#### Sequences Part I: Recurrent Neural Networks (RNNs) COMP 551 Applied Machine Learning

Isabeau Prémont-Schwarz School of Computer Science McGill University

Fall 2024



### Deep Neural Networks

- Neural Networks for Tabular Data
  - MLP
- Neural Networks for Images
  - CNN
- Neural Networks for Sequences
  - The input is a sequence, the output is a sequence, or both are sequences.

English V	$\stackrel{\rightarrow}{\leftarrow}$	French $\checkmark$	Automatic $\checkmark$	Glossary
Lets learn how to translate a sentence that is a sequence of words	×	Apprenons à traduire une phrase qui est une séquence de mots.		

### Learning Objectives

- Tokenization
- Sequence Learning
  - Seq2Vec
  - Vec2Seq
  - Seq2Seq
- Recurrent Neural Networks (RNNs)
  - Vanilla RNN
  - Gated Recurrent Unit (GRU)
  - Backpropagation through time (BPTT)

### • Seq2Vec (sequence classification)

- Seq2Vec (sequence classification)
- Vec2Seq (sequence generation)

- Seq2Vec (sequence classification)
- Vec2Seq (sequence generation)
- Seq2Seq (sequence translation)

### Seq2Vec (sequence classification)

#### Example:



#### Vec2Seq (sequence generation) Example: Caption Generation:



Examples:

Examples:

• Language translation (eg. Fr  $\rightarrow$  En)

Examples:

- Language translation (eg. Fr  $\rightarrow$  En)
- Protein folding (Sequence of proteins  $\rightarrow$  Sequence of angles)

#### Examples:

- Language translation (eg. Fr  $\rightarrow$  En)
- Protein folding (Sequence of proteins  $\rightarrow$  Sequence of angles)
- Time series (eg. Stock price prediction)

Examples:

- Language translation (eg. Fr  $\rightarrow$  En)
- Protein folding (Sequence of proteins  $\rightarrow$  Sequence of angles)
- Time series (eg. Stock price prediction)

Note: We also distinguish between aligned and non-aligned Seq2Seq models.

#### 1D Convolution

• Have a state which keeps track of past information.

- Have a state which keeps track of past information.
- Have an special token < EOS > which designates the end of sequence.

- Have a state which keeps track of past information.
- Have an special token < EOS > which designates the end of sequence.
- Potentially also a token < SOS > which designates the start of sequence.





$$x_t \in \mathbb{R}^d, h_t \in \mathbb{R}^m, o_t \in \mathbb{R}^n$$



Matrix shapes: 
$$U: m \times d, V: m \times m, W: n \times m$$



Matrix shapes:  $U: m \times d, V: m \times m, W: n \times m$ 

$$h_t = f\left(V \cdot h_{t-1} + U \cdot x_t + b_h\right)$$

where f is the non-linear activation function (tanh is often used, but ReLU is also possible)



Matrix shapes:  $U: m \times d, V: m \times m, W: n \times m$ 

$$h_t = f\left(V \cdot h_{t-1} + U \cdot x_t + b_h\right)$$

where f is the non-linear activation function (tanh is often used, but ReLU is also possible)

$$o_t = g\left(W \cdot h + b_o\right)$$

#### Vanilla RNN: example

$$egin{aligned} h_t &= \operatorname{ReLU}\left(V \cdot h_{t-1} + U \cdot x_t + b_h
ight) \ o_t &= \left(W \cdot h_t + b_o
ight), \qquad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1 \end{aligned}$$

#### Vanilla RNN: example

$$egin{aligned} h_t &= \mathsf{ReLU}\left(V \cdot h_{t-1} + U \cdot x_t + b_h
ight) \ o_t &= \left(W \cdot h_t + b_o
ight), \qquad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1 \end{aligned}$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

#### Vanilla RNN: example

$$egin{aligned} h_t &= \mathsf{ReLU}\left(V \cdot h_{t-1} + U \cdot x_t + b_h
ight) \ o_t &= \left(W \cdot h_t + b_o
ight), \qquad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1 \end{aligned}$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1\\0\\0 \end{bmatrix}, x_t = \begin{bmatrix} 1\\0 \end{bmatrix}$$

What is *h<sub>t</sub>*?

$$h_t = \mathsf{ReLU}\left(V \cdot h_{t-1} + U \cdot x_t + b_h
ight)$$

# Vanilla RNN: example $h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$ $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$
$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$h_t = \operatorname{ReLU}(V \cdot h_{t-1} + U \cdot x_t)$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$h_t = \mathsf{ReLU}\left( \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + U \cdot x_t \right)$$

Vanilla RNN: example
$$h_t = \text{ReLU}(V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1\\ 0\\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1\\ 0 \end{bmatrix}$$

$$h_t = \mathsf{ReLU}\left(\begin{bmatrix}0\\0\\1\end{bmatrix} + U \cdot x_t\right)$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$h_t = \mathsf{ReLU}\left(\begin{bmatrix}0\\0\\1\end{bmatrix} + \begin{bmatrix}1&0\\0&1\\-1&0\end{bmatrix} \cdot \begin{bmatrix}1\\0\end{bmatrix}\right)$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1\\ 0\\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1\\ 0 \end{bmatrix}$$

$$h_t = \mathsf{ReLU}\left(\begin{bmatrix}0\\0\\1\end{bmatrix} + \begin{bmatrix}1\\0\\-1\end{bmatrix}\right)$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$h_t = \mathsf{ReLU}\left( \begin{bmatrix} 1\\0\\0 \end{bmatrix} \right) = \begin{bmatrix} 1\\0\\0 \end{bmatrix}$$

Vanilla RNN: example
$$h_t = \text{ReLU} (V \cdot h_{t-1} + U \cdot x_t + b_h)$$
 $o_t = (W \cdot h_t + b_o), \quad x_t \in \mathbb{R}^2, h_t \in \mathbb{R}^3, o_t \in \mathbb{R}^1$ 

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, W = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}, b_h = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, b_o = 0$$

$$h_{t-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, x_t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$o_t = W \cdot h_t = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 1$$

- Vec2Seq (sequence generation)
  - output,  $y_{1:T}$  is generated one token at a time
  - at each step we sample  $y_t$  from the hidden state  $h_t$ and then feed it back to the model to get  $h_{t+1}$

arbitrary-length sequence of vectors  $f_{\theta}: \mathbb{R}^D \to \mathbb{R}^{N_{\infty}C}$ 

*D*: input vector size  $N_{\infty}$ : arbitrary-length sequence of vectors of length *C C*: each output vector size

conditional generative model:

$$p(y_{1:T}|x) = \sum\limits_{h_{1:T}} p(y_{1:T},h_{1:T}|x) = \sum\limits_{h_{1:T}} \prod\limits_{t=1}^T p(y_t|h_t) p(h_t|h_{t-1},y_{t-1},x)$$

with the initial hidden state $p(h_1|h0, y0, x) = p(h_1|x)$ 



• Vec2Seq (sequence generation)  $f_{ heta}: \mathbb{R}^D \to \mathbb{R}^{TC}$ 

conditional generative model:

$$p(y_{1:T}|x) = \sum_{h_{1:T}} p(y_{1:T}, h_{1:T}|x) = \sum_{h_{1:T}} \prod_{t=1}^T p(y_t|h_t) p(h_t|h_{t-1}, y_{t-1}, x)$$

. . . . .

• real-valued output: 
$$\hat{y}_t = W_{hy} h_t$$
  
 $p(y_t|h_t) = \mathcal{N}(y_t|\hat{y}_t, \mathbf{I})$ 

- categorical output:  $\hat{y}_t = ext{softmax}(W_{hy}h_t)$  $p(y_t|h_t) = ext{Categorical}(y_t|\hat{y}_t)$ 



• Vec2Seq (sequence generation)  $f_{\theta} : \mathbb{R}^D \to \mathbb{R}^{TC}$ 

conditional generative model:

$$p(y_{1:T}|x) = \sum_{h_{1:T}} p(y_{1:T}, h_{1:T}|x) = \sum_{h_{1:T}} \prod_{t=1}^T p(y_t|h_t) p(h_t|h_{t-1}, y_{t-1}, x)$$

hidden state:  $p(h_t|h_{t-1},y_{t-1},x) = \mathbb{I}(h_t = f(h_{t-1},y_{t-1},x))$  input-to-hidden weights hidden-to-hidden weights  $h_t = arphi (W_{xh}[x;y_{t-1}]+W_{hh}h_{t-1})$ 



• Vec2Seq (sequence generation)

hidden-to-output weights

$$f_{ heta}: \mathbb{R}^D o \mathbb{R}^{TC}$$

model

$$\hat{y}_t = g(rac{W_{hy}}{W_{hy}}h_t)$$
  
input-to-hidden weights hidden-to-hidden weights  $h_t = arphi(W_{xh}[x;y_{t-1}]+W_{hh}h_{t-1})$ 

RNNs are powerful

- In theory can have unbounded memory and are as powerful as a Turing machine
- In practice, memory size is determined by the size of the latent space and strength of the parameters



Vec2Seq (sequence generation)

conditional generative model:

$$p(y_{1:T}|x) = \sum\limits_{h_{1:T}} p(y_{1:T},h_{1:T}|x) = \sum\limits_{h_{1:T}} \prod\limits_{t=1}^T p(y_t|h_t) p(h_t|h_{t-1},y_{t-1},x)$$

**language modelling:** generating sequences unconditionally (by setting  $x = \emptyset$ ) which is learning joint probability distributions over sequences of discrete tokens, i.e.,  $p(y_1, \ldots, y_T)$ 

#### Example:

character level RNN trained on the book The Time Machine by H. G. Wells (32,000 words and 170k character)

### Output when given prefix "the":

**the** githa some thong the time traveller held in his hand was a glitteringmetallic framework scarcely larger than a small clock and verydelicately made there was ivory in it and the latter than s bettyre tat howhong s ie time thave ler simk you a dimensions le ghat dionthat shall travel indifferently in any direction of space and timeas the driver determinesfilby contented himself with laughterbut i have experimental verification said the time travellerit would be remarkably convenient for the histo

#### See the code here, read more here

• Vec2Seq (sequence generation)

conditional generative model:

 $p(y_{1:T}|x) = \sum\limits_{h_{1:T}} p(y_{1:T},h_{1:T}|x)$ 

#### Example:

CNN-RNN model for image captioning when *x* is embedding by a CNN



<start>

See more here

Giraffes standing

<end>

• Seq2Vec (sequence classification)

$$f_{ heta}: \mathbb{R}^{TD} 
ightarrow \mathbb{R}^{C}$$

predict a single fixed-length output vector given a variable length sequence as input  $y \in \{1, \ldots, C\}$ 

use the final state:

$$\hat{y} = ext{softmax}(Wh_T) \ p(y|x_{1:T}) = ext{Categorical}(y|\hat{y})$$



### Bi-directional RNN:

the hidden states of the RNN depend on the past and future context

gives better results



• Seq2Vec (sequence classification)

$$f_{ heta}: \mathbb{R}^{TD} 
ightarrow \mathbb{R}^{C}$$

 predict a single fixed-length output vector given a variable length sequence as input

$$egin{aligned} h^{
ightarrow}_t &= arphi \left( W^{
ightarrow}_{xh} x_t + W^{
ightarrow}_{hh} h^{
ightarrow}_{t-1} 
ight) \ h^{
ightarrow}_t &= arphi \left( W^{
ightarrow}_{xh} x_t + W^{
ightarrow}_{hh} h^{
ightarrow}_{t+1} 
ight) \end{aligned}$$

$$egin{aligned} h_t &= [egin{aligned} h_t^{
ightarrow}, h_t^{
ightarrow}] \ \overline{h} &= rac{1}{T}\sum_{t=1}^T h_t \end{aligned}$$

 $\hat{y} = rac{ ext{softmax}(War{h})}{p(y|x_{1:T}) = ext{Categorical}(y|\hat{y})}$ 

### Bi-directional RNN:

the hidden states of the RNN depend on the past and future context

gives better results



Seq2Vec (sequence classification)

$$f_{ heta}: \mathbb{R}^{TD} 
ightarrow \mathbb{R}^{C}$$

predict a single fixed-length output vector given a variable length sequence as input

### Example:

Sentiment classification with word level **bidirectional** LSTM trained on a subset of the Internet Movie Database (IMDB) reviews. (20k positive and 20k negative examples)

Prediction examples for two inputs: 'this movie is so great' ⇒ 'positive' 'this movie is so bad' ⇒ 'negative'



• Seq2Seq (sequence translation)

 $f_{ heta}: \mathbb{R}^{TD} o \mathbb{R}^{T'C}$ 

- aligned: T = T'
- unaligned:  $T \neq T'$

• Seq2Seq (sequence translation)

$$f_{\theta}: \mathbb{R}^{TD} \to \mathbb{R}^{TC}$$

**Bi-directional** 

• aligned: T = T'

modify the RNN as:

$$p\left(y_{1:T} \mid x_{1:T}
ight) = \sum_{h_{1:T}} \prod_{t=1}^{T} p\left(y_t \mid h_t
ight) \mathbb{I}\left(h_t = f\left(h_{t-1}, x_t
ight)
ight) \ ext{initial state: } h_1 = f\left(h_0, x_1
ight) = f_0\left(x_1
ight)$$

dense sequence labeling: predict one label per location





Seq2Seq (sequence translation)

$$f_{\theta}: \mathbb{R}^{TD} \to \mathbb{R}^{TC}$$

• aligned: T = T'

modify the RNN as:

$$p\left(y_{1:T} \mid x_{1:T}
ight) = \sum\limits_{h_{1:T}} \prod\limits_{t=1}^{T} p\left(y_t \mid h_t
ight) \mathbb{I}\left(h_t = f\left(h_{t-1}, x_t
ight)
ight)$$

more depth to be more

expressive

input-to-hidden weights hidden-to-hidden weights  $h_t^l = arphi_l \left( W_{xh}^l h_t^{l-1} + W_{hh}^l h_{t-1}^l 
ight)$ 

$$y_t = W_{hy} h_t^L$$



Seq2Seq (sequence translation)

• unaligned:  $T \neq T'$ 

$$f_{\theta}: \mathbb{R}^{TD} \to \mathbb{R}^{T'C}$$

Example: translating English to French

- encode the input sequence to get the context vector, the last state of an RNN,  $c = f_e(x_{1:T})$
- generate the output sequence using an RNN decoder,  $y_{1:T'} = f_d(c)$



source input words

see code here

### Training: Backpropagation through time (BPTT)

unroll the computation graph, then apply the backpropagation

#### Example:

$$egin{aligned} \overline{P} & \overline{P} &$$

 $\begin{array}{c|c} \frac{\partial L}{\partial W_{hx}} \\ \frac{\partial L}{\partial W_{hh}} \\ \frac{\partial L}{\partial W_{hy}} \end{array}$ 



### Training: Backpropagation through time (BPTT)

unroll the computation graph, then apply the backpropagation  $[\operatorname{vec}(W_{hx}); \operatorname{vec}(W_{hh})]$  $egin{aligned} \overline{ extsf{poly}} & \overline{ extsf{h}}_{t} = W_{hx}x_t + W_{hh}h_{t-1} = f\left(x_t, h_{t-1}, w_h
ight) \ \hat{y}_t = W_{hy}h_t = g(h_t, w_y) \end{aligned}$ Example:  $\overset{\text{so}}{\stackrel{\text{so}}{\quad}} L = \frac{1}{T} \sum_{t=1}^{T} \ell\left(y_t, \hat{y}_t\right)$  $rac{\partial f(x_t,h_{t-1},w_h)}{\partial w_h}+rac{\partial f(x_t,h_{t-1},w_h)}{\partial h_{t-1}}rac{\partial h_{t-1}}{\partial w_h}$  $\left\{\begin{array}{c} \frac{\partial L}{\partial W_{hx}}\\ \frac{\partial L}{\partial W_{hh}}\\ \frac{\partial L}{\partial W_{hh}}\end{array}\right\} \frac{\partial L}{\partial w_{h}} = \frac{1}{T} \sum_{t=1}^{T} \frac{\partial \ell(y_{t}, \hat{y}_{t})}{\partial w_{h}} = \frac{1}{T} \sum_{t=1}^{T} \frac{\partial \ell(y_{t}, \hat{y}_{t})}{\partial \hat{y}_{t}} \frac{\partial g(h_{t}, w_{y})}{\partial h_{t}} \frac{\partial h_{t}}{\partial w_{h}} \\ \frac{\partial L}{\partial W_{t}} \end{array}$ expand this expand this recursively  $rac{\partial h_t}{\partial w_h} = rac{\partial f(x_t,h_{t-1},w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{\substack{j=i+1}}^t rac{\partial f(x_j,h_{j-1},w_h)}{\partial h_{j-1}}
ight) rac{\partial f(x_i,h_{i-1},w_h)}{\partial w_h}$ see code here

# Gating and long term memory

### Vanishing and exploding gradients

activations can decay or explode as we go forwards and backwards in time

RNN variations that circumvent this:

- Gated recurrent units (GRU)
  - learns when to update the hidden state, by using a gating unit
- Long short term memory (LSTM)
  - augments the hidden state with a memory cell

# The Gated Recurrent Unit (GRU):



$$egin{aligned} & z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \ & r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \ & \hat{h}_t = \phi(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h) \ & h_t = (1-z_t) \odot h_{t-1} + z_t \odot \hat{h}_t \end{aligned}$$

# Summary

- Recurrent neural networks (RNNs)
  - Vec2Seq (sequence generation)
  - Seq2Vec (sequence classification)
  - Seq2Seq (sequence translation)
  - training with back propagation through time