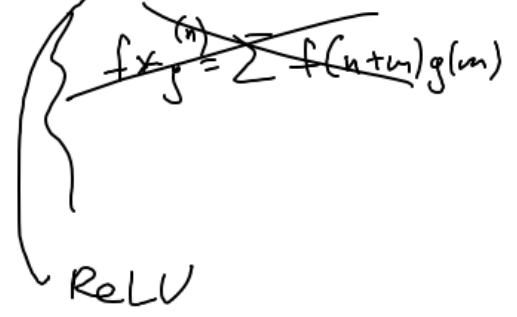
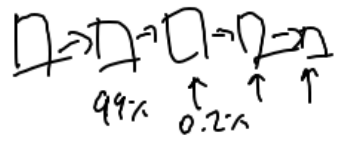


filter of size $k \times k$,
 but in fact it is
 of size $F \times k \times k + 1$
 number of weights in the filter



$\varphi(x) = \max(x, 0)$



Naively, how many operations do we need to calculate the output of the layer?

For each entry of the output: $F k^2$ multiplications
 # entries in the output: $(m-k+1)(n-k+1) \cdot G \approx mnG$

(G channels)

$\rightarrow \text{total} \approx F k^2 mn G$
 $128 \cdot 3^2 \cdot 60 \cdot 60 \cdot 128 \approx 10^3 \cdot 10^4 \cdot 128 \approx 10^9$

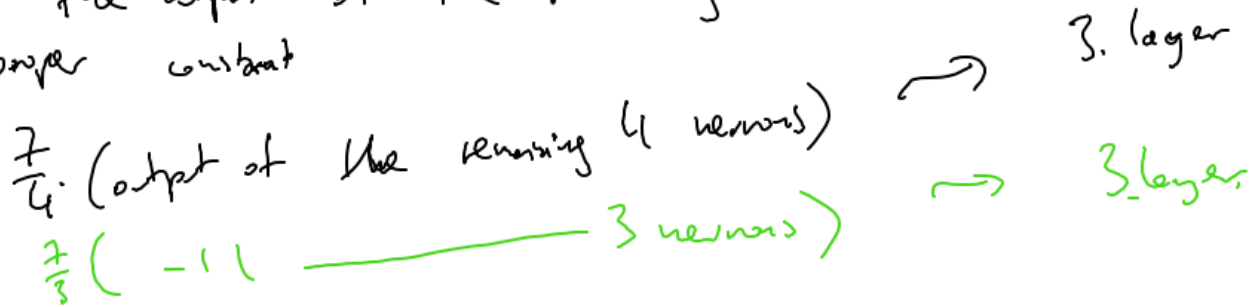
Dropout

cheap way of averaging multiple neural networks
(computationally)

For the prediction phase, we use all neurons

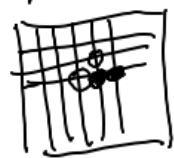


to compensate this turn off, one multiplies the output of the remaining neurons by a proper constant



DeepForest

input



$19^2 = 361$ intersection

predicting moves in Go

transformed



CNN

same padding
without stride
or pooling

3 19x19



probability
the distribution
of the moves

- move
- next
- afterwards

{ white stone,
 black stones
 empty intersections
 ...

the best human chess player (Kasparov)

2017 - program + hardware won with
 2017 - program + -11- won with -11- Go
 Alpha Go

- purely CNN
- same size



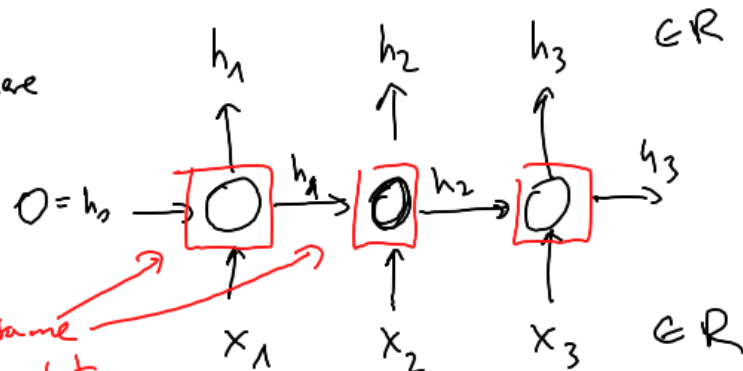
Recurrent Neural Networks (RNN)

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

Simple RNN:

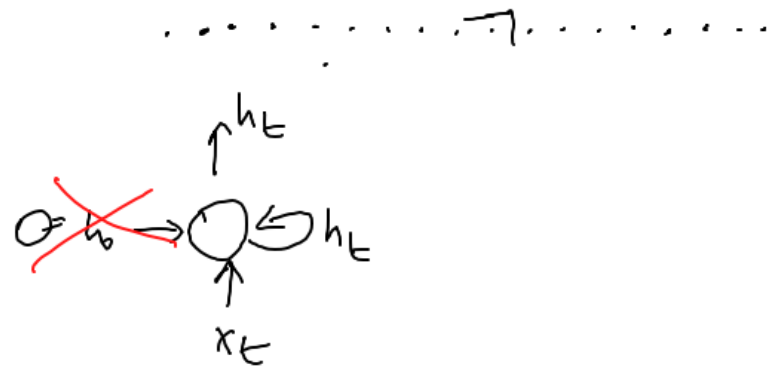
possibly
other layer too

layer



the same
entity, but
at different 'time'

unrolled
picture

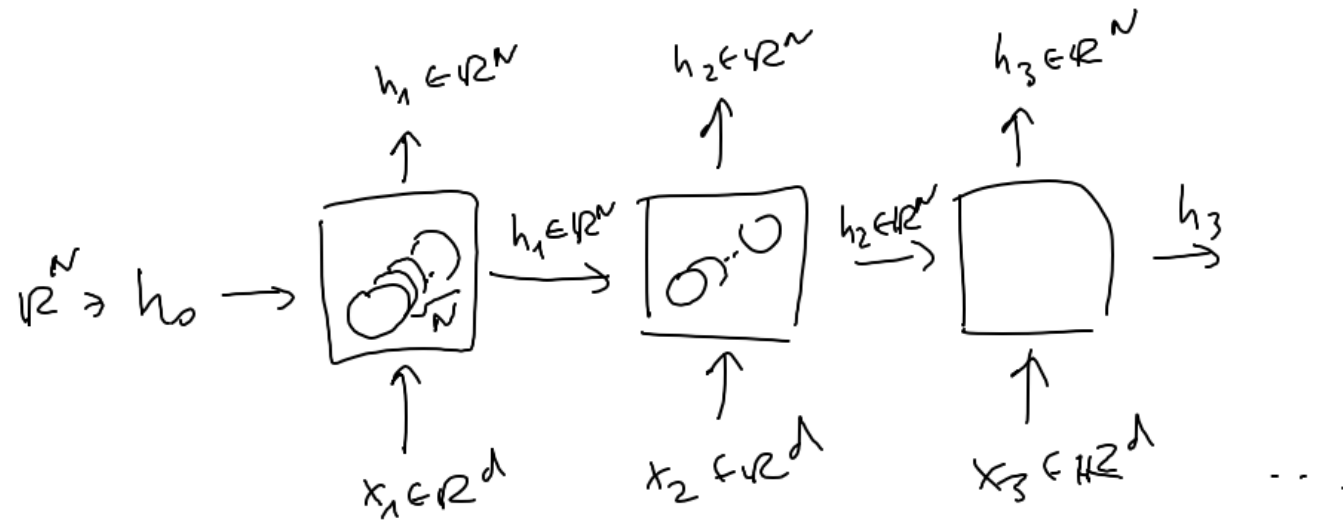


trainable parameters
(weights)

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

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

a more general case: $x_t \in \mathbb{R}^d$, N cells (neurons) in a layer



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

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

↑
↑
↑

[]
[]
[]

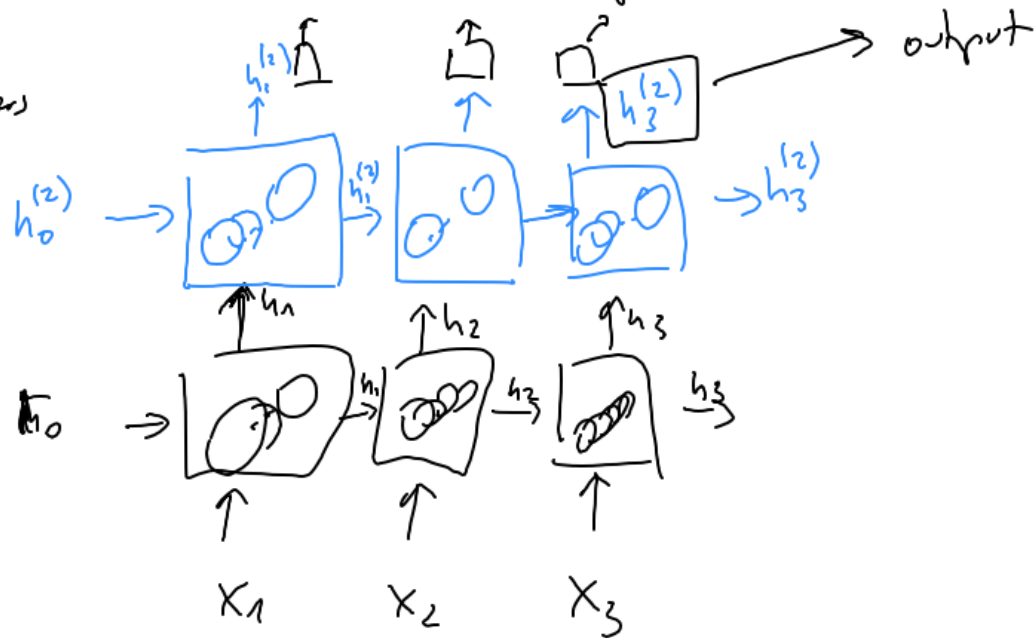
↑
↑
↑

[]
[]
[]

↑
↑

these can be stacked together

possible dense layers



x_1, x_2, x_3, x_4
 feed the network,
 get some output y
 the expected output

he
 she

his
 her

LSTM (Long-Short Term Memory)

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

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

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

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

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

$$h_t = o_t \odot \sigma_h(c_t) \in \mathbb{R}^N$$

if $f_t^{(j)} \approx 1$ and $i_t^{(j)} \approx 0$, then
 $f_t^{(j)} \approx 0$ and $i_t^{(j)} \approx 1$, then

95-97-99

Hoch reiten, Schmidhauben

forget gate

input gate

output gate

$\sigma_f = \text{sigmoid}$

$\sigma_c = \text{tanh}$

memory cell

\odot coordinate-wise multiplication

output

j th coordinate

$c_t^{(j)} \approx c_{t-1}^{(j)}$ (old value is kept)
 $c_t^{(j)} = \tilde{c}_t^{(j)}$ (old value is discarded)