**NUS-ISS**
*Pattern Recognition using Machine Learning System*

# Module 7 - Solving temporal sequential problems using recurrent neural networks
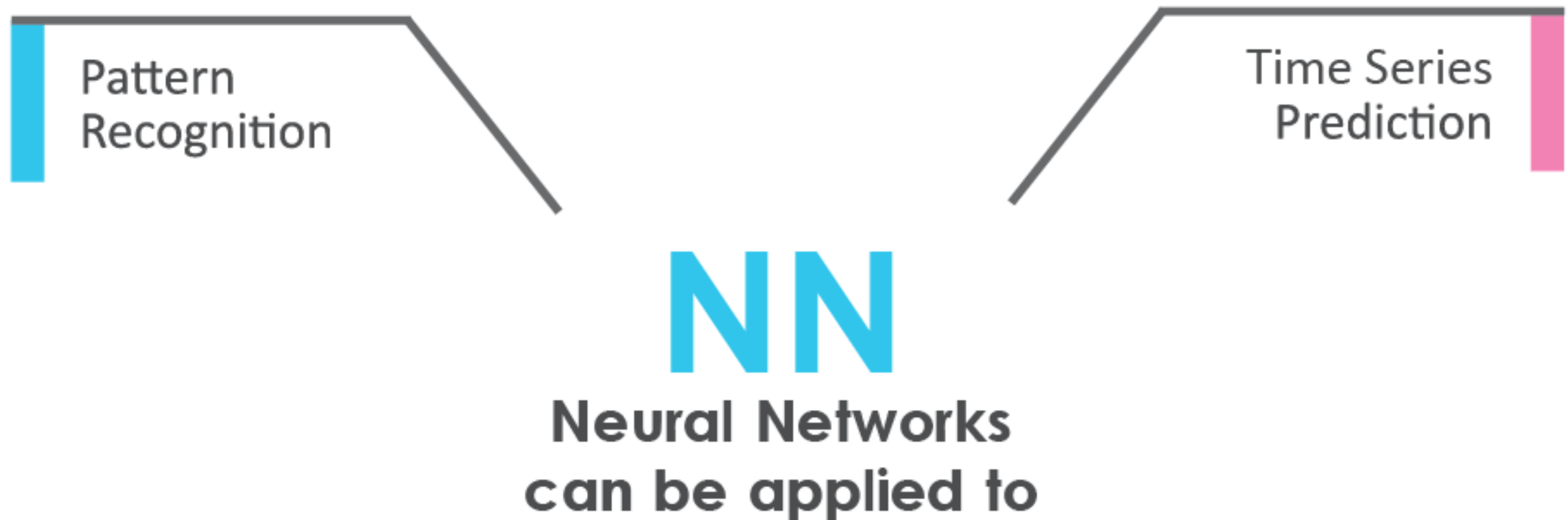
by Dr. Tan Jen Hong

# When time is a factor

# The other application

with time

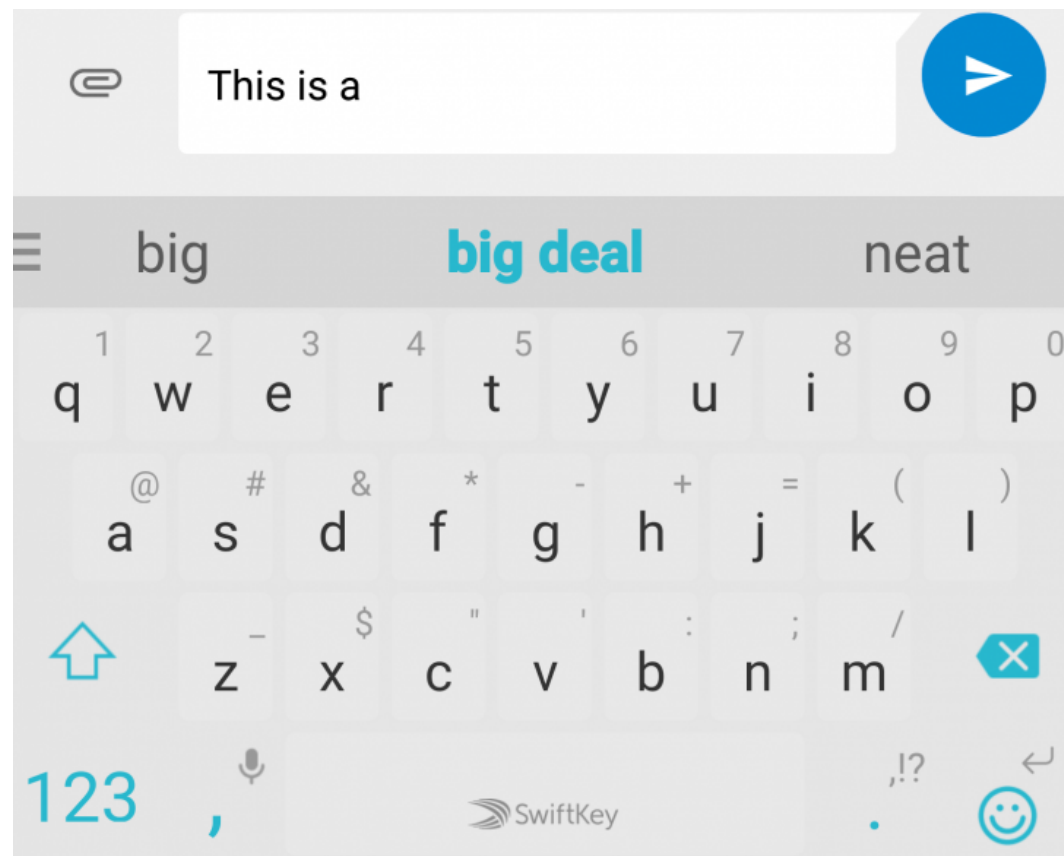Pattern Recognition

Time Series Prediction

## NN
### Neural Networks
### can be applied to

Source: http://www.cvisiontech.com/resources/ocr-primer/ocr-neural-networks-and-other-machine-learning-techniques.html

prumls/m3.2/v1.0 NUS | ISS

# The related inputs
## Not independent

- So far the nets introduced assume inputs are independent from each other

- Implication: the inputs that came before and the inputs that will come after has no relationship

- But for some tasks/problems, this is not true

- E.g. if. you want to predict which word to come in a sentence, you better know what have been typed/ said before

prumls/m3.2/v1.0

# Recurrent neural network

In short, RNN

- Recurrent neural network: a net that perform the same calculation on elements/segments from a sequence

- The output of a current element/ segment depends on the outputs from the previous elements/ segments



'This'    'is'    'a'

prumls/m3.2/v1.0    NUS National University of Singapore | iss INSTITUTE OF SYSTEMS SCIENCE

# CNN vs RNN
## Comparison

| | CNN | RNN / LSTM |
|---|---|---|
| **Usage** | Suitable for spatial data, e.g. image, video | Suitable for temporal data (sequential data), e.g. text, speech |
| **Capability** | Considered more powerful than RNN | Less powerful and slower in calculation |
| **Input** | Take fixed size inputs and generate fixed size outputs | Can handle arbitrary input/output lengths |
| **Nature** | Use local connectivity pattern (through 2D convolution) | Use time series information |

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

# Time step
Prediction

0, 1, 1, 3, 2, 2, 4, 5, 4, 7, 8, 8, 9, **?**

NUS National University of Singapore | iSS INSTITUTE OF SYSTEMS SCIENCE

# Time step
## Make segments

- For time series problem, we don't feed in the entire series into the net, as this will not have the 'recurrent' learning involved

- Instead, we chop the series into segments, each segment with a fixed amount of time steps, called length

prumls/m3.2/v1.0

# makeSteps
The procedure

- Assume we want to have a length of 4, and a distance of 3

|     | segment 2 |     |     | segment 4 |     |
|-----|-----------|-----|-----|-----------|-----|

0,  1,  1,  3,  2,  2,  4,  5,  4,  7,  8,  8,  9

|   segment 1   |     |   segment 3   |     |

- The preprocessed input:

```
0,  1,  1,  3
3,  2,  2,  4
4,  5,  4,  7
7,  8,  8,  9
```

- Assume we want to have a length of 5, and a distance of 3

not sufficient length
to form segment

segment 2

0, 1, 1, 3, 2, 2, 4, 5, 4, 7, 8, 8, 9

segment 1

segment 3

- The preprocessed input:

0, 1, 1, 3, 2

3, 2, 2, 4, 5

4, 5, 4, 7, 8

# Time for exercise

- Write a function that can generate the desired preprocessed input from a 1D series/signals with the below signature

```
> def makeSteps(dat, length, dist):
```

```
0, 1, 1, 3, 2, 2, 4, 5, 4, 7, 8, 8, 9
```

- length 4, distance 3

```
0, 1, 1, 3
3, 2, 2, 4
4, 5, 4, 7
7, 8, 8, 9
```

- length 5, distance 3

```
0, 1, 1, 3, 2
3, 2, 2, 4, 5
4, 5, 4, 7, 8
```

NUS | iSS

# RNN
The working

$$0, \quad 1, \quad 1, \quad 3, \quad 2 \quad \longrightarrow \quad \mathbf{X}_1$$
$$3, \quad 2, \quad 2, \quad 4, \quad 5 \quad \longrightarrow \quad \mathbf{X}_2$$
$$4, \quad 5, \quad 4, \quad 7, \quad 8 \quad \longrightarrow \quad \mathbf{X}_3$$



$\mathbf{x}_1$    rnn    $\mathbf{x}_2$    $\mathbf{x}_3$

W    W    W

$\mathbf{s}_1 = \mathbf{W} \cdot \mathbf{x}_1$    $\mathbf{s}_2 = \mathbf{W} \cdot \mathbf{x}_2$    $\mathbf{s}_3 = \mathbf{W} \cdot \mathbf{x}_3$

$\mathbf{q}_2 = \mathbf{U} \cdot \mathbf{h}_1$    $\mathbf{q}_3 = \mathbf{U} \cdot \mathbf{h}_2$

activation    U    activation    U    activation

$\mathbf{h}_1 = \tanh\left(\mathbf{s}_1 + \mathbf{b}_0\right)$    $\mathbf{h}_2 = \tanh\left(\mathbf{s}_2 + \mathbf{q}_2 + \mathbf{b}_0\right)$    $\mathbf{h}_3 = \tanh\left(\mathbf{s}_3 + \mathbf{q}_3 + \mathbf{b}_0\right)$

NUS National University of Singapore    ISS INSTITUTE OF SYSTEMS SCIENCE

# Dot product
The working

• What is the output of

$$\mathbf{W} \cdot \mathbf{x}$$

• Assume the $\mathbf{W}$ is a $m$ x $n$ matrix, $\mathbf{x}$ is $n$ x 1 matrix, the matrix-vector dot product is

$$\mathbf{W} \cdot \mathbf{x} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \cdots & w_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} w_{11}x_1 + w_{12}x_2 + \cdots + w_{1n}x_n \\ w_{21}x_1 + w_{22}x_2 + \cdots + w_{2n}x_n \\ \vdots \\ w_{m1}x_1 + w_{m2}x_2 + \cdots + w_{mn}x_n \end{bmatrix}$$

• The $m$ determines the size of the product output, i.e. the output feature size

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

# Dot product

The working

- What is the output of

$$\mathbf{W} \cdot \mathbf{x}_1$$

- Assume

$$\mathbf{W} = \begin{bmatrix} 1 & 0 & 0 & 1 & -1 \\ 2 & 0 & 0 & -1 & 1 \end{bmatrix}$$

- and we know

$$\mathbf{x}_1 = \begin{bmatrix} 0 & 1 & 1 & 3 & 2 \end{bmatrix}$$

# Dot product
The working

- What is the output of

$$\mathbf{W} \cdot \mathbf{x}_1$$

- Assume

$$\mathbf{W} = \begin{bmatrix} 1 & 0 & 0 & 1 & -1 \\ 2 & 0 & 0 & -1 & 1 \end{bmatrix}$$

- and we know

$$\mathbf{x}_1 = \begin{bmatrix} 0 & 1 & 1 & 3 & 2 \end{bmatrix}$$

$$\mathbf{W} \cdot \mathbf{x}_1 = \begin{bmatrix} 1 & 0 & 0 & 1 & -1 \\ 2 & 0 & 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \\ 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \times 0 + 0 \times 1 + 0 \times 1 + 1 \times 3 + (-1) \times 2 \\ 2 \times 0 + 0 \times 1 + 0 \times 1 + (-1) \times 3 + 1 \times 2 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

# RNN
The working

- Note: The same U and W for $\mathbf{x}_1$, $\mathbf{x_2}$, $\mathbf{x}_3$

- U and W are changed **only** during updating phase of training

- U is a square matrix



$\mathbf{x}_1$

W

$\mathbf{s}_1 = \mathbf{W} \cdot \mathbf{x}_1$

$+$

activation

$\mathbf{h}_1 = \tanh\left(\mathbf{s}_1 + \mathbf{b}_0\right)$

rnn

$\mathbf{q}_2 = \mathbf{U} \cdot \mathbf{h}_1$

U

$\mathbf{x}_2$

W

$\mathbf{s}_2 = \mathbf{W} \cdot \mathbf{x}_2$

$+$

activation

$\mathbf{h}_2 = \tanh\left(\mathbf{s}_2 + \mathbf{q}_2 + \mathbf{b}_0\right)$

$\mathbf{q}_3 = \mathbf{U} \cdot \mathbf{h}_2$

U

$\mathbf{x}_3$

W

$\mathbf{s}_3 = \mathbf{W} \cdot \mathbf{x}_3$

$+$

activation

$\mathbf{h}_3 = \tanh\left(\mathbf{s}_3 + \mathbf{q}_3 + \mathbf{b}_0\right)$

NUS National University of Singapore | iSS INSTITUTE OF SYSTEMS SCIENCE

# RNN
The problem



$$\mathbf{x}_2$$

rnn

W

$$\mathbf{s}_2 = \mathbf{W} \cdot \mathbf{x}_2$$

$$\mathbf{q}_2 = \mathbf{U} \cdot \mathbf{h}_1 \longrightarrow \oplus$$

U          activation

$$\mathbf{h}_2 = \tanh\left(\mathbf{s}_2 + \mathbf{q}_2 + \mathbf{b}_0\right)$$

- Ideally, we want rnn to have long memories, so that it can connect relationships among data points far before and after a point of interest

- If this is possible, good for language understanding / translation, or how events in stock market correlate to each other

- But, in rnn, the more we perform recurrent operation, mathematically, that is equivalent to adding more layers to the net

- So vanishing gradient comes in ...

# RNN
How to build rnn in Keras

- Use `SimpleRNN` for RNN layers

```python
> from tensorflow.keras.layers import Input
> from tensorflow.keras.layers import SimpleRNN
> from tensorflow.keras.layers import Dense
> from tensorflow.keras.models import Model

> inputs        = Input(shape=(16,64))
> y             = SimpleRNN(32)(inputs)
> y             = Dense(1, activation='sigmoid')(y)

> model         = Model(inputs=inputs,outputs=y)

> model.summary()
```

- How many segments in a single sample? what is the length of each segment? What is the size of the output vector from the RNN layer?

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

# RNN
How to build rnn in Keras

```python
> from tensorflow.keras.layers import Input
> from tensorflow.keras.layers import SimpleRNN
> from tensorflow.keras.layers import Dense
> from tensorflow.keras.models import Model


> inputs        = Input(shape=(16,64))
> y             = SimpleRNN(32)(inputs)
> y             = Dense(1, activation='sigmoid')(y)


> model         = Model(inputs=inputs,outputs=y)


> model.summary()
```

```
_____
Layer (type)                Output Shape              Param #
=================================================================
input_1 (InputLayer)        (None, 16, 64)            0
_____
simple_rnn (SimpleRNN)      (None, 32)                3104        ⟶  ?
_____
dense (Dense)               (None, 1)                 33          ⟶  32 x 1 + 1 = 33
=================================================================
Total params: 3,137                                                        ↑
Trainable params: 3,137                                          the bias of the single
Non-trainable params: 0                                         output neuron
_____
```

# RNN
Parameter calculation

- Let $l_{in}$ denote the number of input feature (the length of each segment)

- Let $l_{out}$ denote the number of the output feature

- The number of parameters:

$$p = \left(l_{out} \times l_{in}\right) + \left(l_{out} \times l_{out}\right) + l_{out}$$

```
_____
Layer (type)              Output Shape              Param #
=============================================================
input_1 (InputLayer)      (None, 16, 64)            0
_____
simple_rnn (SimpleRNN)    (None, 32)                3104
_____
dense (Dense)             (None, 1)                 33
=============================================================
Total params: 3,137
Trainable params: 3,137
Non-trainable params: 0
_____
```

⟶ 32 x 64 + 32 x 32 + 32 = 2048 + 1024 + 32 = 3104

⟶ 32 x 1 + 1 = 33

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

• The working of LSTM in Keras
(according to Francois Chollet)



$$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3$$

$$W_0 \qquad W_0 \qquad W_0$$

$$\mathbf{s}_1 = \mathbf{W}_0 \cdot \mathbf{x}_1 \qquad \mathbf{s}_2 = \mathbf{W}_0 \cdot \mathbf{x}_2 \qquad \mathbf{s}_3 = \mathbf{W}_0 \cdot \mathbf{x}_3$$

$$\mathbf{z}_1 \to (+) \qquad V_0 \to \mathbf{z}_2 = \mathbf{V}_0 \cdot \mathbf{c}_2 \to (+) \qquad V_0 \to \mathbf{z}_3 = \mathbf{V}_0 \cdot \mathbf{c}_3 \to (+)$$

$$\mathbf{c}_1 \qquad \mathbf{c}_2 \quad \mathbf{q}_2 = \mathbf{U}_0 \cdot \mathbf{h}_1 \to (+) \qquad \mathbf{c}_3 \quad \mathbf{q}_3 = \mathbf{U}_0 \cdot \mathbf{h}_2 \to (+)$$

activation   carry   $U_0$   activation   carry   $U_0$   activation

$$\mathbf{h}_1 = \tanh\left(\mathbf{s}_1 + \mathbf{z}_1 + \mathbf{b}_0\right) \qquad \mathbf{h}_2 = \tanh\left(\mathbf{s}_2 + \mathbf{z}_2 + \mathbf{q}_2 + \mathbf{b}_0\right) \qquad \mathbf{h}_3 = \tanh\left(\mathbf{s}_3 + \mathbf{z}_3 + \mathbf{q}_3 + \mathbf{b}_0\right)$$

NUS National University of Singapore   ISS INSTITUTE OF SYSTEMS SCIENCE

# Multiplication
Element-wise

- The output of element-wise multiplication is

$$\mathbf{v} \odot \mathbf{x} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \odot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} v_1 x_1 \\ v_2 x_2 \\ \vdots \\ v_n x_n \end{bmatrix}$$

- The length of $\mathbf{v}$ and $\mathbf{x}$ must be equal

# LSTM

The internal working of
'carry'

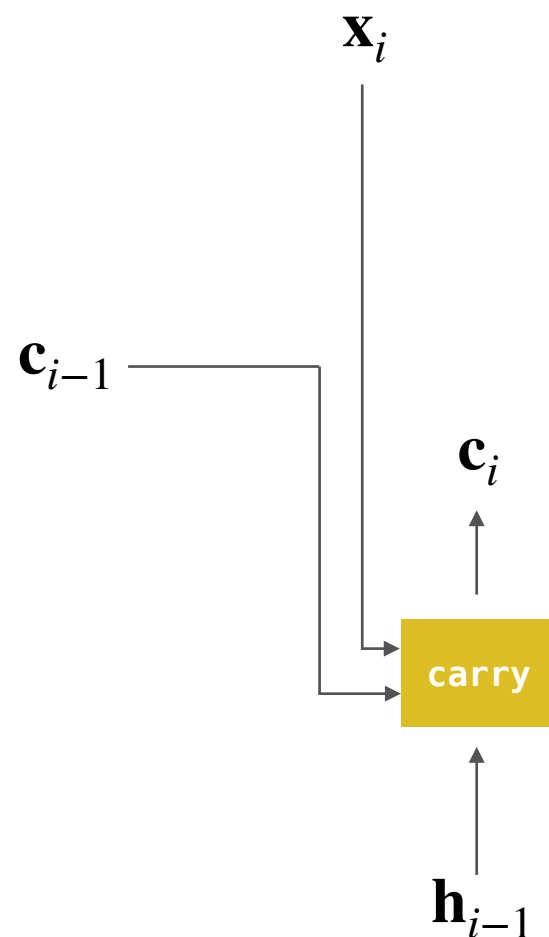$$\mathfrak{w}_i = \sigma \left( \mathbf{U}_\mathfrak{w} \cdot \mathbf{h}_{i-1} + \mathbf{W}_\mathfrak{w} \cdot \mathbf{x}_i + \mathbf{b}_\mathfrak{w} \right)$$

$$\mathfrak{f}_i = \sigma \left( \mathbf{U}_\mathfrak{f} \cdot \mathbf{h}_{i-1} + \mathbf{W}_\mathfrak{f} \cdot \mathbf{x}_i + \mathbf{b}_\mathfrak{f} \right)$$

$$\mathfrak{h}_i = \sigma \left( \mathbf{U}_\mathfrak{h} \cdot \mathbf{h}_{i-1} + \mathbf{W}_\mathfrak{h} \cdot \mathbf{x}_i + \mathbf{b}_\mathfrak{h} \right)$$

- $\sigma$ is the sigmoid function, and let's denote $\odot$ as element-wise multiplication, and we have

$$\mathbf{c}_i = \mathfrak{w}_i \odot \mathfrak{h}_{i-1} + \mathbf{c}_{i-1} \odot \mathfrak{f}_i$$

$\mathbf{x}_i$

$\mathbf{c}_{i-1}$

$\mathbf{c}_i$

carry

$\mathbf{h}_{i-1}$

NUS National University of Singapore | iSS INSTITUTE OF SYSTEMS SCIENCE

# LSTM
How to build lstm in Keras

```
> from tensorflow.keras.layers import Input
> from tensorflow.keras.layers import LSTM
> from tensorflow.keras.layers import Dense
> from tensorflow.keras.models import Model


> inputs      = Input(shape=(16,64))
> y           = LSTM(32)(inputs)
> y           = Dense(1, activation='sigmoid')(y)


> model       = Model(inputs=inputs,outputs=y)


> model.summary()
```

```
_____
Layer (type)              Output Shape          Param #
=======================================================
input_2 (InputLayer)      (None, 16, 64)        0
_____
lstm (LSTM)               (None, 32)            12416      ⟶  ?
_____
dense_1 (Dense)           (None, 1)             33         ⟶  32 x 1 + 1 = 33
=======================================================
Total params: 12,449
Trainable params: 12,449
Non-trainable params: 0
_____
```

the bias of the single
output neuron

# LSTM
Parameters calculation

• Let $l_{\text{in}}$ denote the number of input feature (the length of each segment), $l_{\text{out}}$ denote the number of the output feature

• The number of parameters:

$$p = \left(l_{\text{out}} \times l_{\text{out}} + l_{\text{out}} \times l_{\text{in}} + l_{\text{out}}\right) \times 4$$

```
_____
Layer (type)            Output Shape            Param #
===============================================================
input_2 (InputLayer)    (None, 16, 64)          0
_____
lstm (LSTM)             (None, 32)              12416
_____
dense_1 (Dense)         (None, 1)               33
===============================================================
Total params: 12,449
Trainable params: 12,449
Non-trainable params: 0
_____
```

```
 (32x32 + 32x64 + 32) x 4
= 3104 x 4
= 12416
```

# Stacking
Multiple recurrent layers

- It is possible to stack recurrent layers and make the net deeper

- Note: for LSTM 2 to Dense, only the final output sequence fed into Dense layer

NUS National University of Singapore | iSS INSTITUTE OF SYSTEMS SCIENCE

# LSTM
How to build stacked lstm in Keras

```
> inputs  = Input(shape=(3,5))
> y       = LSTM(7, return_sequences=True)(inputs)
> y       = LSTM(9)(y)
> y       = Dense(1, activation='sigmoid')(y)

> model   = Model(inputs=inputs,outputs=y)
> model.summary()
```

To stack lstm, return_sequences must be set to True

```
_____
Layer (type)                 Output Shape              Param #
=================================================================
input_3 (InputLayer)         (None, 3, 5)              0
_____
lstm_1 (LSTM)                (None, 3, 7)              364
_____
lstm_2 (LSTM)                (None, 9)                 612
_____
dense_2 (Dense)              (None, 1)                 10
=================================================================
Total params: 986
Trainable params: 986
Non-trainable params: 0
_____
```

The output shape is (None, 3, 7) because of returned sequences, else it will simply be (None, 7)

**The number of rows** of the output of a lstm that returns sequences, is always equal to **the number of rows** of the input to the unit

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE