

# STA 221: LECTURE 14

KRISHNA BALASUBRAMANIAN

(UNIVERSITY OF CALIFORNIA, DAVIS)

# Recurrent Neural Network

## TIME SERIES/SEQUENCE DATA

Input:  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$

Each  $\mathbf{x}_t$  is the feature at time step  $t$

Each  $\mathbf{x}_t$  can be an  $d$ -dimensional vector

Output:  $\{y_1, y_2, \dots, y_T\}$

Each  $y_t$  is the output at step  $t$

Multi-class output or Regression output:

$$y_t \in \{1, 2, \dots, L\} \quad \text{or} \quad y_t \in \mathbb{R}$$

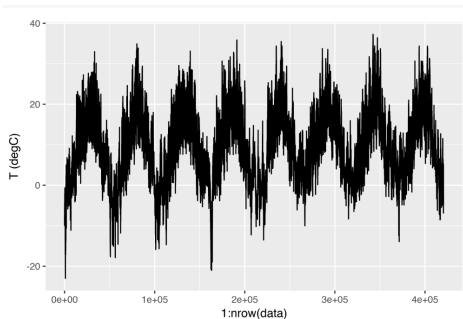
## EXAMPLE: TIME SERIES PREDICTION

Climate Data:

$x_t$ : temperature at time  $t$

$y_t$ : temperature (or temperature change) at time  $t + 1$

Stock Price: Predicting stock price



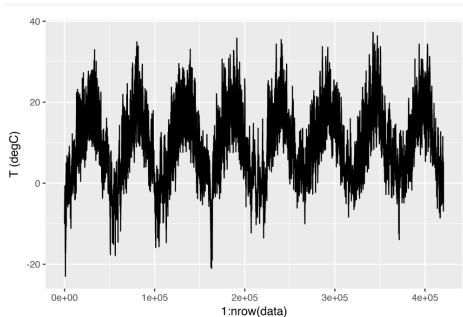
## EXAMPLE: TIME SERIES PREDICTION

Climate Data:

$x_t$ : temperature at time  $t$

$y_t$ : temperature (or temperature change) at time  $t + 1$

Stock Price: Predicting stock price



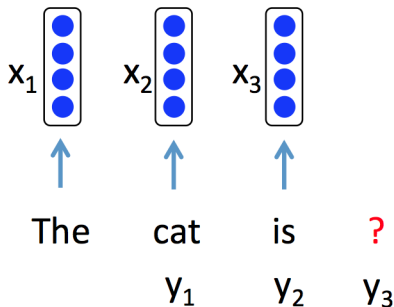
## EXAMPLE: LANGUAGE MODELING

The    cat    is    ?

$\mathbf{x}_t$ : one-hot encoding to represent the word at step  $t$   
([0, ..., 0, 1, 0, ..., 0])

$y_t$ : the next word

$y_t \in \{1, \dots, V\}$   $V$ : Vocabulary size



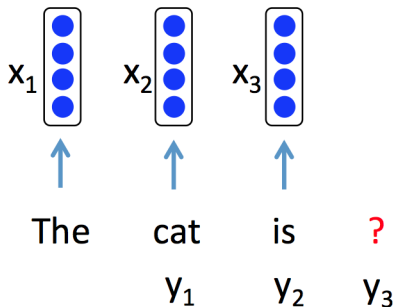
## EXAMPLE: LANGUAGE MODELING

The    cat    is    ?

$\mathbf{x}_t$ : one-hot encoding to represent the word at step  $t$   
([0, ..., 0, 1, 0, ..., 0])

$y_t$ : the next word

$y_t \in \{1, \dots, V\}$     $V$ : Vocabulary size



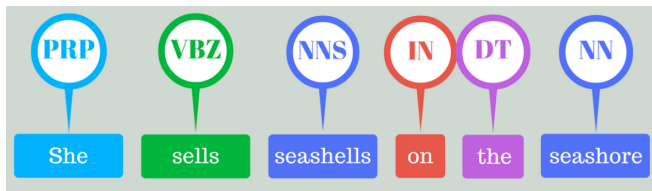
## EXAMPLE: POS TAGGING

Part of Speech Tagging:

Labeling words with their Part-Of-Speech (Noun, Verb, Adjective, ...)

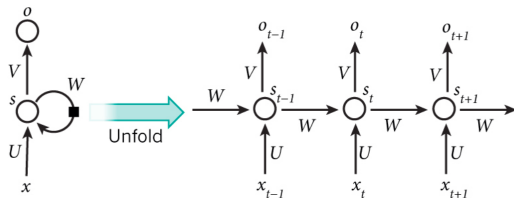
$\mathbf{x}_t$ : a **vector** to represent the word at step  $t$

$y_t$ : label of word  $t$





## RECURRENT NEURAL NETWORK (RNN)



$\mathbf{x}_t$ :  $t$ -th input

$\mathbf{s}_t$ : hidden state at time  $t$  (“memory” of the network)

$$\mathbf{s}_t = f(U\mathbf{x}_t + W\mathbf{s}_{t-1})$$

$W$ : transition matrix  $\mathbf{s}_0$  usually set to be 0

Predicted output at time  $t$ :

$$o_t = \arg \max_i (V\mathbf{s}_t)_i$$

## RECURRENT NEURAL NETWORK (RNN)

Training: Find  $U, W, V$  to minimize empirical loss:

Loss of a sequence:

$$\sum_{t=1}^T \text{loss}(V \mathbf{s}_t, y_t)$$

( $\mathbf{s}_t$  is a function of  $U, W, V$ )

Loss on the whole dataset:

Average loss of all sequence

Solve by Stochastic Gradient Descent (SGD)

## RECURRENT NEURAL NETWORK (RNN)

Training: Find  $U, W, V$  to minimize empirical loss:

Loss of a sequence:

$$\sum_{t=1}^T \text{loss}(V \mathbf{s}_t, y_t)$$

( $\mathbf{s}_t$  is a function of  $U, W, V$ )

Loss on the whole dataset:

Average loss of all sequence

Solve by Stochastic Gradient Descent (SGD)

## RECURRENT NEURAL NETWORK (RNN)

Training: Find  $U, W, V$  to minimize empirical loss:

Loss of a sequence:

$$\sum_{t=1}^T \text{loss}(V \mathbf{s}_t, y_t)$$

( $\mathbf{s}_t$  is a function of  $U, W, V$ )

Loss on the whole dataset:

Average loss of all sequence

Solve by Stochastic Gradient Descent (SGD)

## RNN: TEXT CLASSIFICATION

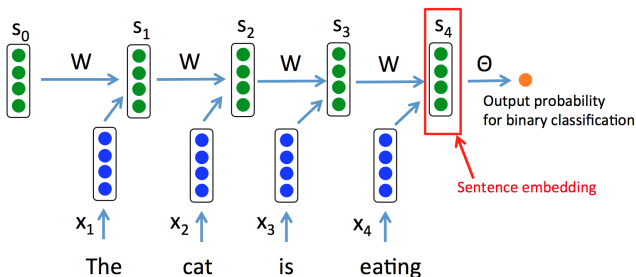
Not necessary to output at each step

Text Classification:

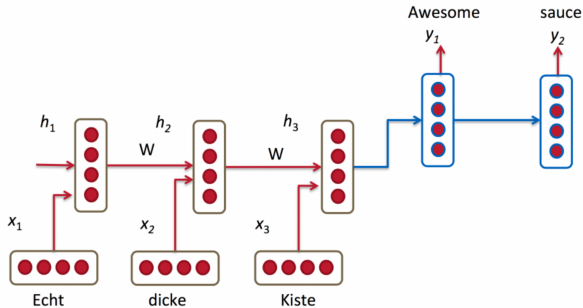
Sentence  $\rightarrow$  category

Output only at the final step

Model: add a fully connected network to the final embedding



# RNN: NEURAL MACHINE TRANSLATION



## PROBLEMS OF CLASSICAL RNN

Hard to capture **long-term dependencies**

Hard to solve (vanishing gradient problem)

Solution:

- LSTM (Long Short Term Memory networks)

- GRU (Gated Recurrent Unit)

- ...