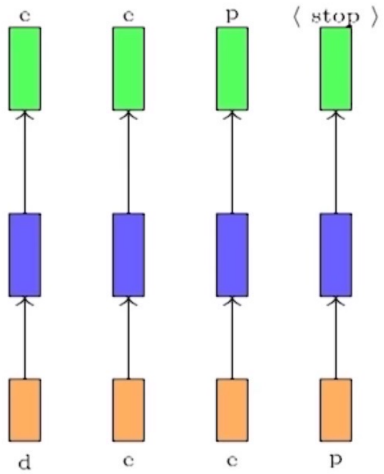# Sequence Learning Problems

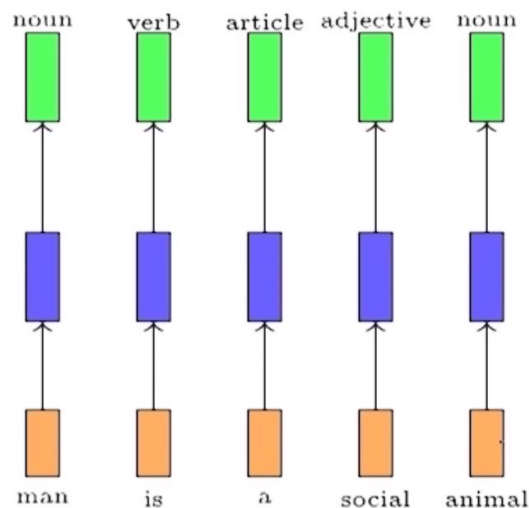1. Successive inputs may not be independent of each other.
2. The input size is not fixed.

## Text Completion
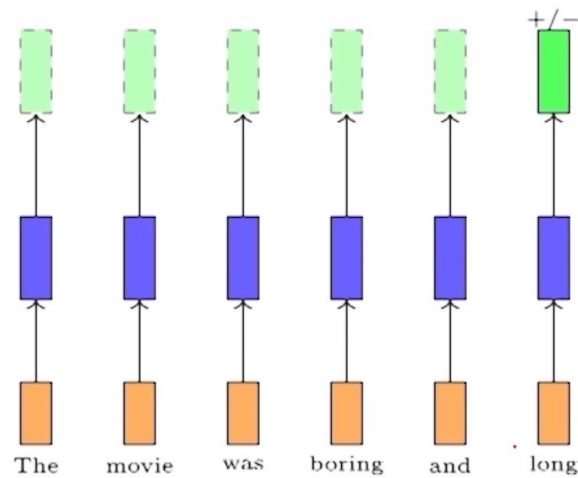


(i) Input is no longer independent.

(ii) Length of word (i/p) is not fixed.

(iii) Each network is performing the same task.

## Part of Speech Tagging

The    movie    was    boring    and    long
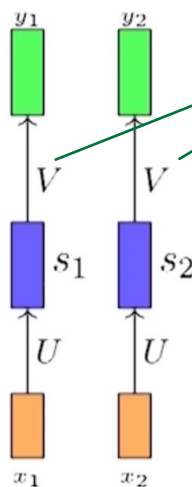
## Objectives :

1. Account for the dependence between inputs.

2. Account for variable size of input

3. Make sure that the function executed at each time step is the same.

bias
↓

$$s_1 = \sigma(U x_1 + b)$$
$$y_1 = 0(V s_1 + c)$$
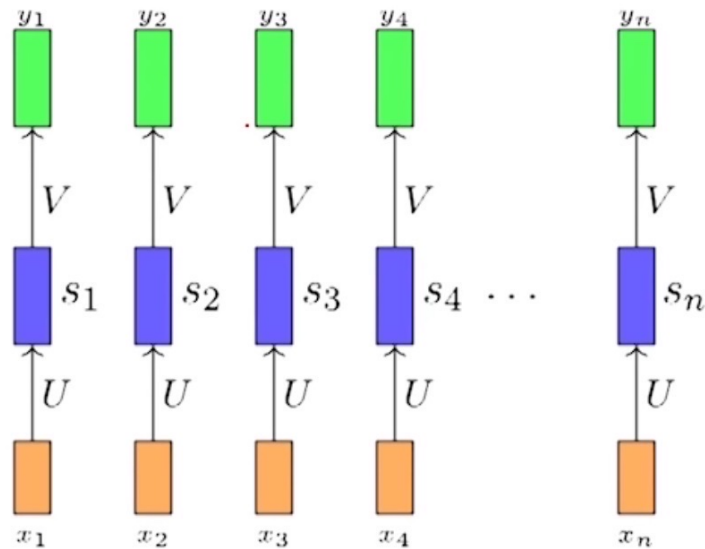← bias

0: output function



Share the same weight accross different time-steps.
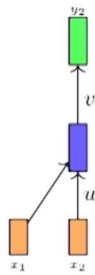
$$s_2 = \sigma(U x_2 + b)$$
$$y_2 = 0(V s_2 + c)$$

Since we want the same function to be implemented at each step, we should share the same network parameters.
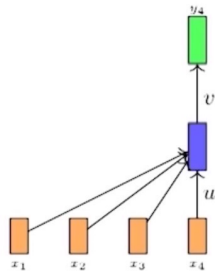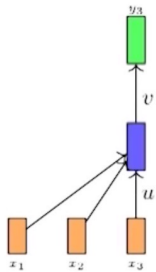
Variable size of the input can be taken care of by replicating the network at each time step.



How do we account for the dependence between the inputs?



Not efficient

<u>Solution</u>: Add a recurrent connection



Recurrent
connection

$$s_2 = \sigma\left(Ux_2 + Ws_1 + b\right)$$
$$y_2 = O\left(Vs_2 + c\right)$$

(Unrolled version)

$$s_i = \sigma\left(Ux_i + Ws_{i-1} + b\right)$$
$$y_i = O\left(Vs_i + c\right)$$

$s_i$: State of the network

# How to train an RNN?

Backpropagation through time.



$$y_1 \qquad y_2 \qquad y_2 \qquad y_2$$

$V \qquad V \qquad V \qquad V$

$W \qquad W \qquad W \qquad W$

$U \qquad U \qquad U \qquad U$

$x_1 \qquad x_2 \qquad x_2 \qquad x_2$

# What is a suitable output function

In this case: softmax



| | Predicted | Predicted | Predicted | Predicted |
|---|---|---|---|---|
| d | 0.2 | 0.2 | 0.2 | 0.2 |
| e | 0.7 | 0.7 | 0.1 | 0.1 |
| p | 0.1 | 0.1 | 0.7 | 0.7 |
| stop | 0.1 | 0.1 | 0.1 | 0.1 |

$V \qquad V \qquad V \qquad V$

$W \qquad W \qquad W$

$U \qquad U \qquad U \qquad U$

d \qquad e \qquad e \qquad e

$\mathcal{L}_1(\theta):$ cross-entropy

$\hat{y}$     $y$    $\mathcal{L}_2(\theta)$       $\mathcal{L}_3(\theta)$    $\cdot\mathcal{L}_4(\theta)$

| | Predicted | True | Predicted | True | Predicted | True | Predicted | True |
|---|---|---|---|---|---|---|---|---|
| d | 0.2 | 0 | 0.2 | 0 | 0.2 | 0 | 0.2 | 0 |
| .e | 0.7 | 1 | 0.7 | 1 | 0.1 | 0 | 0.1 | 0 |
| p | 0.1 | 0 | 0.1 | 0 | 0.7 | 1 | 0.7 | 1 |
| stop | 0.1 | 0 | 0.1 | 0 | 0.1 | 0 | 0.1 | 0 |

$$\min \mathcal{L}(\theta) = \sum_{t=1}^{4} \mathcal{L}_t(\theta)$$

Do backpropagation

$V$   $V$   $V$   $V$

$W$   $W$   $W$

$U$   $U$   $U$   $U$

d   e   e   e

## <span style="color:red">Problem of Vanishing and Exploding gradients</span>

$y_1$   $y_2$   $y_2$   $y_2$

$V$   $V$   $V$   $V$

$W$   $W$   $W$   $W$

$U$   $U$   $U$   $U$
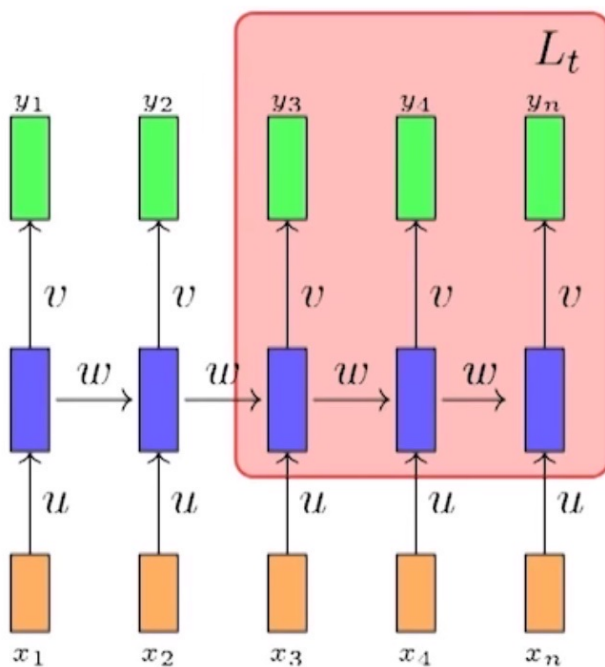
$x_1$   $x_2$   $x_2$   $x_2$

$W$    $W$    $W$

When gradient information is passed backward it's multiplied with $W$.

(1) <u>Exploding Gradients</u>: Gradient value starts to increase very rapidly

roughly when $\|W\| > 1$      $(eg. \ 2^t)$

(2) <u>Vanishing gradients</u>: Gradient value starts to decrease very rapidly.

roughly when $\|w\| < 1$ $\left( eg. \ \frac{1}{2^t} \right)$
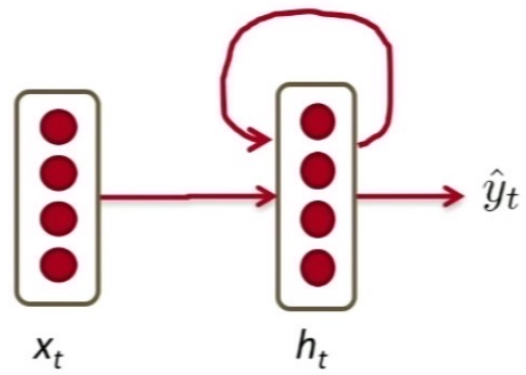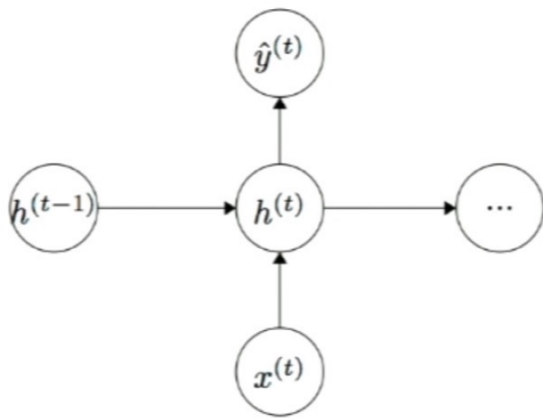
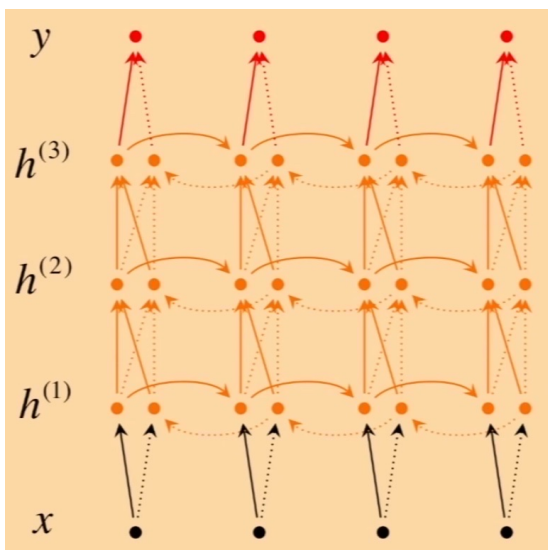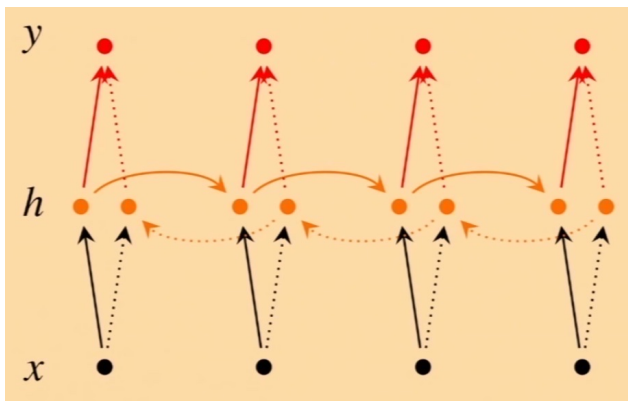One simple solution could be Truncated Backpropagation restrict the product to $T$-terms.



<u>Another solution</u>: Clipping

If gradient > Threshold $T$ then

Gradient = $T$

## Bi-directional RNN





Multiple hidden layers