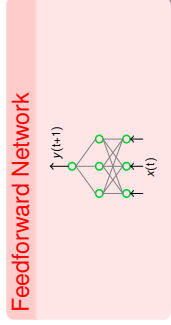


Recurrent Network

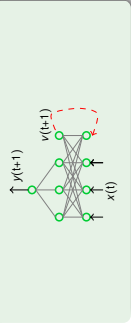
- RNN uses **states** (called self-state) of the **network units** available at **time t** as an input to the other units at **time $t + 1$**
- RNN is suitable for temporal data (like time series)
- Training may involve unfolding and averaging.

Recurrent Networks (RNN)

Feedforward network cannot capture the dependence of $y(t + 1)$ on earlier values of x such as $x(t - 1)$



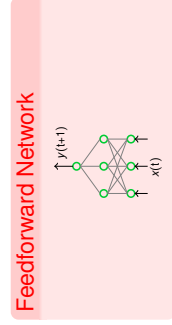
Feedforward Network



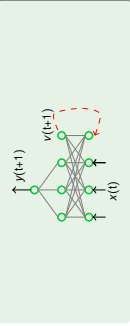
Recurrent Network

Recurrent Networks (RNN)

Feedforward network cannot capture the dependence of $y(t + 1)$ on earlier values of x such as $x(t - 1)$



Feedforward Network

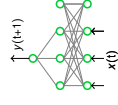


Recurrent Network

- RNN uses **states** (called self-state) of the **network units** available at **time t** as an input to the other units at **time $t + 1$**
- RNN is suitable for temporal data (like time series)

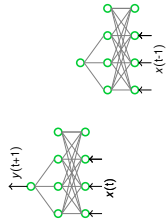
Unfolding RNN in Time

When input at the time t is provided, **what is output at time $(t + 1)$?**



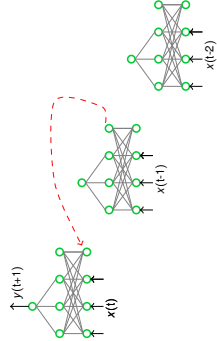
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



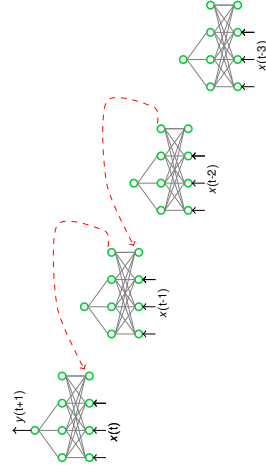
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



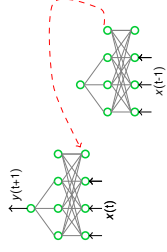
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



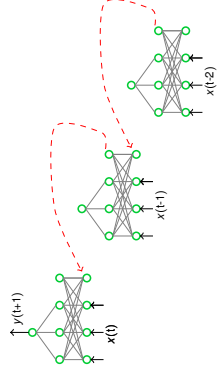
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



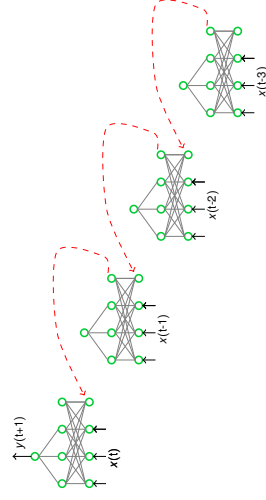
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



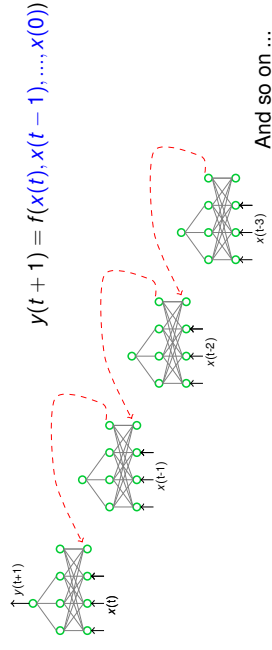
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



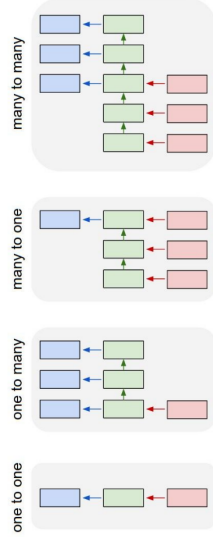
Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?



Various arrangements are possible based on need

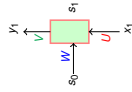
- What should be the next world? **One-to-One**
- Caption the given image? **One-to-Many**
- Segmentation or classification **Many-to-One**
- Translate from one language to other **Many-to-Many**



Output may be different for same input

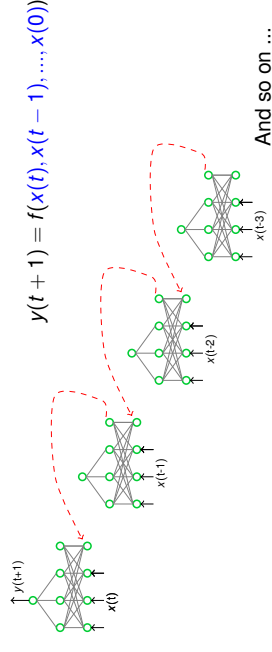
We are optimizing over programs not on functions

Training: A Simplified Version of RNN



Unfolding RNN in Time

When input at the time t is provided, what is output at time $(t + 1)$?

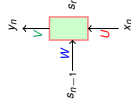


By this way RNN incorporates history of the network in output

Model and an application of RNN ¹

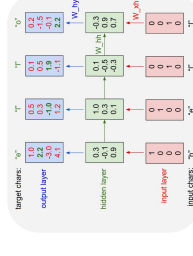
```

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



¹karpathy.github.io: Minimal character-level Vanilla RNN model

Training: A Simplified Version of RNN

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 $s_{t-1}, s_{t-2}, \dots, s_1$

Exploding and Vanishing Gradient

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

Focus on a single term $\frac{\partial s_j}{\partial s_{j-1}}$

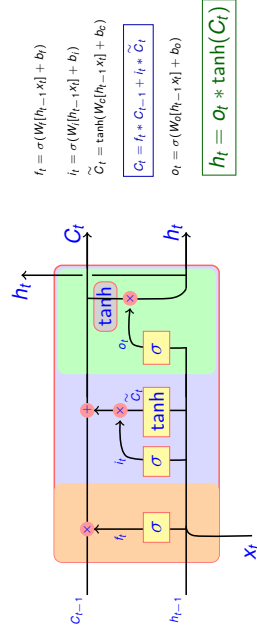
it can be shown that it is upper-bounded by some constant (which depends on W) let's call that value c

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

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

LSTM: Long Short Term Memory ²

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



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

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

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

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

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

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

²Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 8.8 (1997): 1738-1760.

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 ∂^+)

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

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

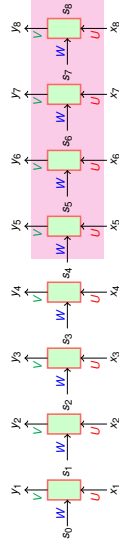
How to Avoid Exploding and Vanishing Gradient

1 Clipping

Try **normalizing** value of c keeping it some range $\bar{T}_i \leq c \leq Th$

2 Truncated Backpropogation

At any step look for only for last k timestamps



3 LSTM

Special circuits for handling long term dependencies

Example-01 [1/2]

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

```
from random import randint
from numpy import array
import keras
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)]

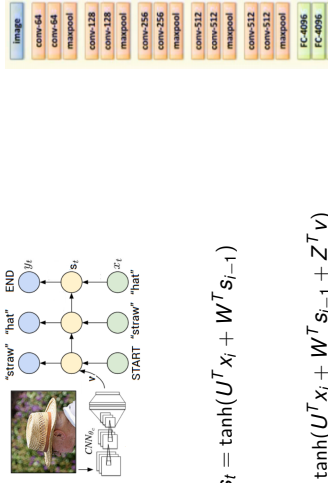
def one_hot_encode(sequence, n_features):
    encoding = list()
    for value in sequence:
        vector = [0 for _ in range(n_features)]
        vector[value] = 1
        encoding.append(vector)
    return array(encoding)

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
```


Example-05: Image Captioning ⁷

Thank You!



Before:

$$s_t = \tanh(U^T x_t + W^T s_{t-1})$$

Now:

$$s_t = \tanh(U^T x_t + W^T s_{t-1} + Z^T v)$$

⁷ Karpathy, Andrej, and Li Fei-Fei. "Deep visual-semantic alignments for generating image descriptions." Proceedings of the IEEE conference on computer vision and pattern recognition, 2015.