

# Applied Deep Learning for NLP

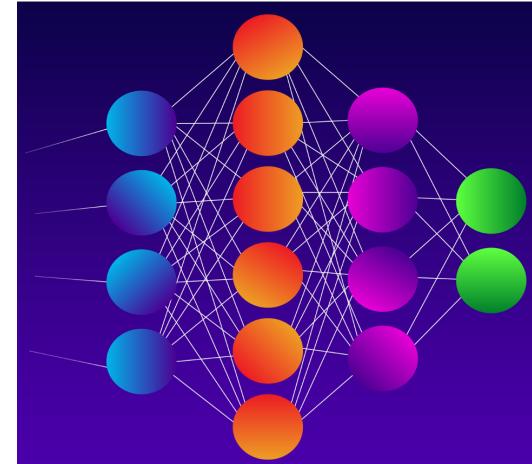
## Week 4 - RNNs and LSTMs

Juan Carlos Medina Serrano

Technische Universität München  
Hochschule für Politik  
Political Data Science

Munich, 12. November 2020

political  
data  
science



## Task: Text Classification

We have seen examples of text classifications in the coding sessions. It is the most widely used NLP task in the industry.

**Binary** classification (sigmoid) and **Multiclass** classification (softmax)

*Examples:* spam detection, misinformation detection, authorship attribution, sentiment analysis, email categorization (personal, social,...), news categorization (sports, politics, entertainment,...)

Loss: Cross-entropy loss (most of the time, sometimes others like Hinge loss)

$$L_{binary} = -(y * \log(p) + (1 - y) * \log(1 - p))$$

$$L_{multiclass} = - \sum_{c=1}^M y_c \log(p_c)$$

Key to understand cross-entropy:

- ▶  $\log(1) = 0$ ,  $\log(0.5) = -0.3$ ,  $\log(0.2) = -0.7$ .
- ▶  $y$  is a one-hot vector, with only 1 for the correct class and zero for the rest

## Task: Text Classification

Positives: Class we care about detecting, for spam, spam is the positive class!

Negatives: The other class, normally the less interesting

True positives, true negatives: Model made the right decision

False positives: Model predict positive, but it was wrong.

False negatives: Model predict negative, but it was wrong.

### METRICS:

- ▶ **Accuracy** Fraction of times the model makes the correct predictions as compared to the total number of predictions N.

$$\text{Accuracy} = \frac{\text{correct}}{N}$$

- ▶ **Precision**. Shows how precise is the model. Given all the predicted positive cases, how many are really positive.

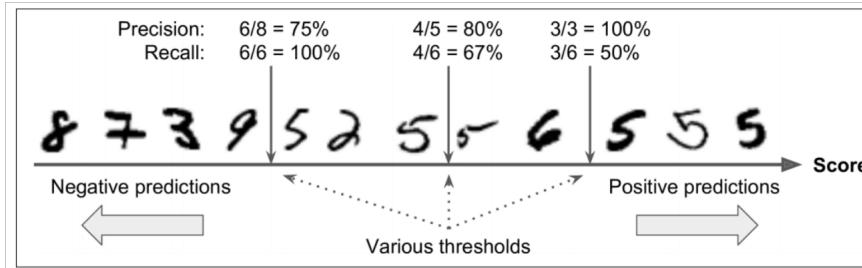
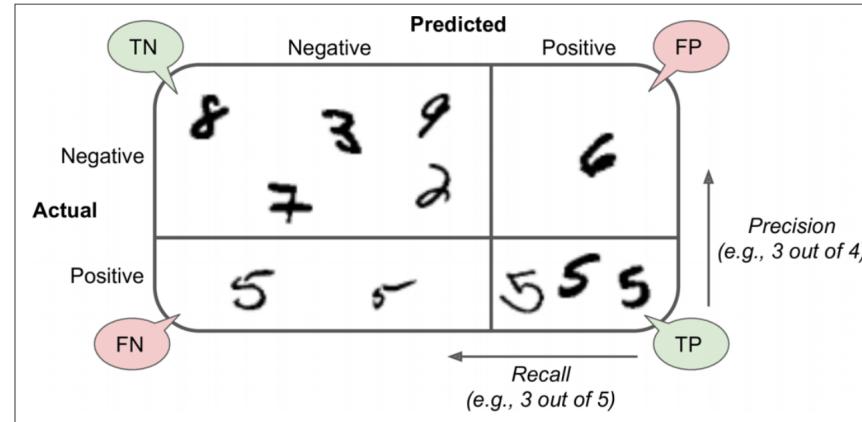
$$\text{Precision} = \frac{TP}{TP + FP}$$

- ▶ **Recall**. Shows how many of the positive cases your model can detect from all real positive cases.

$$\text{Recall} = \frac{TP}{TP + FN}$$

## Task: Text Classification

Confusion Matrix:



## Task: Language Modeling

Predict what word comes next given a text corpus.

Given a sequence of words  $x^{(1)}, x^{(2)}, \dots, x^{(t)}$  compute the probability of the next word  $x^{(t+1)}$

$$P(x^{(t+1)}|x^{(1)}, \dots, x^{(t)})$$

Can be also interpreted as assigning a probability to a piece of text

$$P(x^{(1)}, \dots, x^{(t)}) = P(x^{(1)}) * P(x^{(2)}|x^{(1)}) * \dots * P(x^{(t)}|x^{(1)}, \dots, x^{(t-1)})$$

## Task: Language Modeling

Before neural networks counting n-grams was the main technique. The probability of word  $w$ :

$$P(w | "Harry \ opened \ the") = \frac{\text{count}("Harry \ opened \ the \ w")}{\text{count}("Harry \ opened \ the")}$$

Disadvantage: Need to store all n-grams. What about never seen n-grams?

How to evaluate if a good language model is good?

**PERPLEXITY** Inverse probability of the corpus of size  $T$

$$\prod_{t=1}^T \left( \frac{1}{P(x^{(t+1)} | x^{(1)}, \dots, x^{(t)})} \right)^{\frac{1}{T}}$$

Lower perplexity is better!

## Task: Text Generation

Language models are the basic blocs for other NLP tasks.

In text generation, we use the probabilities calculated from the language model to generate new text.

Next word prediction consists of taking a **sample** from the probability distribution (Like when you are texting on the phone, the word suggestions are the ones with the highest probabilities)

## Language model with FNN

A fixed-window language model

output distribution

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

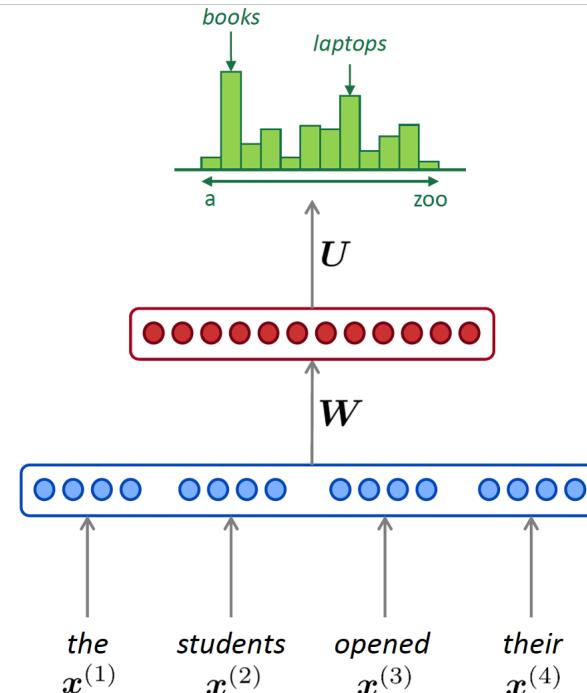
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$

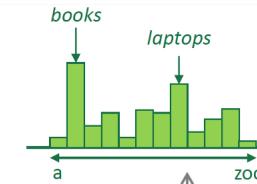


## Recurrent Neural Networks (RNNs)

Apply the same weights matrices repeatedly

**output distribution**

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}h^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$



**hidden states**

$$h^{(t)} = \sigma(\mathbf{W}_h h^{(t-1)} + \mathbf{W}_e e^{(t)} + \mathbf{b}_1)$$

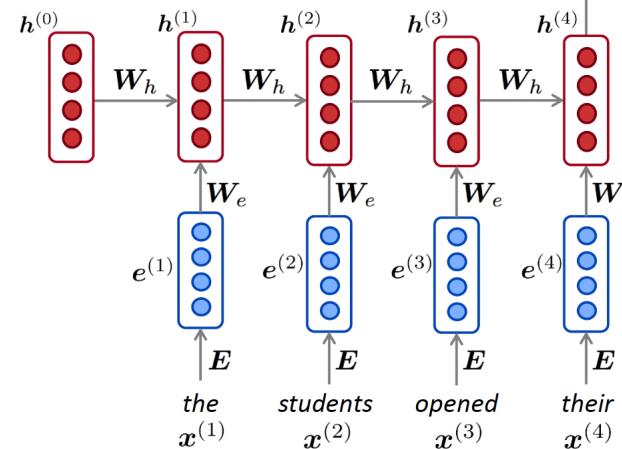
$h^{(0)}$  is the initial hidden state

**word embeddings**

$$e^{(t)} = \mathbf{E}x^{(t)}$$

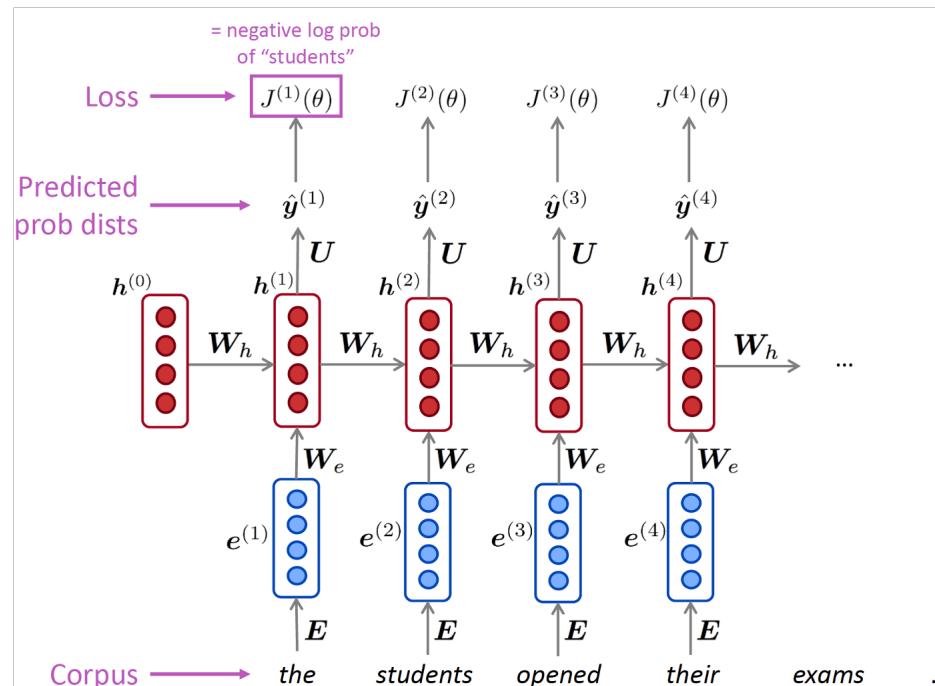
words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$



## Training an RNN

$$Loss = L(\theta) = \frac{1}{T} \sum_{t=1}^T L^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{w=1}^W -y_w^{(t)} \log(\hat{y}_w^{(t)}) = J(\theta)$$



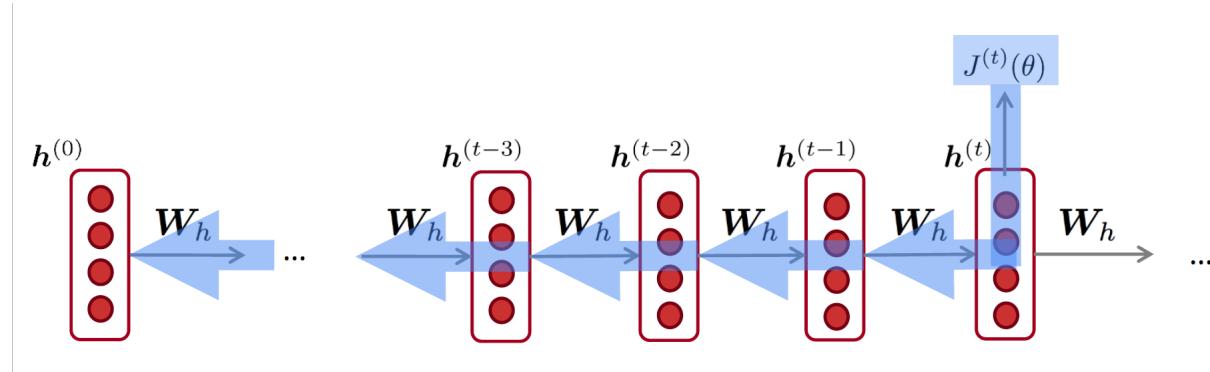
## Training an RNN

In practice, we don't calculate the loss of the whole corpus!

We use batches of sentences, compute the loss and update weights

## Training an RNN

**Backpropagation** needs to happen also through time and through the network



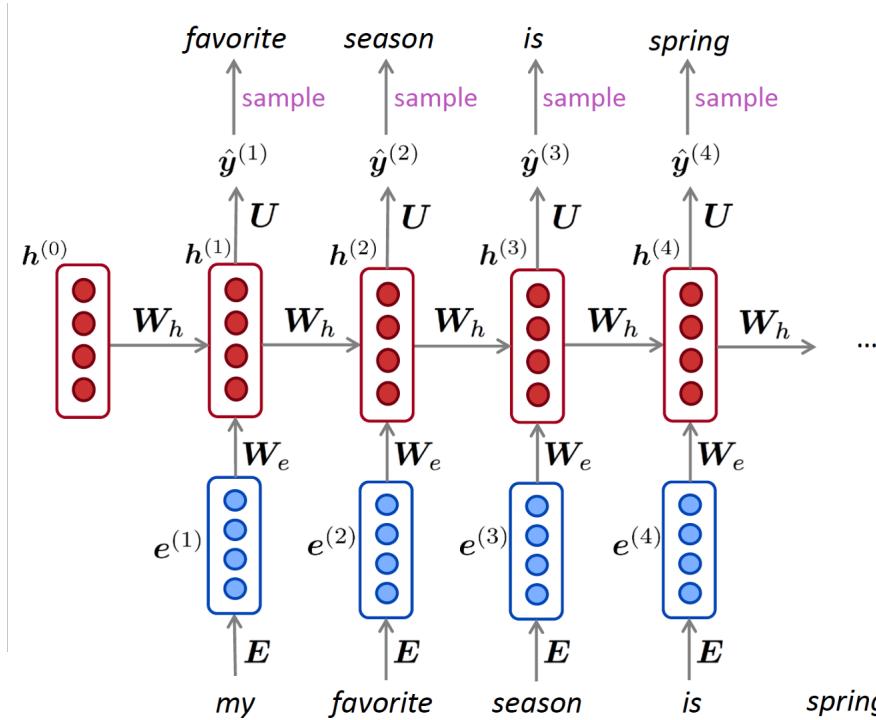
Problems: 1) exploding gradients, 2) vanishing gradients, 3) information is lost after many time steps.

Against exploding gradients, apply **gradient clipping**.

For vanishing gradients? Batch normalization and similar approaches are not as effective

## Text generation with RNN

Sampled output is next step's input:



## LSTM

Long-short term memory (LSTM) is a cell that improves the simple RNN hidden state layer. We now have a **cell state** to remember information from the past and **gates** to decide what is more important, present or past.

Old information can be kept! Solves also the vanishing gradient problem

## LSTM

**Sigmoid function:** all gate values are between 0 and 1

$$\begin{aligned}\mathbf{f}^{(t)} &= \sigma(\mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f) \\ \mathbf{i}^{(t)} &= \sigma(\mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i) \\ \mathbf{o}^{(t)} &= \sigma(\mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o)\end{aligned}$$

All these are vectors of same length  $n$

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c)$$

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \circ \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \circ \tilde{\mathbf{c}}^{(t)}$$

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \circ \tanh \mathbf{c}^{(t)}$$

Gates are applied using element-wise product

## LSTM

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Sigmoid function:** all gate values are between 0 and 1

$$\begin{aligned} f^{(t)} &= \sigma \left( \mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f \right) \\ i^{(t)} &= \sigma \left( \mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i \right) \\ o^{(t)} &= \sigma \left( \mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o \right) \end{aligned}$$

All these are vectors of same length  $n$

$$\tilde{\mathbf{c}}^{(t)} = \tanh \left( \mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c \right)$$

$$\mathbf{c}^{(t)} = f^{(t)} \circ \mathbf{c}^{(t-1)} + i^{(t)} \circ \tilde{\mathbf{c}}^{(t)}$$

$$\mathbf{h}^{(t)} = o^{(t)} \circ \tanh \mathbf{c}^{(t)}$$

Gates are applied using element-wise product

## LSTM

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Sigmoid function:** all gate values are between 0 and 1

$$\begin{aligned} f^{(t)} &= \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f) \\ i^{(t)} &= \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i) \\ o^{(t)} &= \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o) \end{aligned}$$

$$\tilde{c}^{(t)} = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c)$$

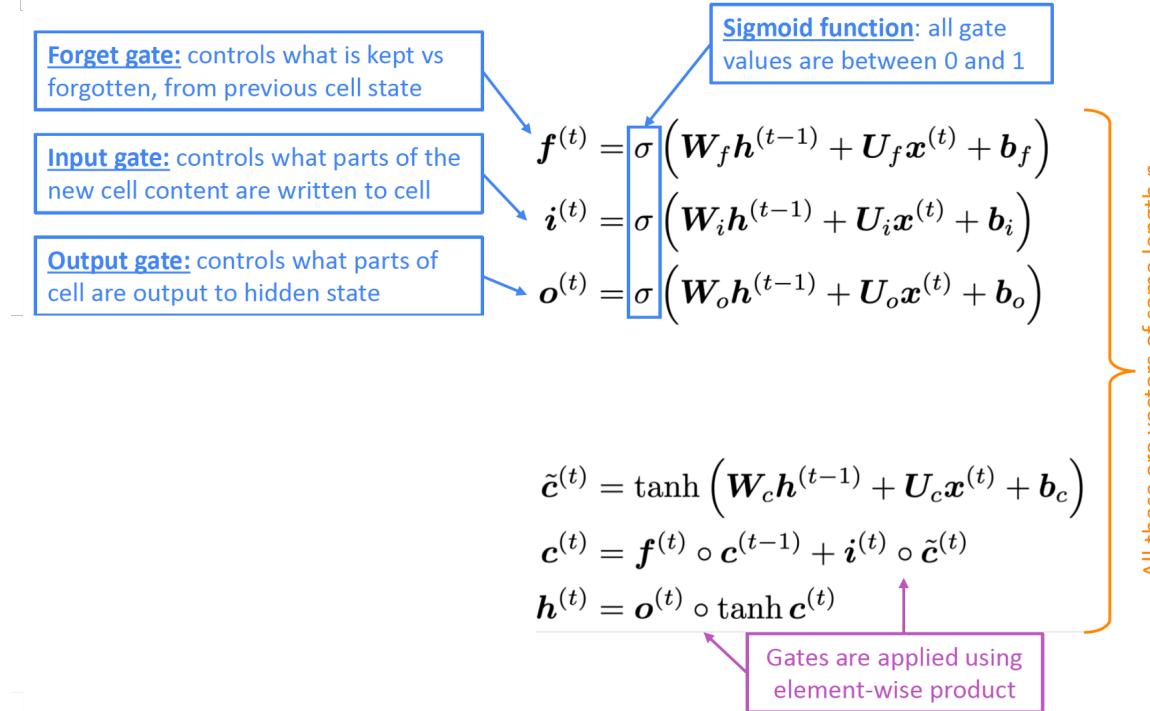
$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

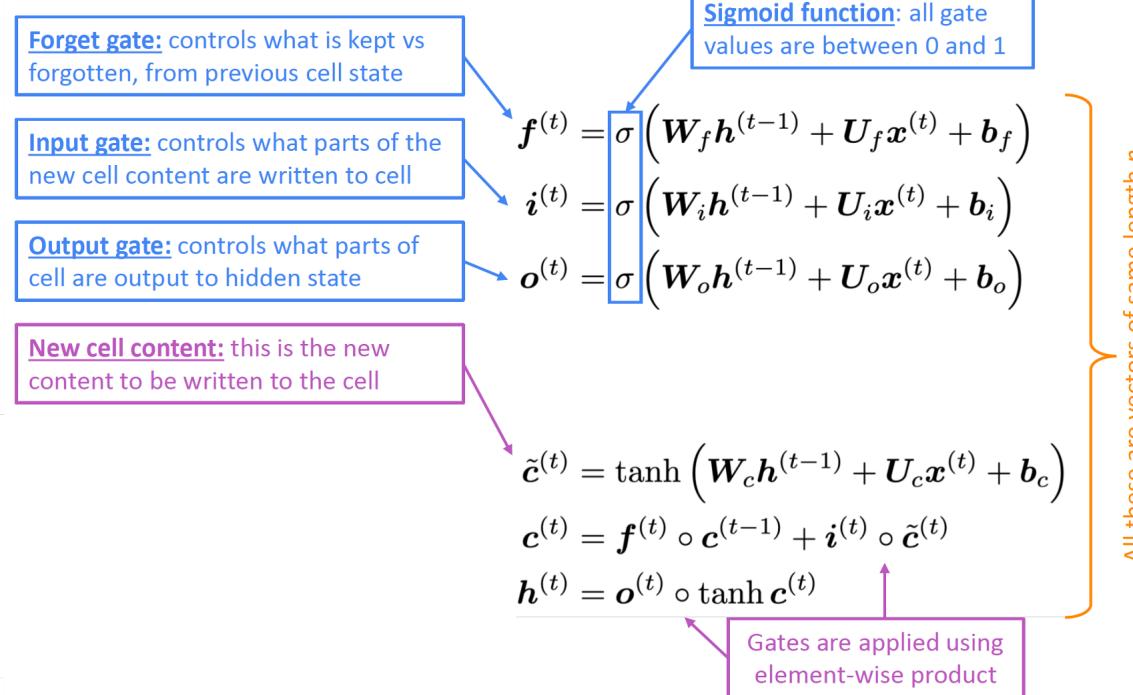
Gates are applied using element-wise product

All these are vectors of same length  $n$

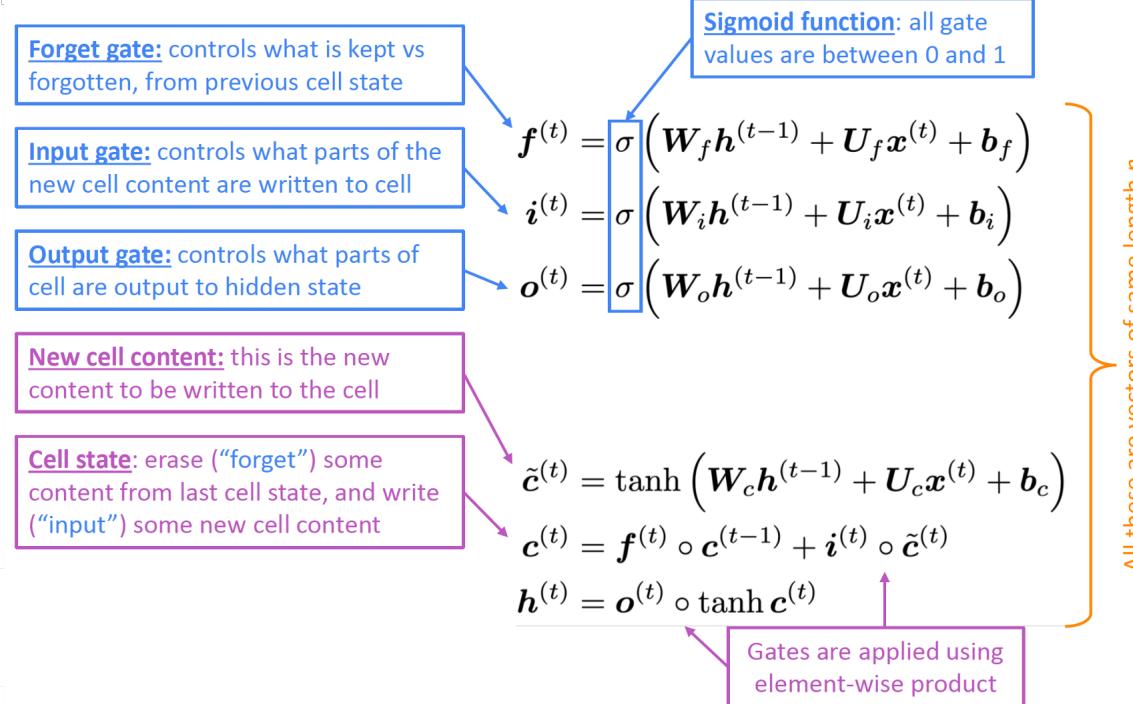
## LSTM



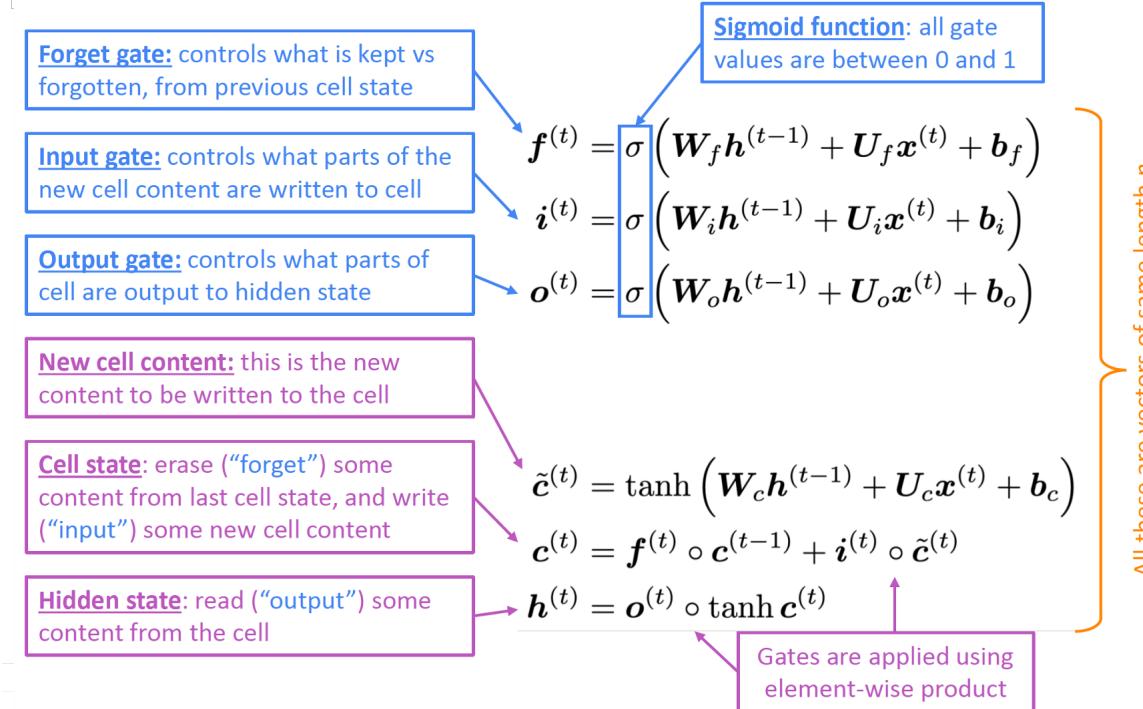
## LSTM



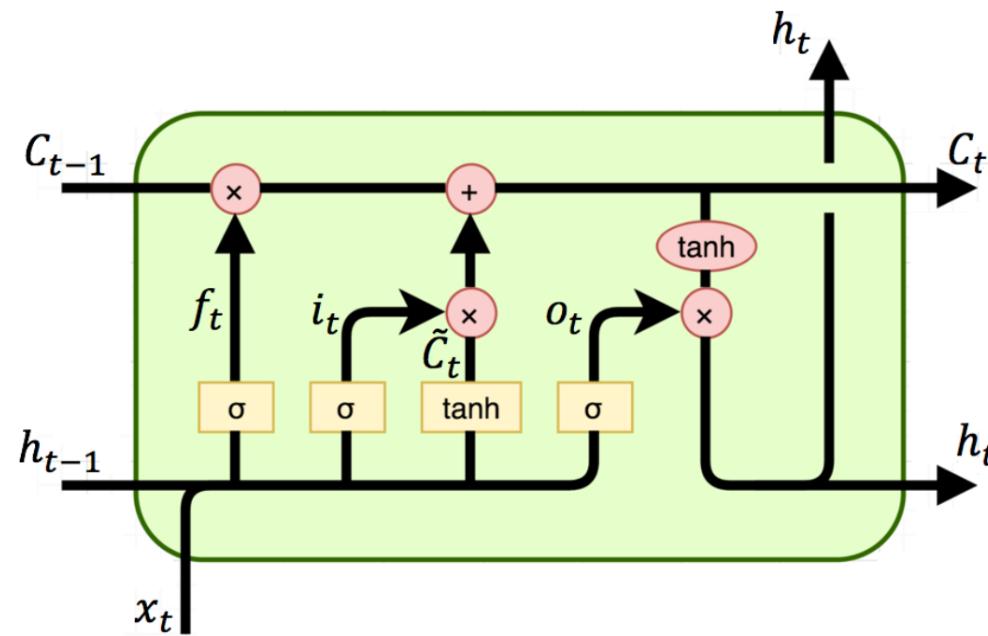
## LSTM



## LSTM



## LSTM



## RNN Cell

output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}h^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma(\mathbf{W}_h h^{(t-1)} + \mathbf{W}_e e^{(t)} + \mathbf{b}_1)$$

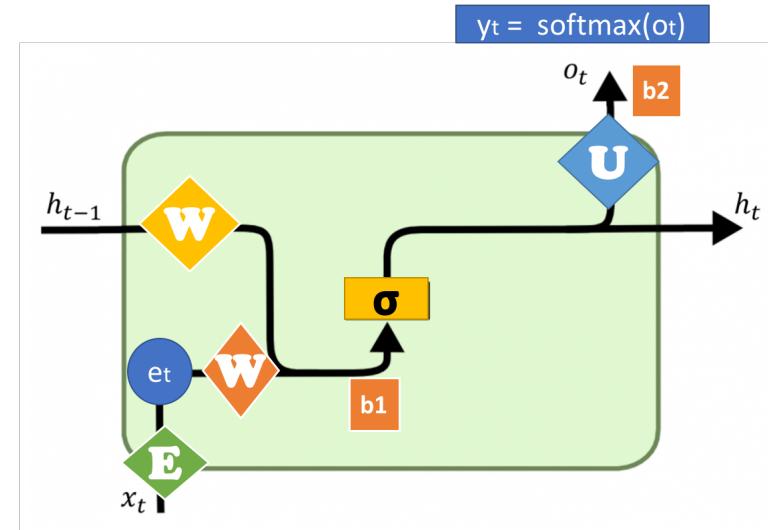
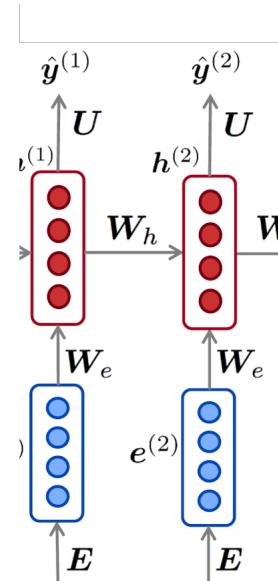
$h^{(0)}$  is the initial hidden state

word embeddings

$$e^{(t)} = \mathbf{E}x^{(t)}$$

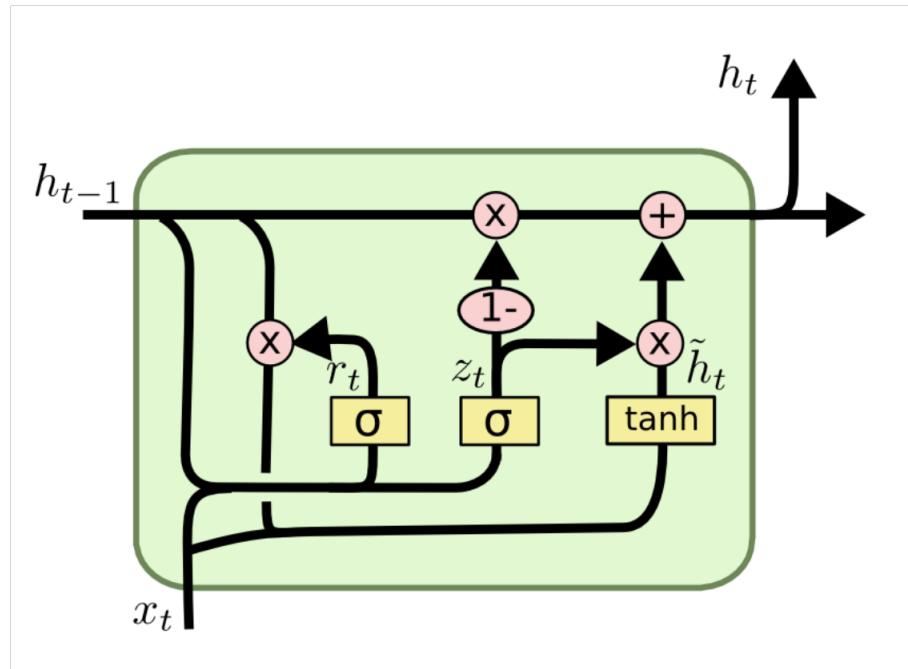
words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$



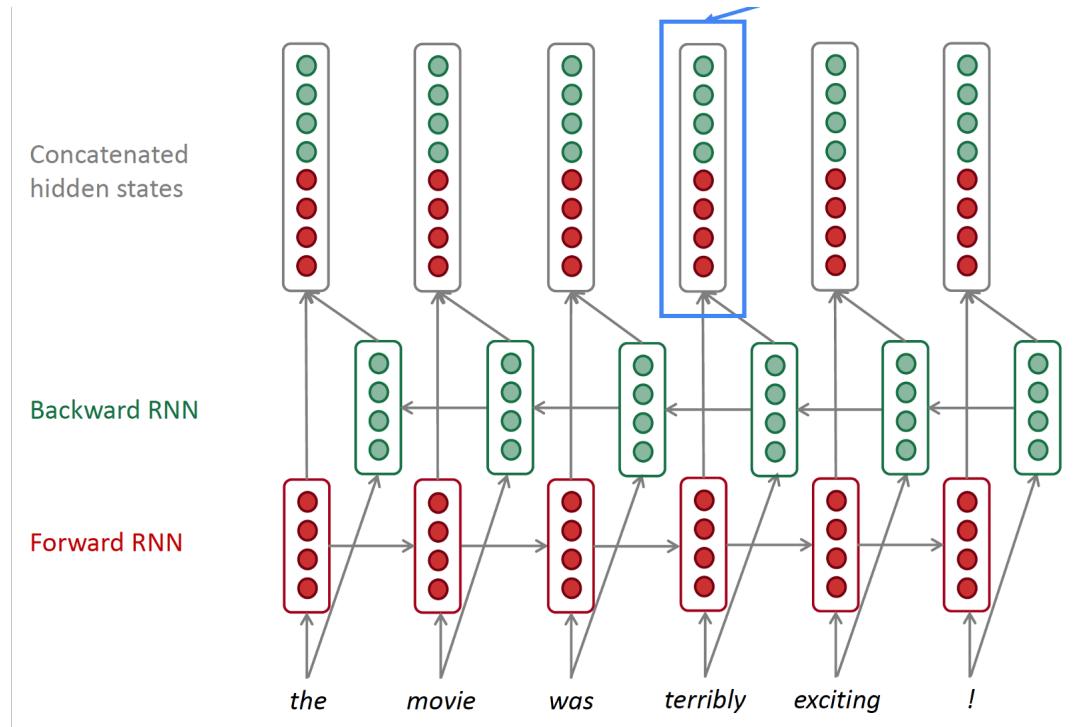
## GRU

LSTMs are complicated, with many parameters more than RNNs. Gated Recurrent Unit (GRUs) simplifies the LSTM cell, without taking out its benefits. No more cell state, the hidden state keeps the information from the past



## Bidirectional RNNs

Get information from past and future. Only for tasks where the complete sentence is given. (Not to generate text)



## Multi-layer RNNs

RNNs can also have more hidden layers! (Normally 2 to 4 are used, otherwise too complex and hard to train)

