

Long short-term memory (LSTM) networks for the analysis of time series

(CIV6540 - Probabilistic Machine Learning for Civil Engineers)

Professor: James-A. Goulet

Département des génies civil, géologique et des mines
Polytechnique Montréal



Chapter 10 – Goodfellow et al. (2016)
Deep Learning, MIT Press















Olah (2015)
Understanding LSTM Networks, (blog)





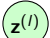

Module #7c Outline

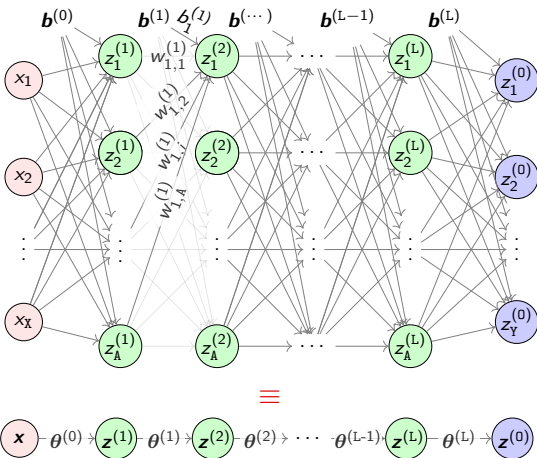
Introduction
Feedforward NN
(review)
Recurrent NN
LSTM
Example

Topics organization

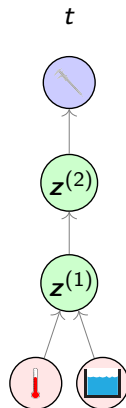
<i>Background</i>	1	Revision probability & linear algebra	
	2	Probability distributions	
<i>Machine Learning Basics</i>	0	Introduction	
	3	Bayesian Estimation	$p(A B) = \frac{p(B A)p(A)}{p(B)}$
	4	MCMC sampling & Newton	 
<i>Supervised learning</i>	5	Regression	
	6	Classification	 
	7	LSTM networks for time series	
<i>Unsupervised learning</i>	7	State-space model for time-series	
<i>Decision Making & RL</i>	8	Decision Theory	
	9	AI & Sequential decision problems	

Feedforward NN Representation

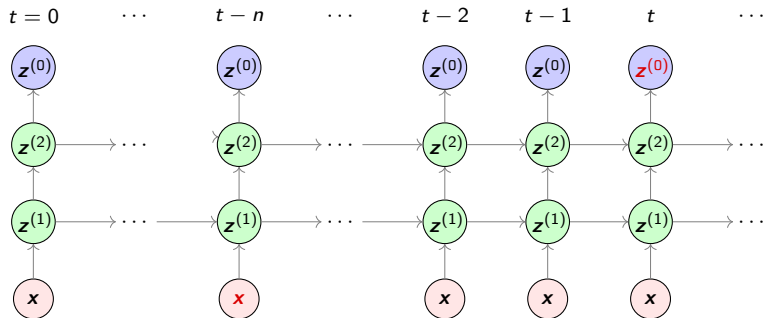
-  \mathbf{x} : Input layer
-  $\mathbf{a}^{(l)}$: Activation units
-  $\mathbf{z}^{(l)}$: Hidden units
-  $\mathbf{z}^{(0)}$: Output layer
- θ : parameters



Example - Modelling the displacement of a dam

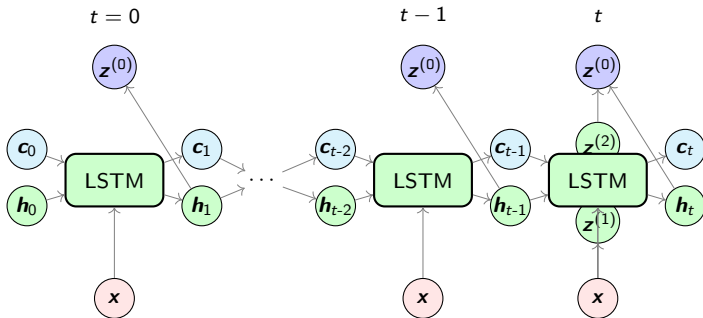


Recurrent neural networks



Vanishing gradient for long-term dependencies: $\frac{\partial z_t^{(0)}}{\partial x_{t-n}} \rightarrow 0$

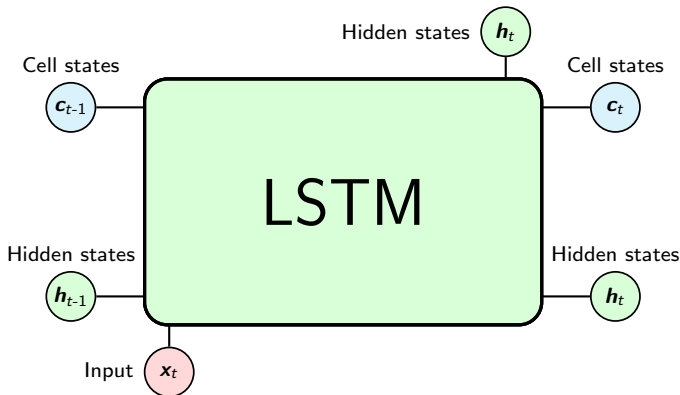
LSTM network



The key behind the LSTM is the cell state (memory)

$$c_t = [c_1, c_2, \dots, c_C]_t^T \in \mathbb{R}^C$$

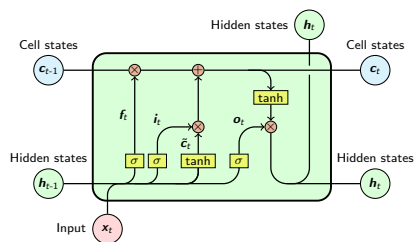
The LSTM cell



LSTM gates: **forget** irrelevant past information from c_{t-1} ,
write new relevant information in c_t ,
 and **rely on a subset** of c_t to make predictions h_t

[Adapted from J. Leon, Beerware]

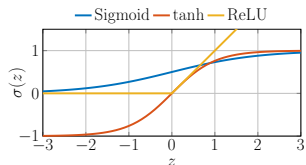
The forget gate



$$\mathbf{c}_{t-1} = [c_1, c_2, \dots, c_C]_{t-1}^T \in \mathbb{R}^C$$

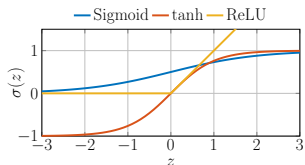
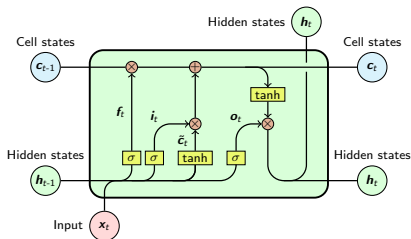
$$\mathbf{f}_t = [f_1, f_2, \dots, f_C]_t^T \in (0, 1)^C$$

$$\mathbf{f}_t \odot \mathbf{c}_{t-1} = [c_1, 0, \dots, c_C]_t^T$$



$$\mathbf{f}_t = \sigma(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_f)$$

The input and candidate gates



$$\tilde{\mathbf{c}}_t = [\tilde{c}_1, \tilde{c}_2, \dots, \tilde{c}_C]_t^T \in (-1, 1)^C$$

$$\mathbf{i}_t = [i_1, i_2, \dots, i_C]_t^T \in (0, 1)^C$$

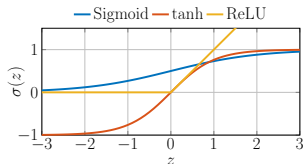
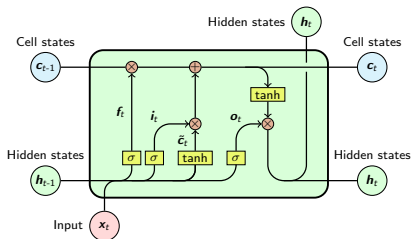
$$\mathbf{i}_t \odot \tilde{\mathbf{c}}_t = [\tilde{c}_1, 0, \dots, \tilde{c}_C]_t^T$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_i \cdot [\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_i)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_c \cdot [\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_c)$$

$$\mathbf{c}_t = \underbrace{\mathbf{f}_t \odot \mathbf{c}_{t-1}}_{\text{forget}} + \underbrace{\mathbf{i}_t \odot \tilde{\mathbf{c}}_t}_{\text{write}}$$

The output gate



$$\mathbf{c}_t = [c_1, c_2, \dots, c_C]_t^T \in \mathbb{R}^C$$

$$\mathbf{o}_t = [o_1, o_2, \dots, o_C]_t^T \in (0, 1)^C$$

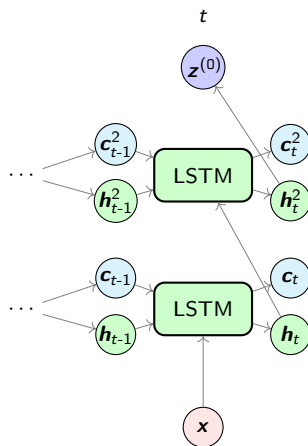
$$\mathbf{o}_t \odot \mathbf{c}_t = [c_1, 0, \dots, c_C]_t^T$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o \cdot [\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_o)$$

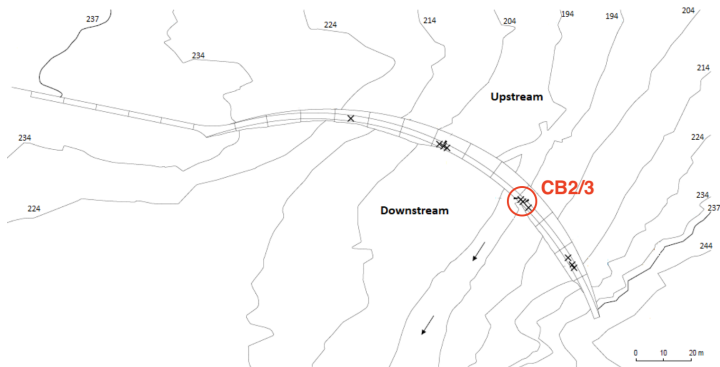
$$\mathbf{h}_t = \underbrace{\mathbf{o}_t \odot \tanh(\mathbf{c}_t)}$$

Attention on a subset of the cell state

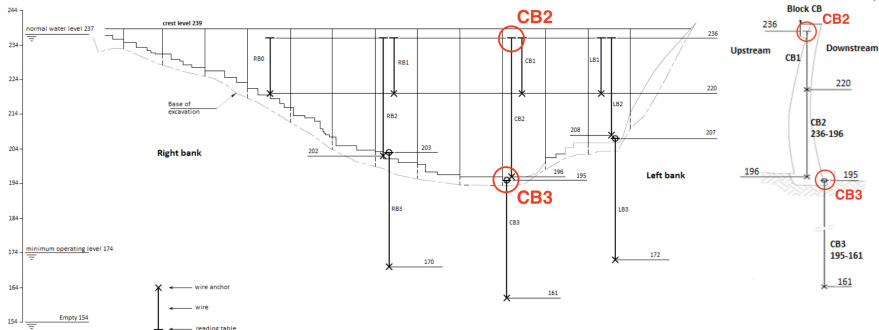
Multi-layers LSTM network



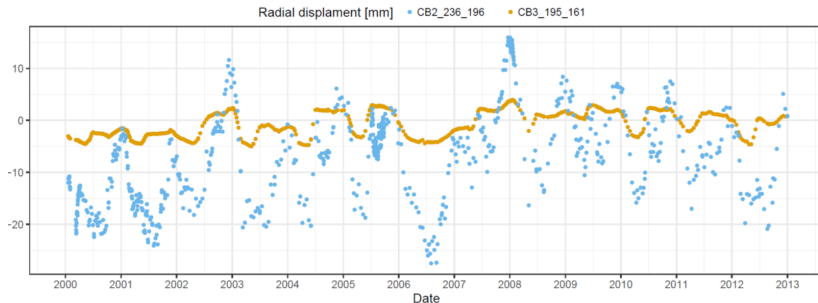
Dam structural health monitoring



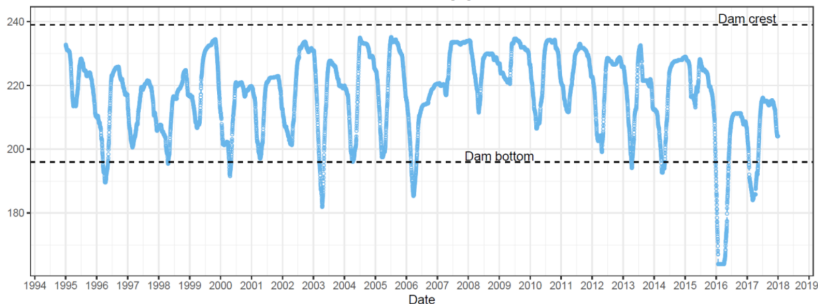
Dam structural health monitoring



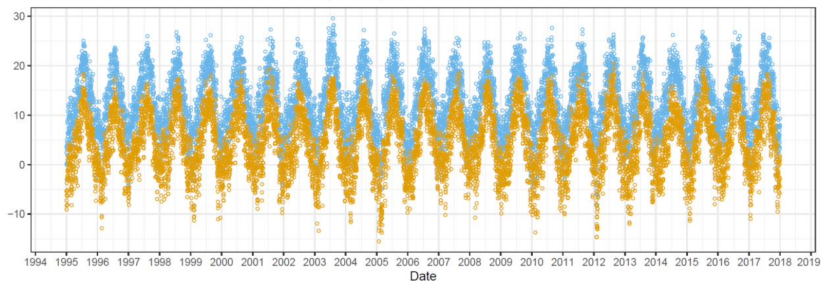
Data – Pendulum (~ 1 datapoint /1.5 week) ⓘ



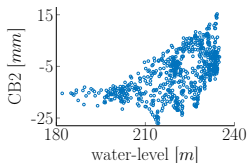
Data – Reservoir level



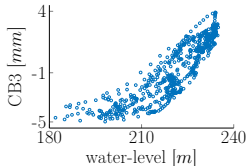
Data – Temperature 🌡️



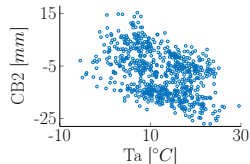
Relationships between datasets



Water level v.s.
CB2

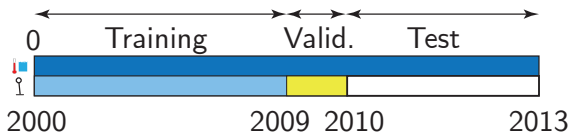


Water level v.s.
CB3

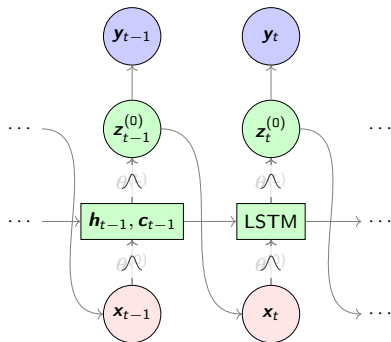


Temp. v.s.
CB2 & CB3

Compare the predictive capacity of LSTM and BDLM



LSTM architecture with teacher forcing



We use **TAGI** instead of backprop to train the model...

Two sources of uncertainty:

- Parameters (θ)

$$\theta = \mathcal{N}(\theta; \mu_\theta, \Sigma_\theta)$$

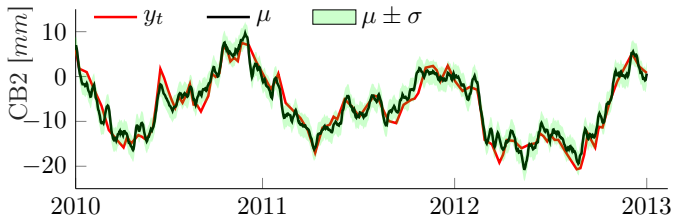
- Prediction error (w)

$$y = z^{(0)} + w$$

Comparison of LSTM & BDLM – CB2

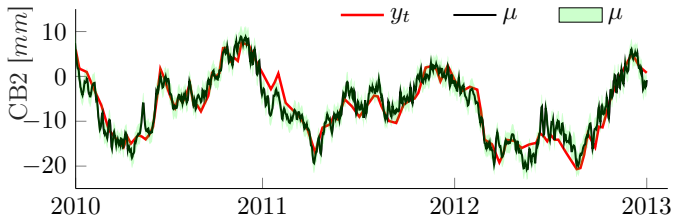
BDLM

($LL = -63.9$)



LSTM

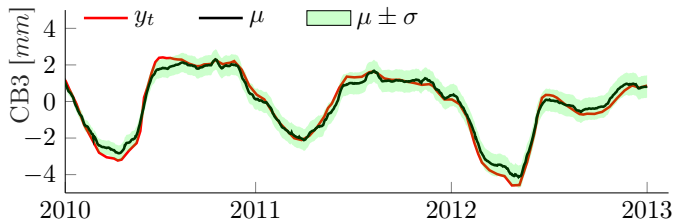
($LL = -65.9$)



Comparison of LSTM & BDLM – CB3

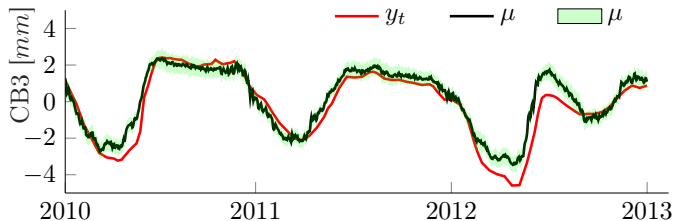
BDLM

($LL = -12.1$)

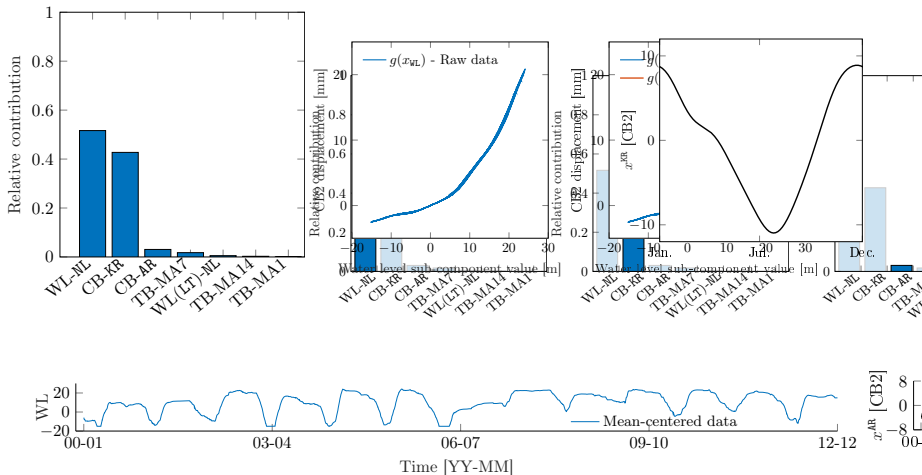


LSTM

($LL = -30.2$)








CB2 – BDLM's model structure & interpretation








LSTM v.s. BDLM...

BDLM

- ▶ Good accuracy 
- ▶ Uncertainties 
- ▶ Non-stationary data 
- ▶ Interpretable 
- ▶ Manual setup 

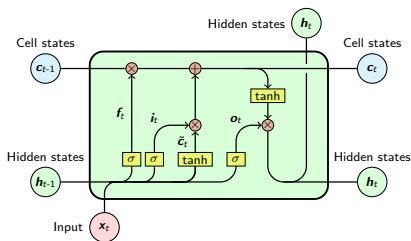
LSTM

- ▶ Good accuracy 
- ▶ Uncertainties *
- ▶ Stationary data 
- ▶ Poor interpretability 
- ▶ Black-box setup 



Deka, Vuong, Goulet, Côté and Miquel, (2022). *Dam Behaviour Prediction Using an Ensemble of Bayesian Dynamic Linear Model and Bayesian LSTM Networks*, ICOLD workshop

Summary



$$\mathbf{c}_t = [c_1, c_2, \dots, c_C]_t^T \in \mathbb{R}^C$$

$$\mathbf{h}_t = [h_1, h_2, \dots, h_C]_t^T \in (-1, 1)^C$$

f_t : Forget gate $\in (0, 1)^C$

i_t : Input gate $\in (0, 1)^C$

$\tilde{\mathbf{c}}_t$: Candidate values $\in (-1, 1)^C$

$$\mathbf{c}_t = \underbrace{f_t \odot \mathbf{c}_{t-1}}_{\text{forget}} + \underbrace{i_t \odot \tilde{\mathbf{c}}_t}_{\text{write}}$$

o_t : Output gate $\in (0, 1)^C$

$$\mathbf{h}_t = \underbrace{o_t \odot \tanh(\mathbf{c}_t)}$$

Attention on a subset of the cell state