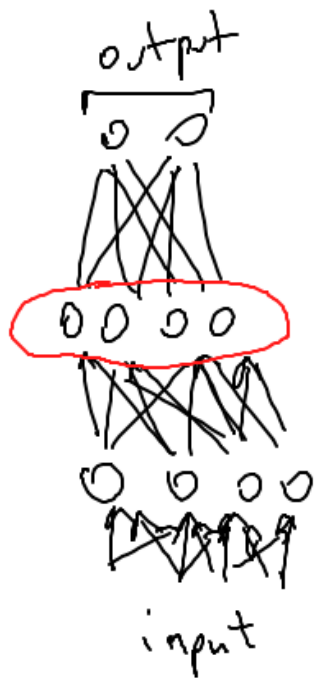
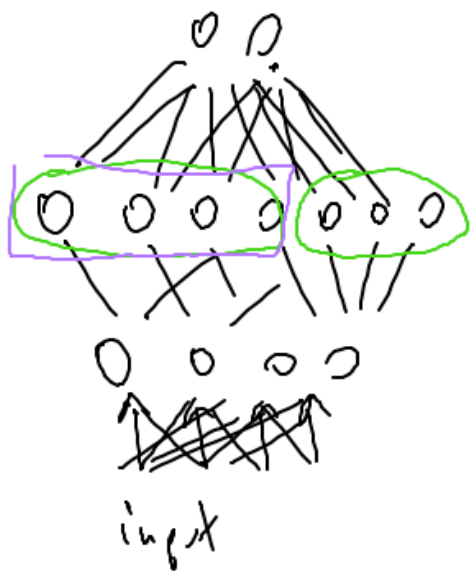


Dropout



or maybe ?



Let us say that we temporarily turn off the 4 neurons in .

Then during the training we need to compensate that fact, and multiply the output from the remaining 3 neurons by $\frac{7}{3}$.

At the end we will have the sum of outputs of all 7 neurons,

this sum = (the sum of 4) + (the sum of 3) = $\frac{7}{4} (\quad) \cdot \frac{4}{7} + \frac{7}{3} \cdot (\quad) \cdot \frac{3}{7}$

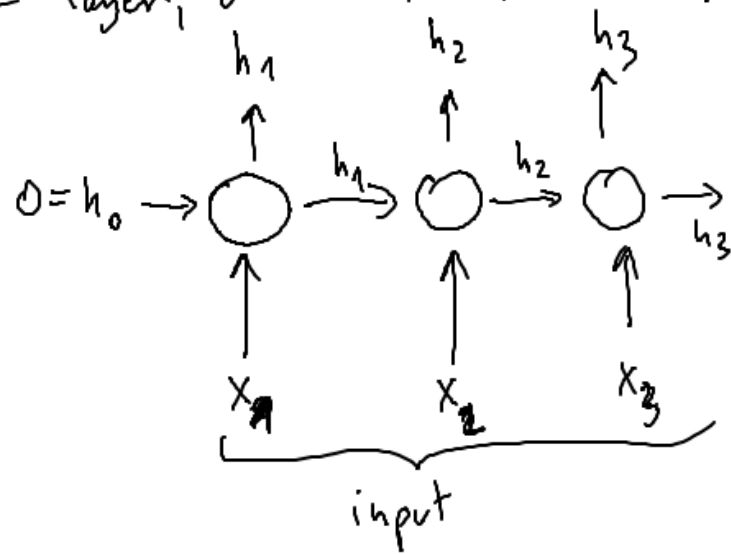
In the dropout, we randomly select some of the neurons in a layer that are going to be turned off (for that particular batch)

Usually we turn off 25% - 50%.



Recurrent Neural Networks (RNN)

one layer, one cell, input = sequence of length 3



$$h_t \in \mathbb{R}$$

$\in \mathbb{R}$

trainable parameters (weights)

Simple RNN cell:
$$h_t = \sigma(W \cdot x_t + U \cdot h_{t-1} + b)$$

$W, U, b \in \mathbb{R}$

↑ inputs to the cell

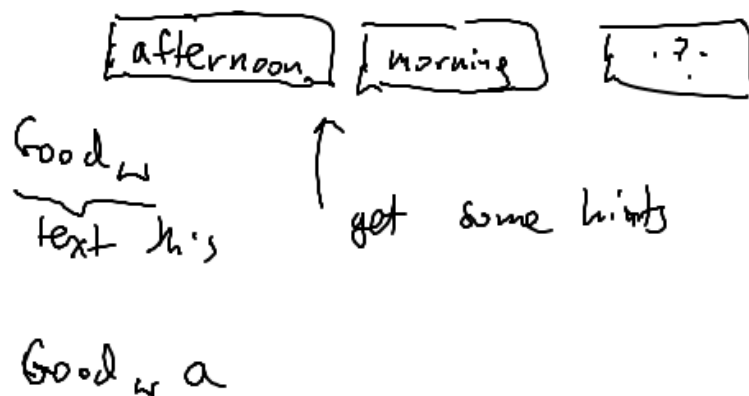
e.g.

$$h_1 = \sigma(W \cdot x_1 + U \cdot h_0 + b)$$
$$h_2 = \sigma(W \cdot x_2 + U \cdot h_1 + b)$$
$$h_3 = \sigma(W \cdot x_3 + U \cdot h_2 + b)$$

One can treat h_3 as ~~the~~ a prediction for the fourth element of the sequence (x_1, x_2, x_3)

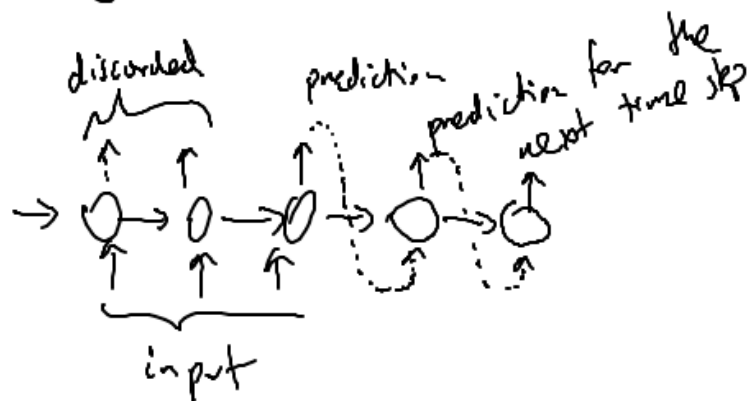
Some applications:

- Texting on a smartphone

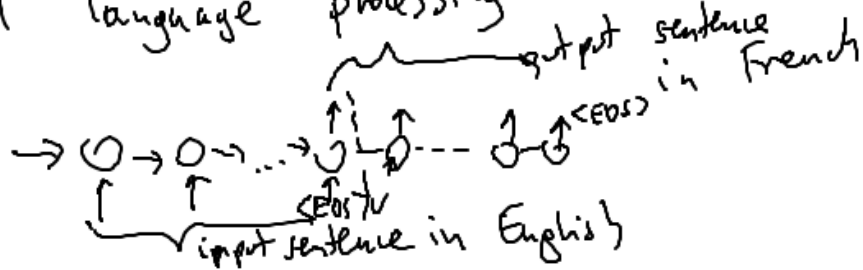


- Predicting some time series

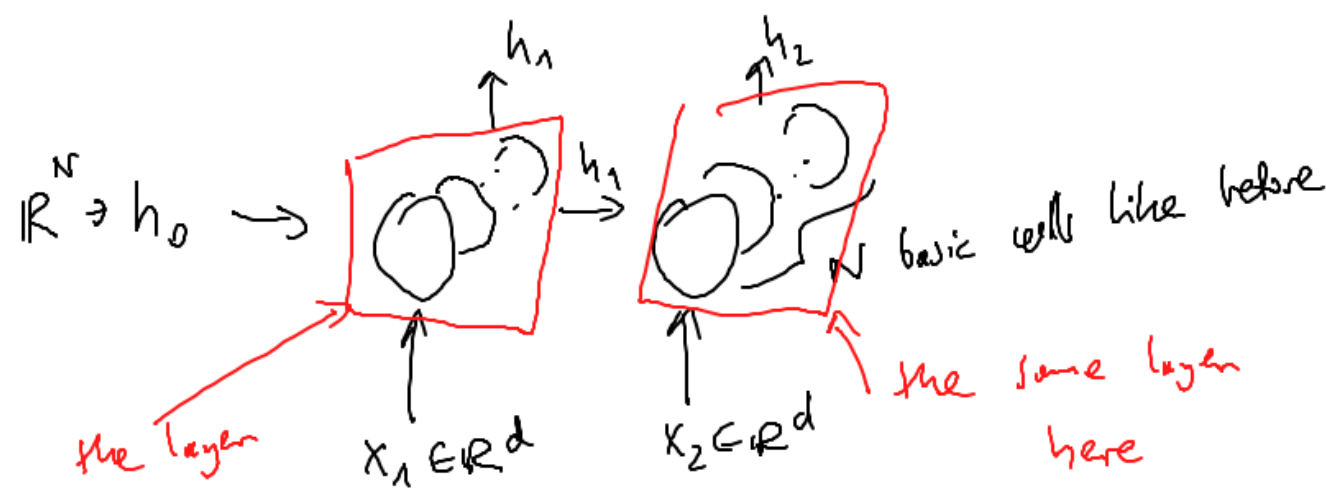
tidal



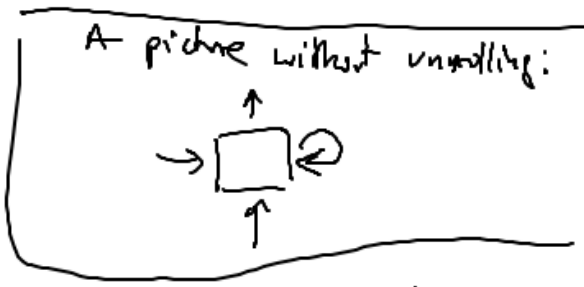
- * In natural language processing



Generalisation to a multidimensional case with multiple neurons in a layer



← one ^{single RNN} layer
 ← picture with unrolling




$$h_t = \sigma \left(\underbrace{W \cdot x_t + U \cdot h_{t-1} + b}_{\in \mathbb{R}^N} \right)$$

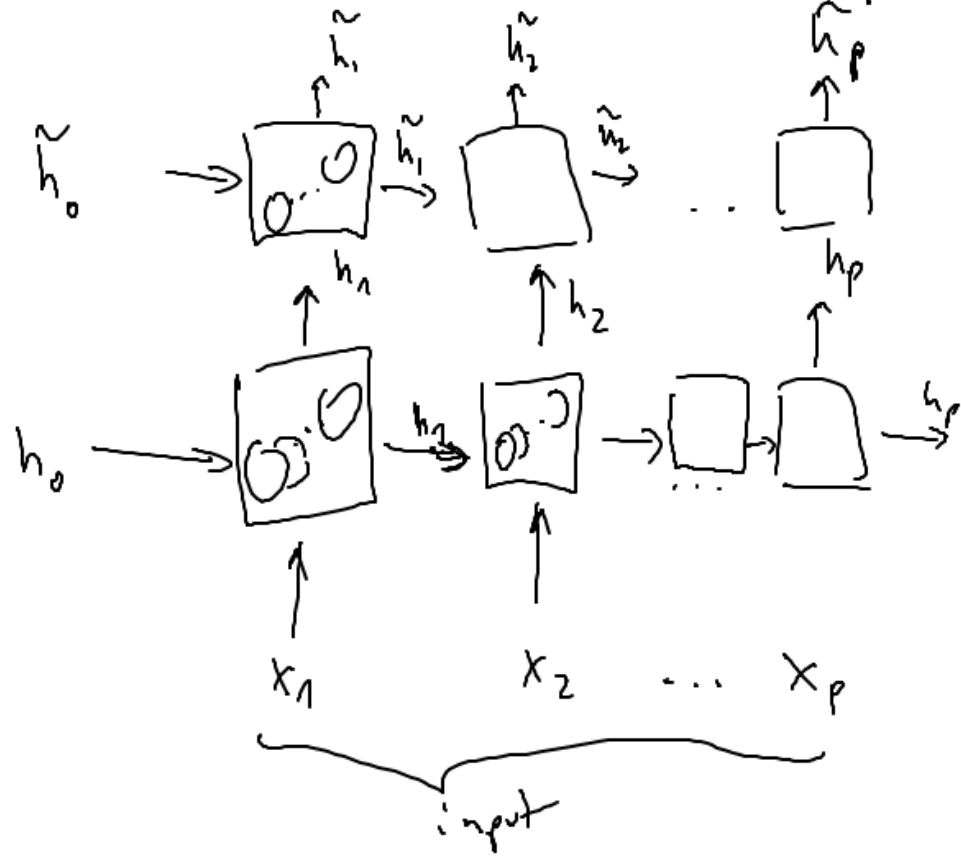
Annotations for the equation above:

- W : matrix $N \times d$
- x_t : $d \times 1$
- U : matrix $N \times N$
- h_{t-1} : $N \times 1$
- b : \mathbb{R}^N

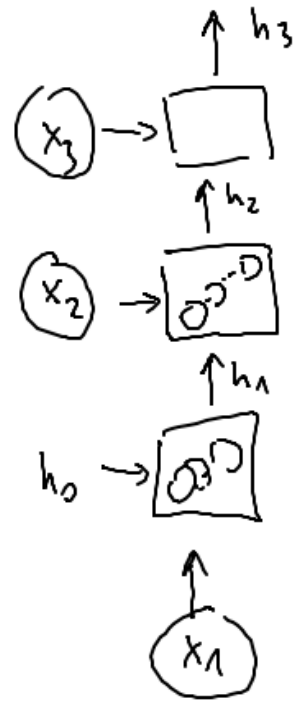
$x_t \in \mathbb{R}^d$

$$\begin{bmatrix} \sigma(\dots) \\ \sigma(\dots) \\ \dots \\ \sigma(\dots) \end{bmatrix} \in \mathbb{R}^N$$

One can have multiple layers:  eg. dense layer



Other way of looking at a single RNN layer:



When a sequence has 3 elements,
then the RNN layer is similar
to a classical NN with 3 layers.

Usual problems with
vanishing/exploding gradients.

LSTM (Long-Short Term Memory)

$$f_t = \sigma_g \left(\underbrace{W_f}_{N \times d} \cdot \underbrace{x_t}_{R^d} + \underbrace{U_f}_{N \times N} \cdot h_{t-1} + \underbrace{b_f}_{R^N} \right)$$

$$i_t = \sigma_g \left(\underbrace{W_i}_{N \times d} \cdot \underbrace{x_t}_{R^d} + \underbrace{U_i}_{N \times N} \cdot h_{t-1} + \underbrace{b_i}_{R^N} \right)$$

$$o_t = \sigma_g \left(\underbrace{W_o}_{N \times d} \cdot \underbrace{x_t}_{R^d} + \underbrace{U_o}_{N \times N} \cdot h_{t-1} + \underbrace{b_o}_{R^N} \right)$$

$$\tilde{c}_t = \sigma_c \left(\underbrace{W_c}_{N \times d} \cdot \underbrace{x_t}_{R^d} + \underbrace{U_c}_{N \times N} \cdot h_{t-1} + \underbrace{b_c}_{R^N} \right)$$

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

$$h_t = o_t \cdot \sigma_h(c_t)$$

$$\left\{ \begin{array}{l} x_t \in \mathbb{R}^d, \\ f_t, i_t, o_t, h_t \in \mathbb{R}^N \end{array} \right.$$

$d=N=1 \rightarrow 12$ weights

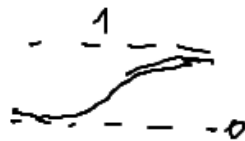
95-97 - 99

Hochreiter, Schmidhuber

forget gate

input gate

output gate



Standard choice:

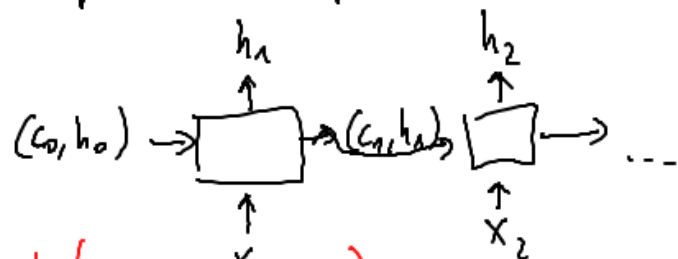
$$\sigma_g(x) = 1/(1+e^{-x})$$

$$\sigma_c(x) = \sigma_h(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

\odot = element-wise product

← the memory cell

← the output and the hidden state



$$4(N \cdot d + N^2 + N)$$