



# Deep Learning

## Lecture 8: Deep Learning Architectures (Part I)

University of Agder,  
Kristiansand, Norway

Prepared by: Hadi Ghauch<sup>\*</sup>, Hossein S. Ghadikolaei<sup>†</sup>

<sup>\*</sup> Telecom ParisTech

<sup>†</sup> Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>

April 2019

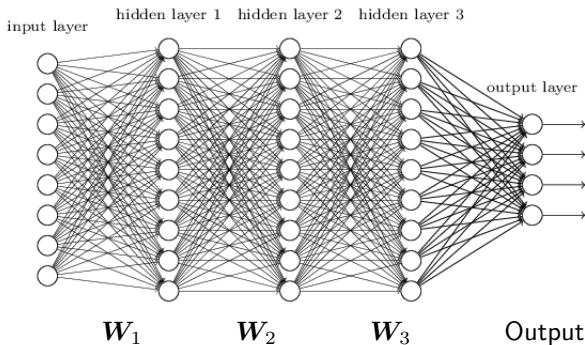
# Outline

1. Recurrent Neural Networks
2. Long Short Term Memory Networks

# Recap: Deep Neural Networks (DNNs)

## Multi-layer Perceptron:

- Output from some perceptrons are used in the inputs to other perceptrons

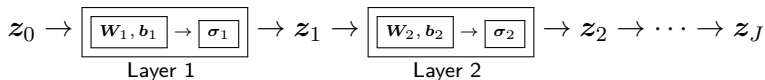


## Feedforward Deep Neural Network (DNN):

- Multi-layer Perceptron with 'many' hidden layers
- Only forward connections: from input to hidden layer or output

# Recap: Deep Neural Networks (DNNs)

DNN with  $J$ -layers: cascade of  $J$  filters + activation



In compact form:  $z_j = \sigma_j(W_j z_{j-1} + b_j)$

$z_j \in \mathbb{R}^{d_j}$ : output to layer  $j$

$W_j \in \mathbb{R}^{d_j \times d_{j-1}}$ : **Weights** for layer  $j$  (drop bias for simplicity)

$\sigma_j : \mathbb{R}^{d_j} \mapsto \mathbb{R}^{d_j}$ : **Activation func** for layer  $j$

# Outline

1. Recurrent Neural Networks

2. Long Short Term Memory Networks

# Motivation and Setting

Feedforward NNs assume input/output have **fixed dimension**:

e.g. image classification, regression

assume **samples of training set are i.i.d.**

these assumptions restrictive when dealing with sequential data

**Feedforward NN ill-equipped** when **data is sequential**: temporal/causal dependence in data: e.g., classifying frames of a movie, or words in sentence

knowing previous samples critical to predict next sample

**Recurrent NNs**: family of NNs designed to learn sequential data

training set **time-dependent samples** (not i.i.d.):  $\{ (x_t, y_t) \}_{t=1}^T$

sample  $(x_t, y_t)$  depends on previous ones

**classifying images/frames of a movie**: time index  $t$  is frame number

consecutive frames correlated (i.i.d. assumption does not hold)

learning such relation critical to classify other frames

**Intuition**: Share parameters from different parts of model

enables learning sequences not seen in training and improve generalization

Recurrent NNs implement parameter sharing by design:

each output depends on previous outputs

# Recurrent NN

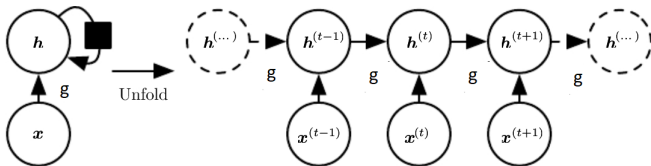
Motivation: **dynamical system**,  $s_t = g(s_{t-1}, x_t)$ ,  $x_t$  external input  
states are **recursively generated by**  $g()$

**Goal:** Use recurrent NN learn states of dynamical system  
not suitable for feedforward NNs

**Recurrent NNs:** designed specifically to learn states in a recursive setting  
similar to feed-forward NNs used to learn non-recursive func

several designs used. start with **plain recurrent NNs**:  $h_t = g(h_{t-1}, x_t)$   
set state the dynamical system equal to hidden unit,  $h_t$

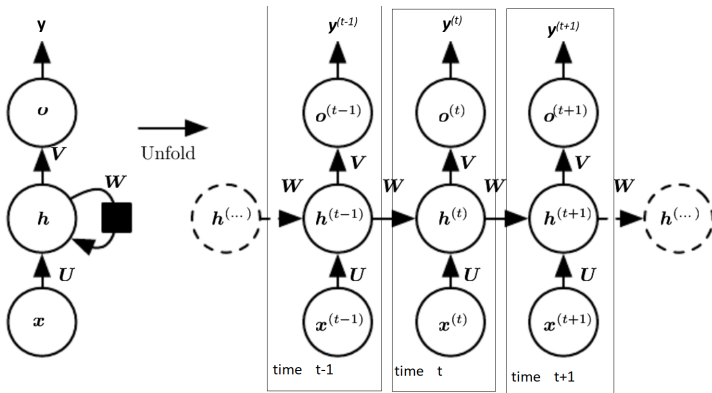
Model system by unfolding the states: **start with unfolded model**



Simple **toy architecture**. Next, move more realistic architecture.

# Prevalent design Recurrent NNs

So far no learning: now lets add some **parameters to learn**



Single input single output at each time: at time  $t$ , one input-output pair  $(x_t, y_t)$ .

Goal: Learn a **same set parameters/weights**,  $V, W, U$ , that approximates **all samples**,  $\{x_t, y_t\}_{t=1}^T$

parameters/weights shared across all samples: notion of **parameter sharing**



# Mathematical Model of Recurrent NNs

## Model for sample $t$ :

$\mathbf{x}_t \in \mathbb{R}^n$ : Input vector at time  $t$

$\mathbf{h}_t \in \mathbb{R}^n$ : Hidden layer/unit at time  $t$

## Inputs for $\mathbf{h}_t$ :

linear combination of prev hidden layer:  $\mathbf{W}\mathbf{h}_{t-1}$ , where  $\mathbf{W} \in \mathbb{R}^{n \times n}$

linear combination of input:  $\mathbf{U}\mathbf{x}_t$ , where  $\mathbf{U} \in \mathbb{R}^{n \times n}$

**what happens in  $\mathbf{h}_t$ ?** sum them then apply non-linear activation:  $\psi(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t)$

$\psi: \mathbb{R}^n \mapsto \mathbb{R}^n$ : vector-valued vector func (applied element-wise)

## Output for $\mathbf{h}_t$ :

linearly combine output of  $\mathbf{h}_t$ :  $\mathbf{V}\mathbf{h}_t$ , where  $\mathbf{V} \in \mathbb{R}^{n \times n}$

then apply non-linear activation:  $\phi(\mathbf{V}\mathbf{h}_t)$

$\phi: \mathbb{R}^n \mapsto \mathbb{R}^n$ : vector-valued vector func (applied element-wise)

$\hat{\mathbf{y}}_t \in \mathbb{R}^n$ : **predicted output vector** (by RNN) at time  $t$

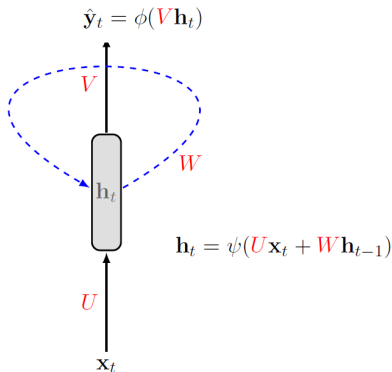
$\mathbf{y}_t \in \mathbb{R}^n$ : **true output vector** at time  $t$

$\{\mathbf{x}_t, \mathbf{y}_t\}_{t=1}^T$ : training set of  $T$  samples (recall the time index  $t$ )

assume all dimensions are same ( $= n$ ) for simplicity

**Goal:** Learn RNN parameters,  $\mathbf{W}, \mathbf{V}, \mathbf{U}$  to minimize the loss b/w  
predicted output by RNN,  $\hat{\mathbf{y}}_t$ , and true output,  $\mathbf{y}_t$

# Plain Recurrent NNs (folded form)



snapshot of input-output at time  $t$

Source: Trivedi, Kodor, 2017

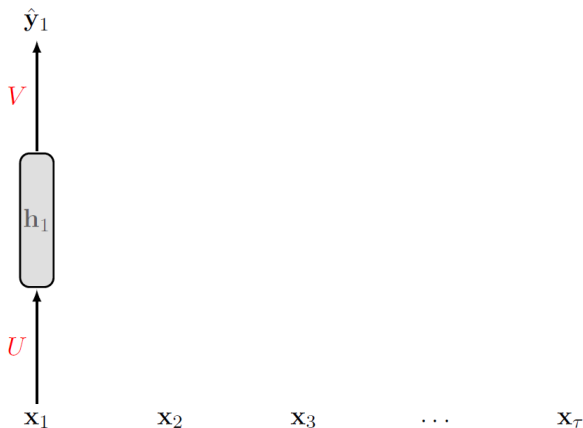
## Plain Recurrent NNs (folded form):

$\phi()$  ,  $\psi()$ : vector-valued vector function (non-linear activation)

$U$ ,  $W$ ,  $V$ : weight matrices of appropriate size

# Plain Recurrent NNs (unfolded form)

Previous Recurrent NN in unfolded form

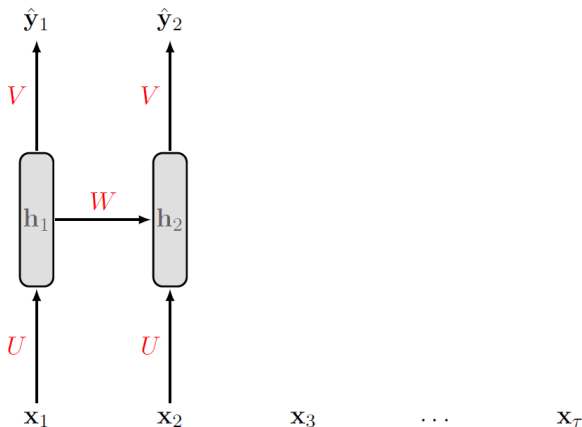


Source: Trivedi, Kodor, 2017

snapshot of input-output relation at time  $t = 1$

# Plain Recurrent NNs (unfolded form)

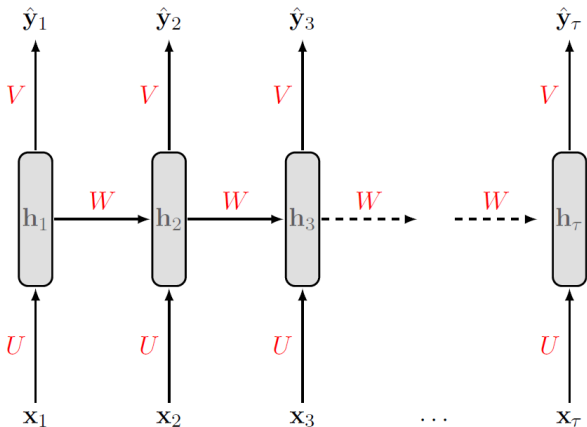
Previous Recurrent NN in unfolded form



Source: Trivedi, Kodor, 2017

snapshot of input-output relation at time  $t = 2$

# Plain Recurrent NNs (unfolded form)



Parameters (weights) the same or all time samples (weights indep of  $t$ ):

learn set of **parameters (weights) shared for all time samples**

By learning  $W, U, V$  uncover **time-dependence/correlation** among **time-dependent samples** of training set

# Training Recurrent NN

**Express mathematically input-output relation** at time  $t$ :

$$\mathbf{a}_t = \mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + b:$$

linearly combine  $\mathbf{x}_t$  and  $\mathbf{h}_{t-1}$

$$\mathbf{h}_t = \psi(\mathbf{a}_t):$$

apply non-linearity at input of hidden layer

$$\mathbf{o}_t = \mathbf{V}\mathbf{h}_t + c:$$

linearly combine output of hidden layer

$$\hat{\mathbf{y}}_t = \phi(\mathbf{o}_t):$$

apply non-linearity just before output

$\hat{\mathbf{y}}_t$  is prediction that RNN outputs for sample  $\mathbf{x}_t$  (drop  $b, c$  wlog):

$$\hat{\mathbf{y}}_t := \phi[ \mathbf{V}\psi(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t) ]$$

# Training Recurrent NN

## Training Recurrent NNs:

Loss for sample  $t$ :  $f_t(\mathbf{y}_t, \hat{\mathbf{y}}_t)$  (diff b/w true and predicted value)

**negative log-likelihood** loss is a common choice

total loss function,  $f$  is sum of individual losses:  $f := \sum_t f_t$

## Empirical risk minimization:

$$\arg \min_{\mathbf{U}, \mathbf{W}, \mathbf{V}} \sum_{t=1}^T f_t \left( \underbrace{\mathbf{y}_t}_{\text{true}}, \underbrace{\phi[ \mathbf{V} \psi(\mathbf{W} \mathbf{h}_{t-1} + \mathbf{U} \mathbf{x}_t) ]}_{\text{RNN prediction}} \right) := f \quad (1)$$

How to optimize  $\mathbf{U}, \mathbf{W}, \mathbf{V}$  ?

Compute  $\nabla_{\mathbf{W}} f, \nabla_{\mathbf{U}} f, \nabla_{\mathbf{V}} f$  using same 'tricks' as backpropagation called **backpropagation trough time (BPTT)** for RNN training

# BPTT

1. gradient of  $V$ :

$$\nabla_V f = \sum_t (\nabla_{o_t} f)(\mathbf{h}_t)^T$$

2. gradient of  $W$ :

$$\nabla_W f = \sum_t [ \text{diag}(\mathbf{1} - \mathbf{h}_t \circ \mathbf{h}_t) (\nabla_{\mathbf{h}_t} f) (\mathbf{h}_{t-1})^T ]$$

$$(\nabla_{\mathbf{h}_t} f) = \begin{cases} (\mathbf{V})^T (\nabla_{o_t} f), & \text{for } t = T \\ (\mathbf{W})^T \text{diag}(\mathbf{1} - \mathbf{h}_{t+1} \circ \mathbf{h}_{t+1}) (\nabla_{\mathbf{h}_{t+1}} f) + (\mathbf{V})^T (\nabla_{o_t} f), & \text{for } t < T \end{cases}$$

$(\nabla_{\mathbf{h}_t} f)$  depends on  $(\nabla_{\mathbf{h}_{t+1}} f)$ :

start hidden layers with largest index,  $\mathbf{h}_T$ , and work backwards (like BackProp)

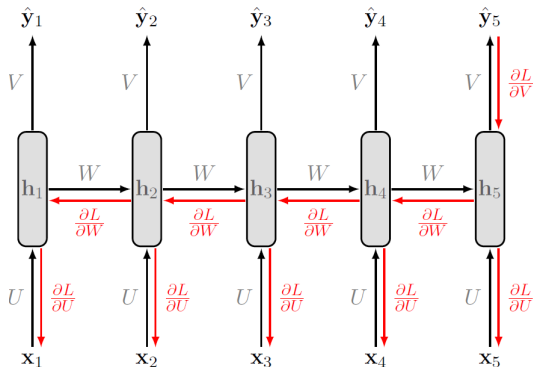
3. gradient of  $U$ :

$$\nabla_U f = \sum_t [ \text{diag}(\mathbf{1} - \mathbf{h}_t \circ \mathbf{h}_t) (\nabla_{\mathbf{h}_t} f)(\mathbf{x}_t)^T ]$$

see [Goodfellow, 2016], Chap 10, for detailed derivations of BPTT.



# BPTT



Source: Trivedi, Kodor, 2017

**BPTT** = BP in the unrolled RNN

sample  $(x_5, \hat{y}_5)$ : compute  $\nabla_V f$ , then  $\nabla_W f$ , then  $\nabla_U f$

sample  $(x_4, \hat{y}_4)$ : compute  $\nabla_V f$ , then  $\nabla_W f$ , then  $\nabla_U f$

sample  $(x_3, \hat{y}_3)$ : compute  $\nabla_V f$ , then  $\nabla_W f$ , then  $\nabla_U f$

## Some variants on Plain Recurrent NN:

- many inputs ( $x_t$ ), single output ( $y_t$ ). Used in sentiment analysis
- **bidirectional Recurrent NNs**: plain RNNs capture time-dependence in chronological order:  $h_1, \dots, h_t$ . Bidirectional RNNs also capable of capturing reverse chronological order:  $g_t, \dots, g_1$ .  
widely applied in handwriting and speech recognition, bioinformatics
- **Encoder-decoder Recurrent NN**: designed to handle input and output sequences that have different lengths.  
widely used in language translation (e.g., one sentence may have different number of words in different languages)
- **Deep Recurrent NN**: more hidden layers between input and output  
depth increase approximation capacity (but make optimization harder)
- **Deep Recursive NN**: computational graph is a tree  
used for Natural Language processing and computer vision

# Outline

1. Recurrent Neural Networks

2. Long Short Term Memory Networks

# Practical Problems for Deep RNNs:

## Challenges for learning long term dependencies:

many layers (matrix multiplications) needed to capture long-term dependence  
i.e., **deep RNNs** needed

**challenge for deep RNNs:** output becomes too small after multiplication with many matrices

illustrate with a simple RNN (w/out nonlinearity),  $\mathbf{W} \in \mathbb{R}^{n \times n}$

$$\mathbf{h}_t = \mathbf{W}\mathbf{h}_{t-1} \Rightarrow \mathbf{h}_T = \underbrace{\mathbf{W}\mathbf{W} \dots \mathbf{W}}_{T \text{ times}} \mathbf{h}_1 = \mathbf{W}^T \mathbf{h}_1 \Rightarrow \|\mathbf{h}_T\|_2 \leq (\lambda_{\max}[\mathbf{W}])^T \|\mathbf{h}_1\|_2$$

if  $\lambda_{\max}[\mathbf{W}] < 1 \Rightarrow \|\mathbf{h}_T\|_2 \rightarrow 0, T \rightarrow \infty$

output,  $\mathbf{h}_T$ , becomes too small, so does its gradient **vanishing gradient**

depends on spectral radius of  $\mathbf{W}$

problem peculiar to RNNs (due to **loops** b/w input and output)

**some fixes to vanishing gradient:** skip connections, leaky hidden units, echo state networks; see [Goodfellow, 2016], Chap 10

**Exploding gradient:**  $\lambda_{\min}[\mathbf{W}] > 1 \Rightarrow \|\mathbf{h}_T\|_2 \rightarrow \infty, T \rightarrow \infty$

output,  $\mathbf{h}_T$ , blows up: so does the gradient of  $\mathbf{h}_T$  (rarely occurs)

# Long Short Term Memory (LSTM):

## Long Short Term Memory (LSTM) Networks:

among most efficient implementations of RNNs

alleviate vanishing grad: add more paths in NN where gradient well defined

**Idea:** additional loops in NN where gradient does not vanish

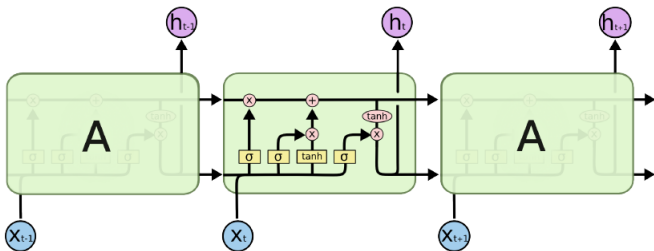
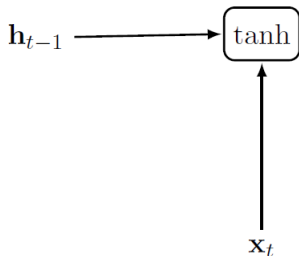


Figure: Chris Olah

basic unit of LSTM is a memory cell (shaded green). Let us zoom-in to understand it

# LSTM

Build LSTM gradually step by step.



$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t)$$

$\tanh()$ : non-linear activation  
(element-by-element func)

$\tilde{\mathbf{c}}_t$ : output of non-linear activation  
(tanh)

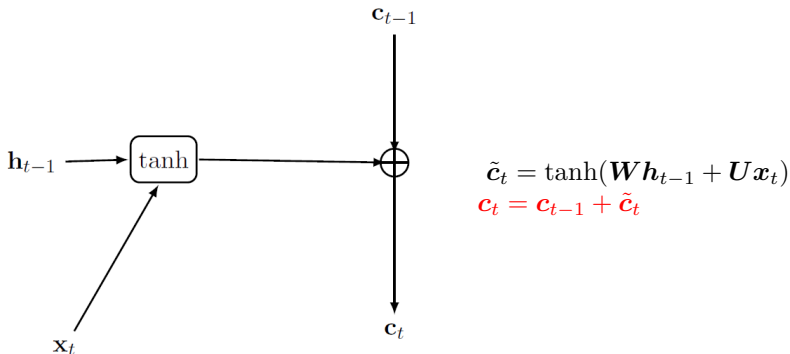
$\mathbf{U}, \mathbf{V}$ : matrices of appropriate size

snapshot of input-output at time  $t$

source: Trivedi, Kodori, 2017

so far it looks like first layer of a plain RNN

# LSTM



snapshot of input-output at time  $t$

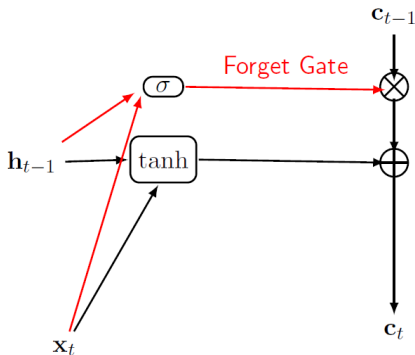
source: Trivedi, Kodori, 2017

$\tilde{c}_t$ : output of non-linear activation ( $\tanh$ )

$c_t$ : running sum of output of non-linear activation, over  $t$   
sum increases with  $t$  and may blow-up. how to prevent it?

scale down previous output  $c_{t-1}$ : **forget gate**

# LSTM



$$\begin{aligned}\tilde{c}_t &= \tanh(W h_{t-1} + U x_t) \\ c_t &= f_t \odot c_{t-1} + \tilde{c}_t \\ f_t &= \sigma_f(W_f h_{t-1} + U_f x_t)\end{aligned}$$

$W_f, U_f$  : weight matrices

snapshot of input-output at time  $t$   
source: Trivedi, Kodor, 2017

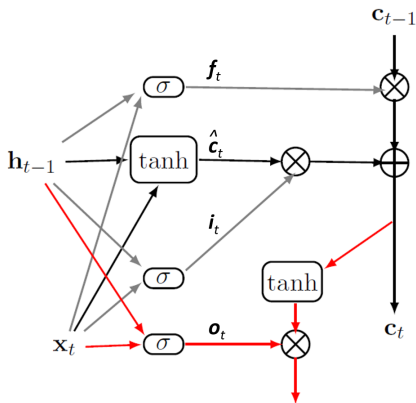
## LSTM with forget gate:

multiply (element by element) previous output  $c_{t-1}$ , by forget factor  $f_t$ :  
less importance on previous samples

$W_f, U_f$ : two additional weights to be learned



# LSTM



$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t)$$

$$c_t = f_t \circ c_{t-1} + \tilde{c}_t \circ i_t$$

$$f_t = \sigma_f(W_f h_{t-1} + U_f x_t)$$

Define input/output gates:

$$i_t = \sigma_i(W_i h_{t-1} + U_i x_t)$$

$$o_t = \sigma_o(W_o h_{t-1} + U_o x_t)$$

$$h_t = o_t \circ \tanh(c_t)$$

$W_i, U_i, W_o, U_o$ : weight matrices

snapshot of input-output at time  $t$   
source: Trivedi, Kodor, 2017

fully-fledged LSTM

# Unpacking the fully fledged LSTM:

Information added/deleted from cell via input, forget or output gates  
each gate modeled by two matrices  $W_x, U_x$  ( $x = \{o, i, f\}$ ), that are learned from training

**Cell state:**  $c_t = f_t \odot c_{t-1} + \tilde{c}_t \odot i_t$

think of  $c_t$  as state of the cell

updated as func of info to be deleted ( $f_t$ ), and to be added ( $i_t$ )

**Forget gate:**  $f_t = \sigma_f(W_f h_{t-1} + U_f x_t)$

$f_t$  multiplies previous output,  $c_{t-1}$ , to reduce importance of previous samples by learning  $W_f, U_f$  LSTM throw away some unnecessary information

**Input gate:**  $i_t = \sigma_i(W_i h_{t-1} + U_i x_t)$

what information from input to store in cell ?

by multiplying output of non-linearity,  $\tilde{c}_t$

intuition: learning  $W_i, U_i$  LSTM decides what information to add to cell (store)

**Output gate:**  $o_t = \sigma_o(W_o h_{t-1} + U_o x_t)$

pass a "filtered version" of output (same intuition as input)

intuition: learning  $W_o, U_o$  LSTM decides what information to pass to output

**Update cell state:**  $c_t = f_t \odot c_{t-1} + \tilde{c}_t \odot i_t$

update state taking into account forget gate  $f_t$ , and input  $i_t$

## Related issues LSTM:

**Training LSTM:** similar to Recurrent NNs

in addition to learning  $\mathbf{W}, \mathbf{U}$  for Recurrent NNs

for LSTM the weight matrices,  $\mathbf{W}_i, \mathbf{U}_i, \mathbf{W}_o, \mathbf{U}_o, \mathbf{W}_f, \mathbf{U}_f$ , also learned.

BPTT may still be used (derivations more complex)

**Gated Recurrent Unit:** is a related implementation of LSTM

see [Goodfellow, 2016], Chapter 10

# Conclusions

Several DL architectures

- Recurrent NNs: motivated by learning sequential data  
math model of plain recurrent NN, training with BPTT, some variants  
challenge of long term dependencies
- LSTM: practical solution to bypass challenge of long term dependencies  
solution to vanishing gradient  
discussed math model for LSTM, role of gates: input, output, forget

## Some references

J. Wu, "Introduction to convolutional neural networks ", May 1, 2017, available on arxiv

A. Goodfellow, Y. Bendgio, A Courville, "Deep Learning", MIT Press, Chap. 8, 9

Y. LeCun, L Bottou, G. Orr, K-B Muller, "Efficient Backprop", Neural Networks: Tricks of the trade

S. Trivedi, R. Kondor, CMSC 35246 Deep Learning, Spring 2017

C. M. Bishop, "Pattern Recognition and Machine Learning", Springer-Verlag, 2006

S. Haykin, "Neural Networks and Learning Machines", 3rd, 2009, Chap 4

Check Tutorial by Ruslan Salakhutdinov:  
<http://www.cs.cmu.edu/~rsalakhu/>



# Deep Learning

## Lecture 8: Deep Learning Architectures (Part I)

University of Agder,  
Kristiansand, Norway

Prepared by: Hadi Ghauch<sup>\*</sup>, Hossein S. Ghadikolaei<sup>†</sup>

<sup>\*</sup> Telecom ParisTech

<sup>†</sup> Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>

April 2019