

Introduction to Machine Learning: Deep Supervised Learning

Stefania Sarno

January 6, 2021

Outline

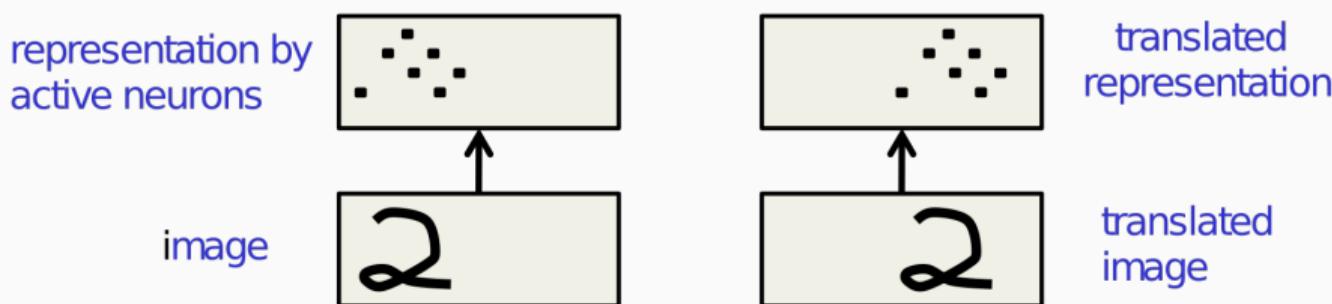
1. Convolutional Neural Networks (CNNs)
2. Recurrent Neural Networks (RNNs)
3. Long-Short Term Memory (LSTM)

What are CNNs?

- CNNs are a class of deep neural network that represents the deep learning framework for computer vision.
- They are also known as shift invariant or space invariant artificial neural networks (SIANN), based on shared-weights architecture and translation invariance characteristics.
- They are powerful tools for image segmentation, object detection, object recognition

Why Object Recognition is difficult?

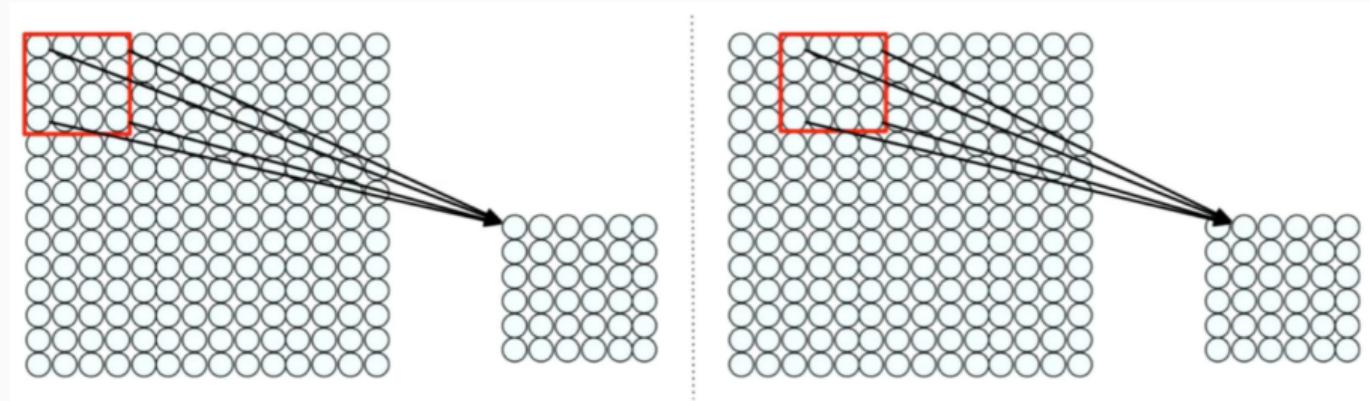
- Segmentation, Lighting, Deformation, Affordances
- Viewpoint (translation): changing the viewpoint would completely change the activation of the input layer in a multilayer perceptron



- Fully connected networks do not have spatial information and have plenty of parameters

How can we preserve spatial information?

Connect patches of inputs to neurons in hidden layers



We want to have patches that are able to detect particular aspects of the image

Feature Extraction with Convolution

We want the following

1. Apply a filter to extract spatial information
2. Apply different filters in each convolutional layer
3. Learn each filter using shared parameters

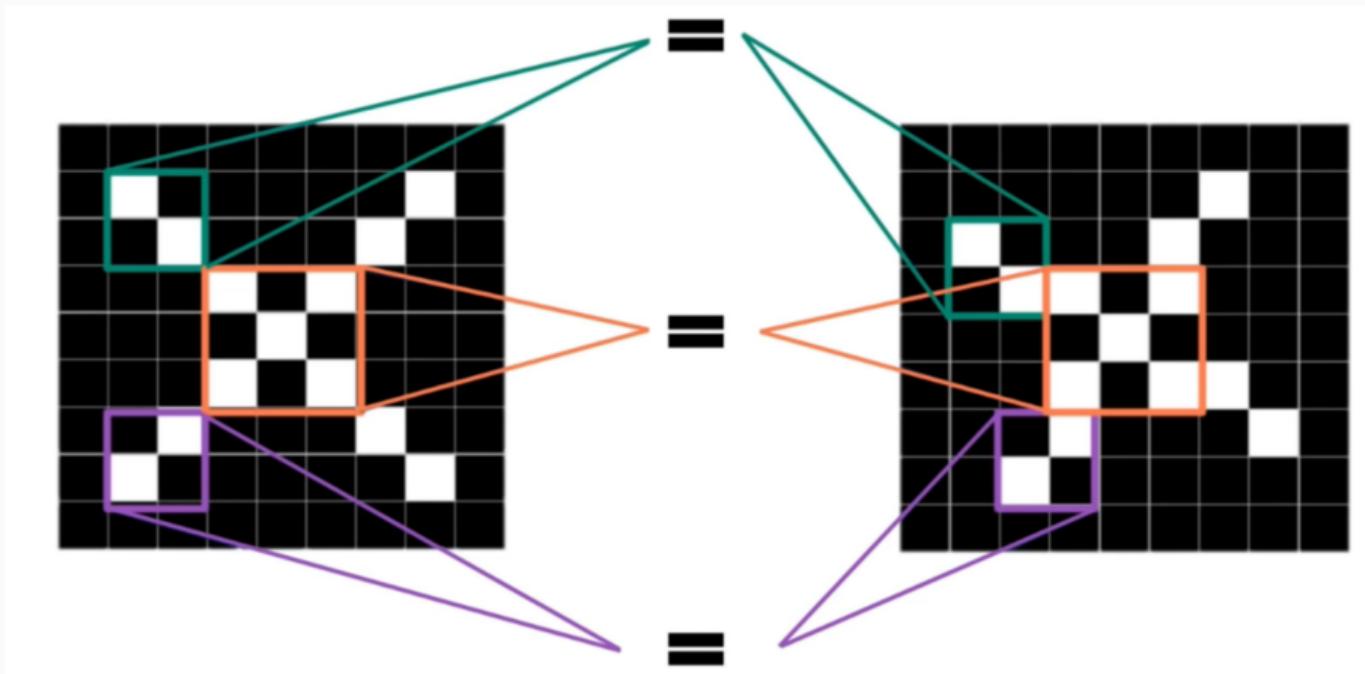
Relevant Features of a X

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	1	-1	-1	1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

Relevant Features of a X

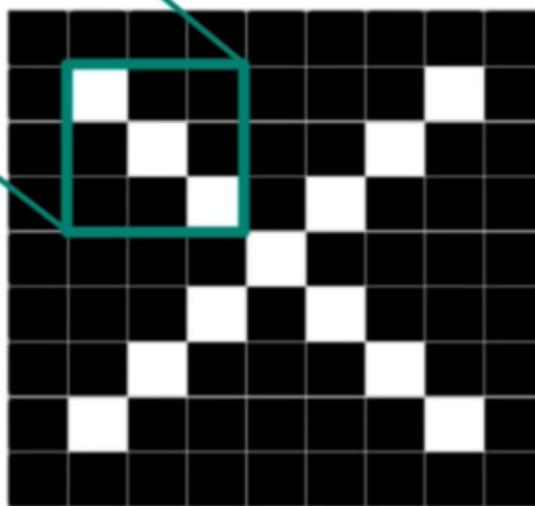


Relevant Features of a X

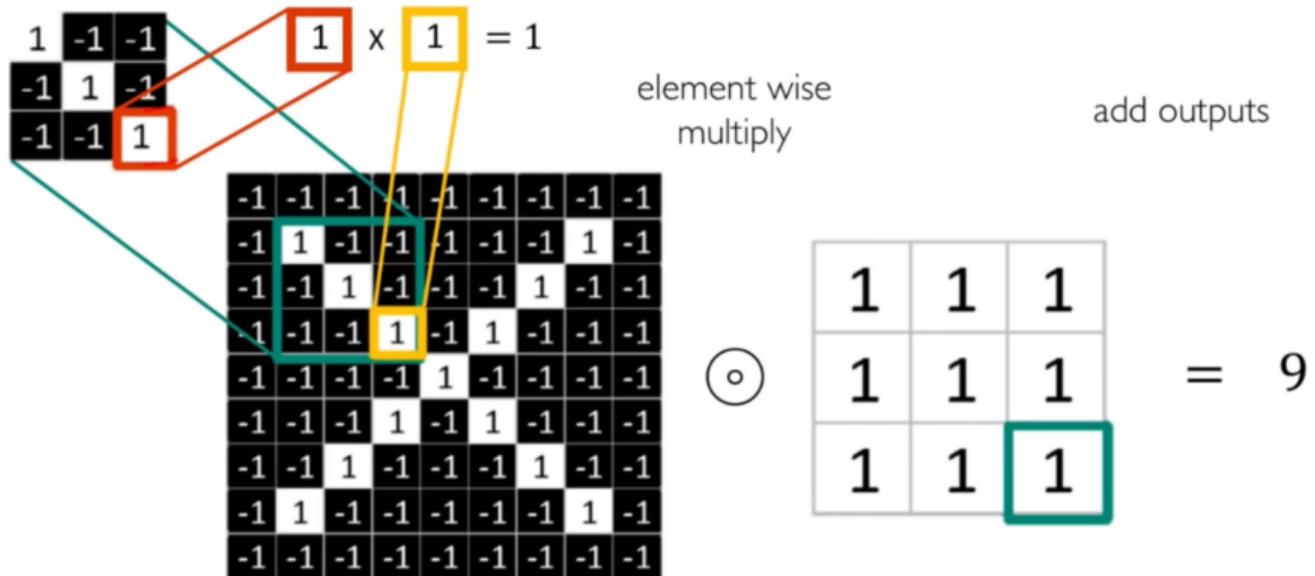
1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

-1	-1	1
-1	1	-1
1	-1	-1



The Convolution Operation



The Convolution Operation

1 x1	1 x0	1 x1	0	0
0 x0	1 x1	1 x0	1	0
0 x1	0 x0	1 x1	1	1
0	0	1	1	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

filter



4		

feature map

The Convolution Operation

1	1	1	0	0
0	1	1	1	0
0	0	1 <small>×1</small>	1 <small>×0</small>	1 <small>×1</small>
0	0	1 <small>×0</small>	1 <small>×1</small>	0 <small>×0</small>
0	1	1 <small>×1</small>	0 <small>×0</small>	0 <small>×1</small>



1	0	1
0	1	0
1	0	1

filter



4	3	4
2	4	3
2	3	4

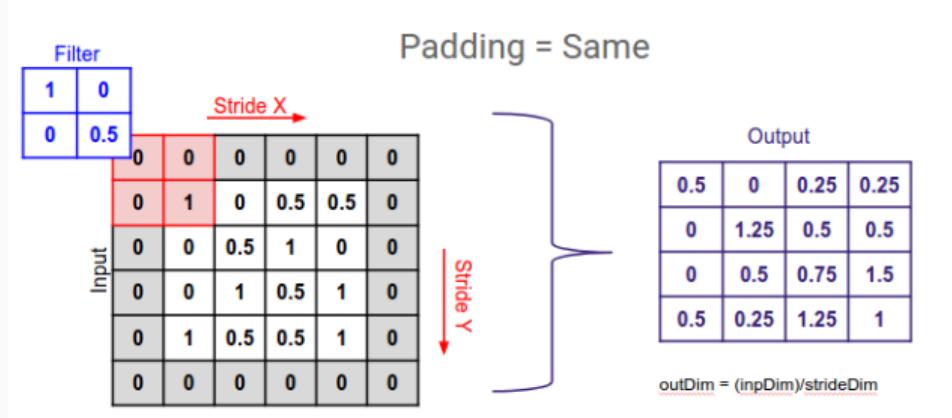
feature map

Padding before convolution

After convolution The dimension of the output O image (either height or width is):

$$O = \frac{W - K - 2P}{S} + 1 \quad (1)$$

where W = input dimension, K = filter size, P = padding, S = stride.



Without padding

1. Shrinking outputs
2. Loosing information on corners of the image

The Feature Maps



Original



Sharpen



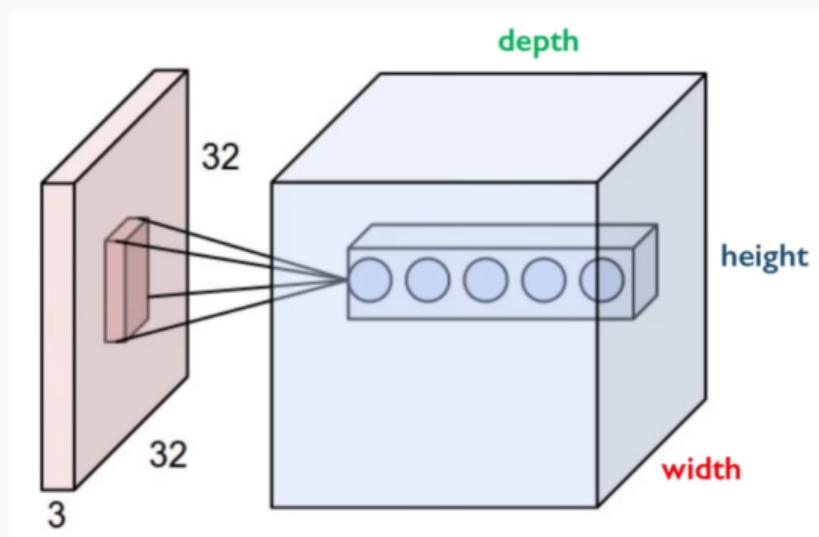
Edge Detect



"Strong" Edge
Detect

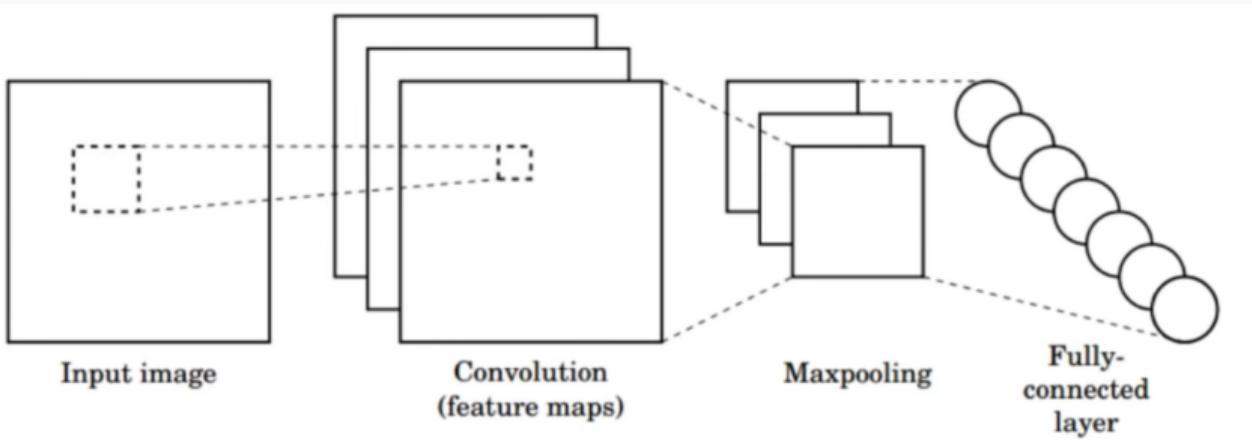
Spatial Arrangement of the Output Volume

- width and height depend on the input volume and on the filter size
- depth represents the number of weights



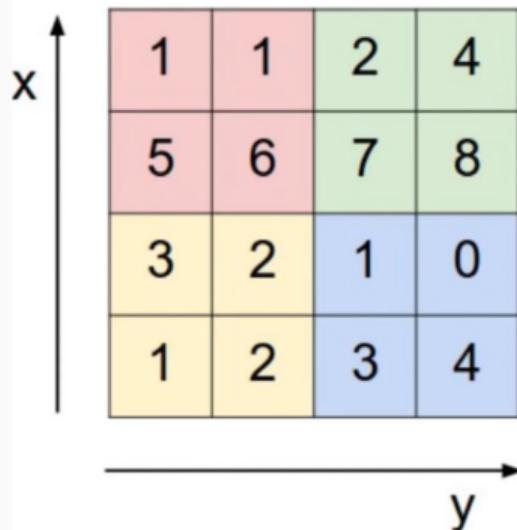
CNN for Classification

1. Multiple Convolutional Layers:
 - Apply a convolutional filter
 - Apply a non-linear function (ReLU)
 - Perform Maxpooling
2. Multiple Fullyconnected Layers
3. A softmax output



Pooling

- Reduce dimensionality
- Preserve spatial invariance



max pool with 2x2 filters
and stride 2

A 2x2 grid representing the output of max pooling. The values are:

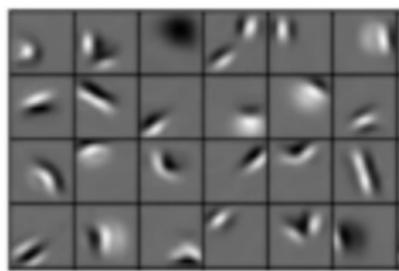
6	8
3	4

The output grid is also color-coded by value: 6 is pink, 8 is light green, 3 is yellow, and 4 is light blue.

- Reduced dimensionality
- Spatial invariance

Feature Map in CNNs

Low level features



Edges, dark spots

Mid level features



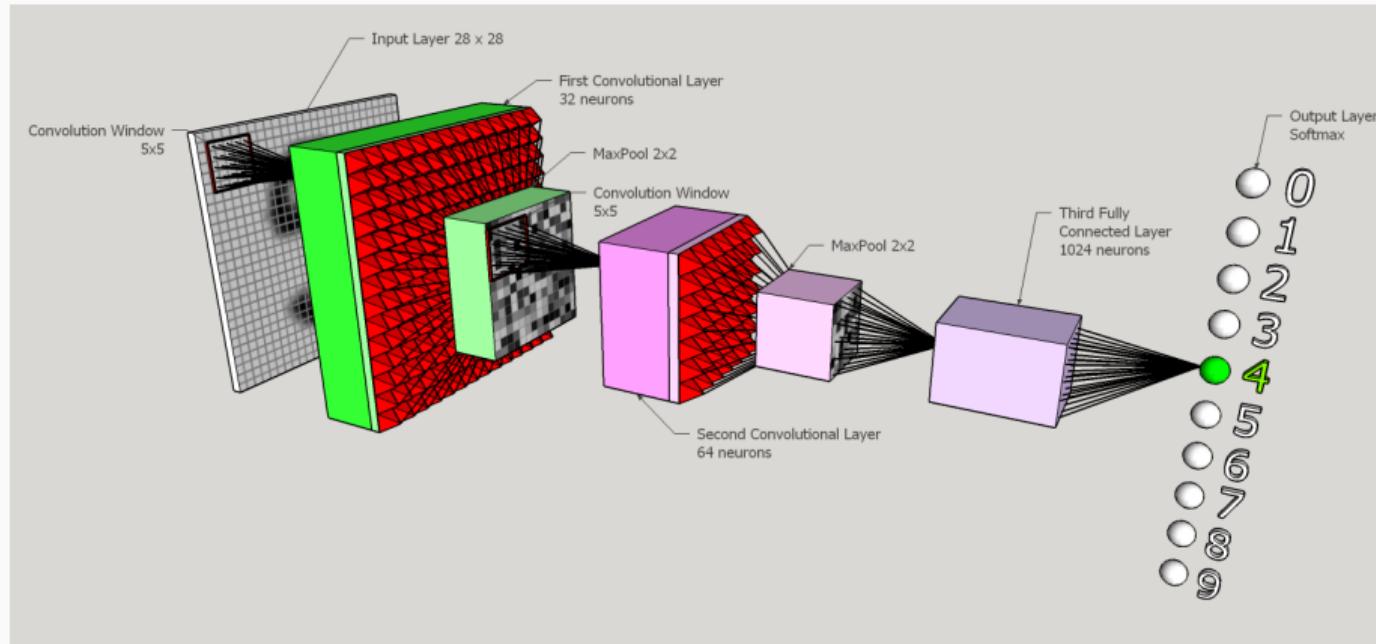
Eyes, ears, nose

High level features



Facial structure

A CNN for the MNIST dataset



Why do we need RNNs?

To model sequential data (language model).

- Machine translation (text to text)

Traduction

Français

Anglais

Arabe

Détecter la langue



J'aime les réseaux de neurones performants.



43/5000

Désactiver la traduction instantanée



Anglais

Français

Arabe

Traduire

I like high-performance neural networks.



Suggérer une modification

Why do we need RNNs?

To model sequential data (language model).

- Image captioning (image to text)



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

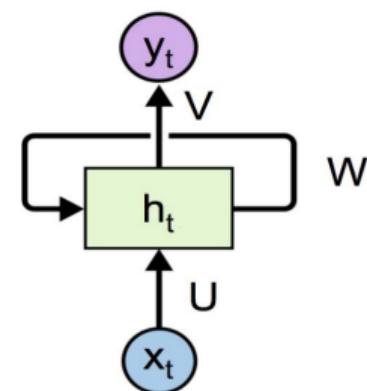
What are RNNs?

Vanilla Recurrent Networks

- The network is defined by three sets of parameters: U, W, V
- The parameters are shared across time

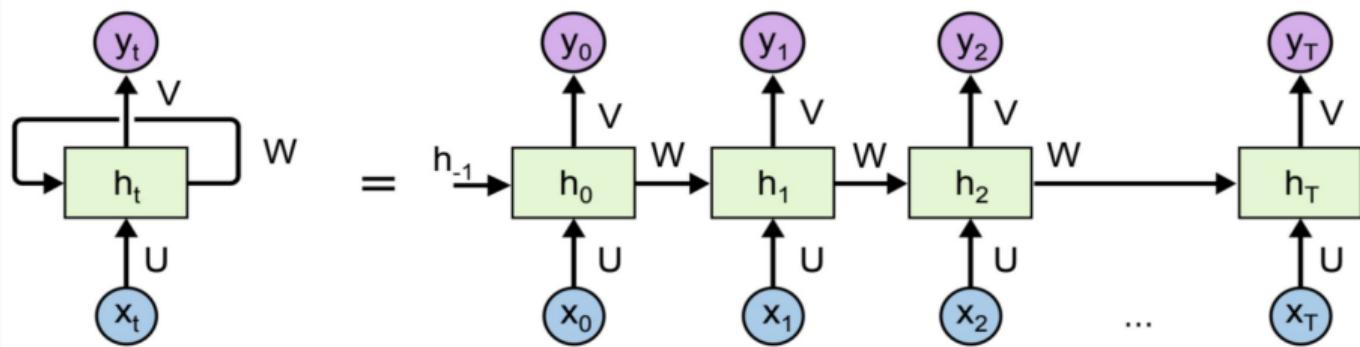
$$h_t = \tanh(Ux_t + Wh_{t-1})$$
$$y_t = f(Vh_t)$$

Output
Intermediate state
Input



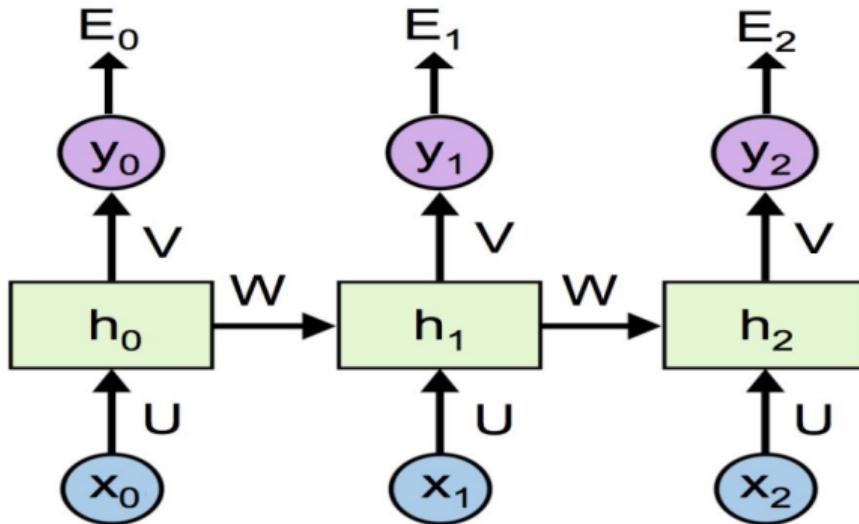
The unrolled network

Vanilla Recurrent Networks



Training RNNs

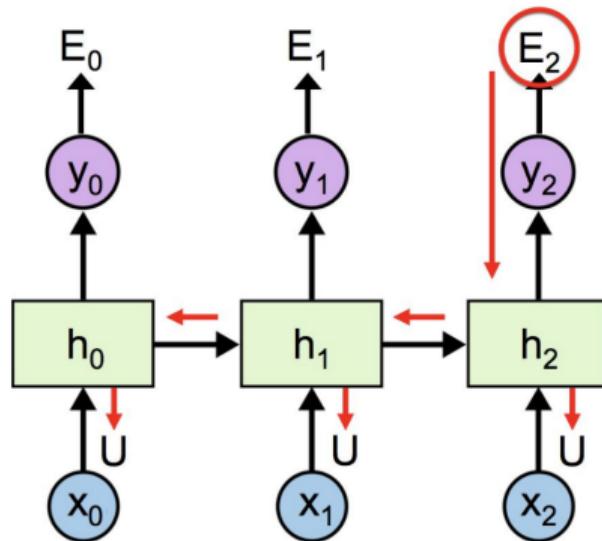
- The loss is: $E = \sum_{t=0}^T E_t$
- Do normal backpropagation using the loss E



Backpropagation through time (BPTT)

$$\frac{\partial E_2}{\partial U} = \frac{\partial E_2}{\partial h_2} \left(x_2^T + \frac{\partial h_2}{\partial h_1} \left(x_1^T + \frac{\partial h_1}{\partial h_0} x_0^T \right) \right)$$

Image from Christopher Olah's blog



Vanishing/Exploding gradient problem

$$\begin{aligned}\mathbf{h}_t &= \theta\phi(\mathbf{h}_{t-1}) + \theta_x \mathbf{x}_t \\ \mathbf{y}_t &= \theta_y \phi(\mathbf{h}_t)\end{aligned}$$

$$\frac{\partial E}{\partial \theta} = \sum_{t=1}^S \frac{\partial E_t}{\partial \theta}$$

$$\frac{\partial E_t}{\partial \theta} = \sum_{k=1}^t \frac{\partial E_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \frac{\partial \mathbf{h}_k}{\partial \theta}$$

[Yoshua Bengio et al]

Vanishing/Exploding gradient problem

$$\frac{\partial E_t}{\partial \theta} = \sum_{k=1}^t \frac{\partial E_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \frac{\partial \mathbf{h}_k}{\partial \theta}$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} = \prod_{i=k+1}^t \theta^T \text{diag}[\phi'(\mathbf{h}_{i-1})]$$

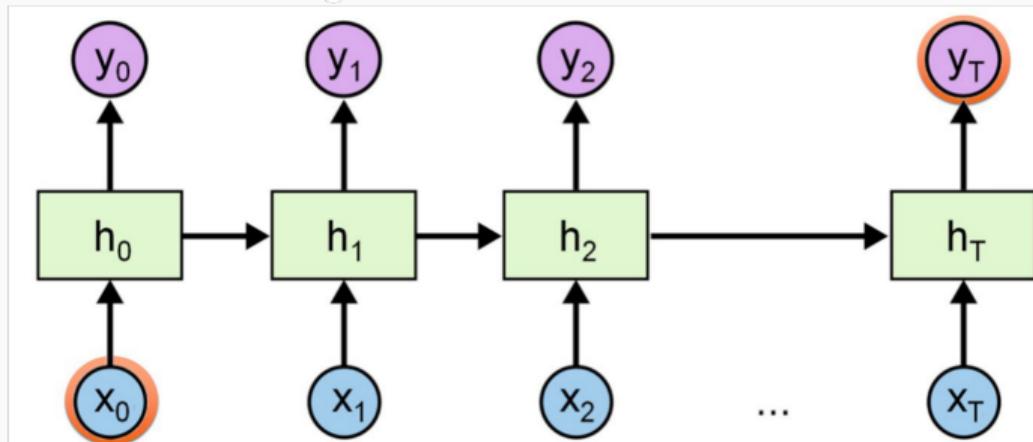
$$\left\| \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right\| \leq \|\theta^T\| \|\text{diag}[\phi'(\mathbf{h}_{i-1})]\| \leq \gamma_\theta \gamma_\phi$$

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \right\| \leq (\gamma_\theta \gamma_\phi)^{t-k}$$

The gradient problem is problematic

Long term dependencies
de Montréal
des algorithmes

Who went to Paris?



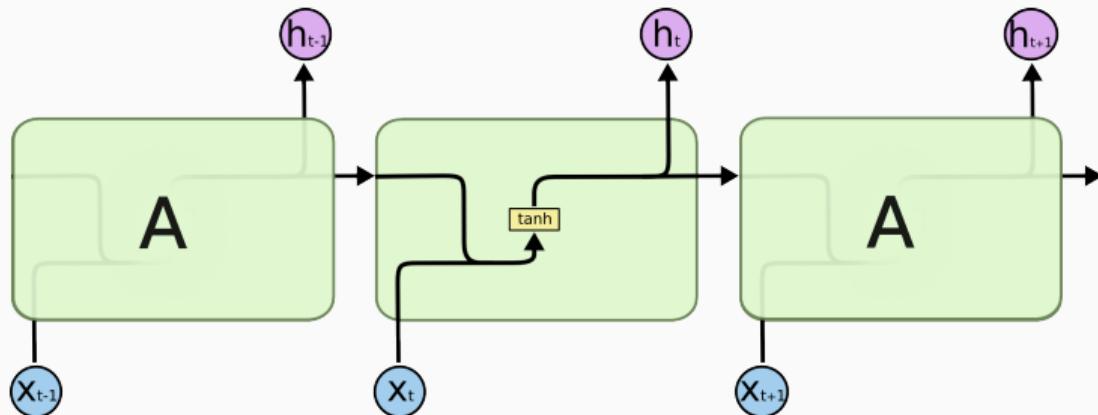
Yoshua

went

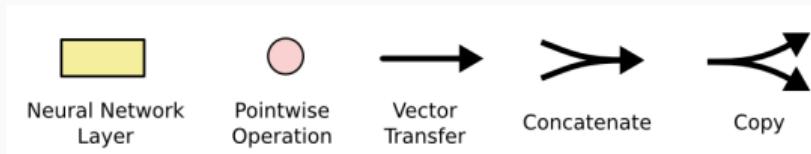
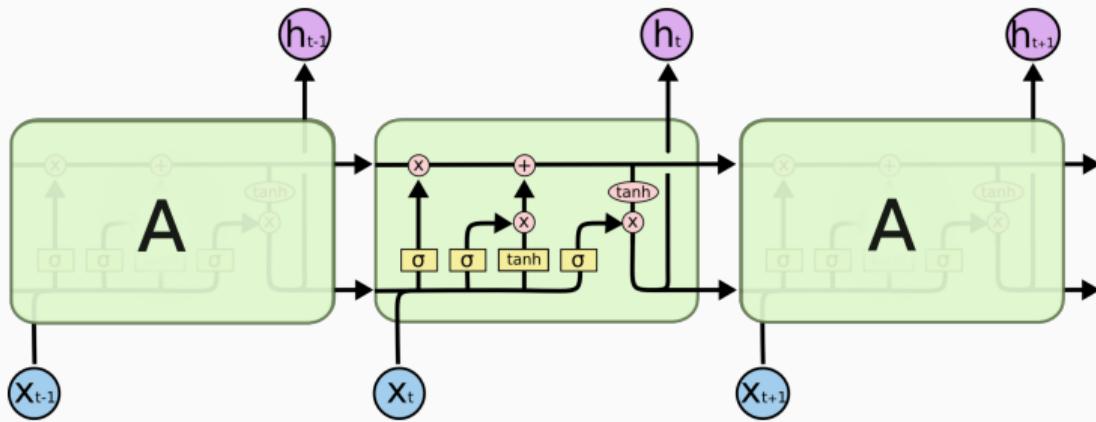
to

Paris

Long short-term memory (LSTM)

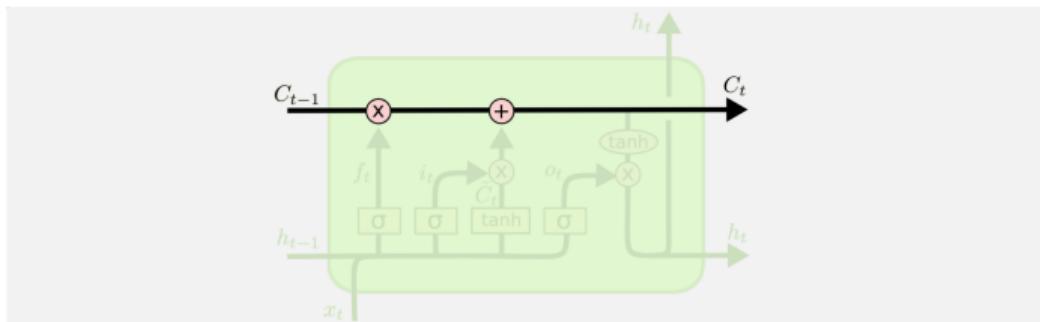


Long short-term memory (LSTM)



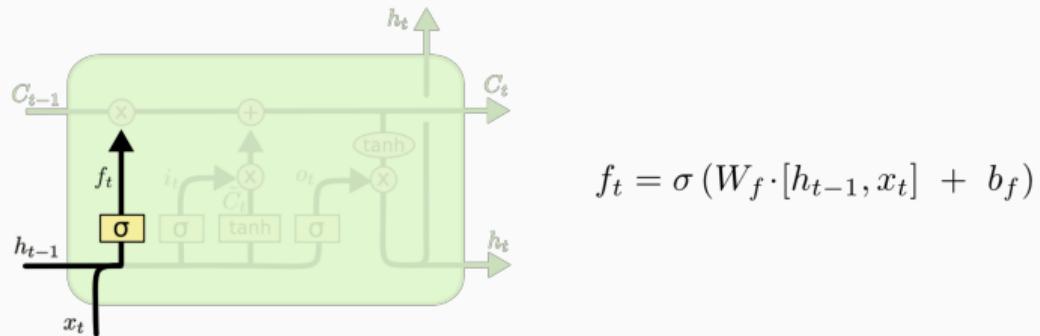
Long short-term memory (LSTM)

Introduce a cell state



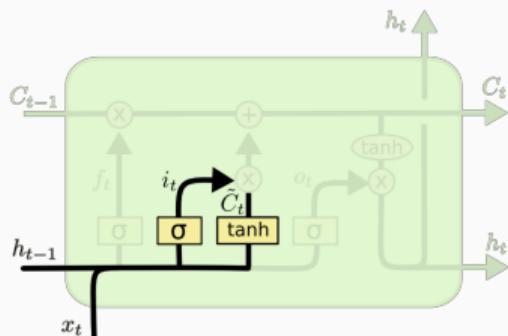
Long short-term memory (LSTM)

Introduce a forget gate (how much to keep in the cell state)



Long short-term memory (LSTM)

Introduce an input gate (how much of the input is added to the cell state)

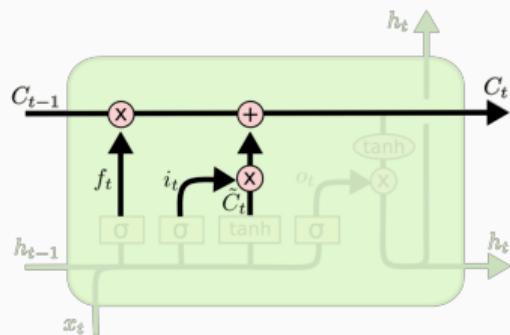


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Long short-term memory (LSTM)

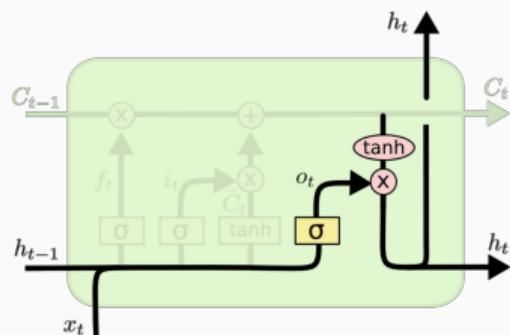
Update the cell state



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Long short-term memory (LSTM)

Produce an output



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$