

# Text Classification and Neural Networks

Jelke Bloem

Text Mining  
Amsterdam University College

25-03-2025

# Announcements

- Assignment 2 due today
- Transformers reading assignment on Friday
- Assignment 3 coming up

# Overview

- 1 Text Classification
- 2 Logistic Regression
- 3 Neural Networks
- 4 Neural Language Models
- 5 Recurrent Neural Networks
- 6 Towards Transformer Models
- 7 Extras

## **Text Classification**

# Task definition

- We are given a **training set**  $\{X, Y\}$  of data pairs  $(x, y)$ , where  $x$  is a text document and  $y$  is the class the document belongs to.
- Each  $y \in \mathcal{Y}$ , where  $\mathcal{Y} = \{c_1, c_2, \dots, c_k\}$  are the distinct (finite and enumerable) classes we have. If  $|\mathcal{Y}| = k = 2$ , we have a binary classification task.
- Using a *learning method*, our goal is to learn a **classifier**, or a classification function  $\gamma$  that maps documents to classes:

$$\gamma : \mathcal{X} \rightarrow \mathcal{Y}$$

- The fact that we use annotated data to learn makes this a form of *supervised learning*. Note that a “document” can be anything really: words, text sequences, longer texts.

# Examples

task	$x$	$y$
language ID	text	{english, mandarin, greek, ...}
spam classification	email	{spam, not spam}
authorship attribution	text	{jk rowling, james joyce, ...}
genre classification	novel	{detective, romance, gothic, ...}
sentiment analysis	text	{postive, negative, neutral, mixed}

*Credit: David Bamman (UC Berkeley).*

# Text representation

Our text documents  $X$  can be **represented** in many ways:

- Pre-computed features (e.g., the length of the document or the average length of the words it contains).
- A selection of words (e.g., only stopwords for language detection).
- Words in isolation (so called “bag of words”, or unigram model).
- Conjunctions of words (e.g., bigrams).
- Higher-order features (e.g., PoS).
- Word embeddings.

## Example: Authorship Attribution

- Texts of unknown origin
- A group of potential authors
- Known texts by those authors



## Example: Authorship Attribution

A classification problem where each potential author is a class label, and the known texts are labeled training data

## Example: Authorship Attribution

A classification problem where each potential author is a class label, and the known texts are labeled training data

- Stylometry
- Forensic linguistics
- Historical linguistics
- Translation studies

# Example: Authorship Attribution

## Known cases

- The Federalist Papers
  - ▶ Historical essays in support of American constitution by 3 authors
  - ▶ But who wrote which of the 85 essays?
- Unabomber Manifesto
- Bot detection
- Authorship of pseudonymously published book
- Authorship of historical scientific texts

## Example: Authorship Attribution

Features for classification

# Example: Authorship Attribution

## Features for classification

- Word features (e.g. linking words)
- N-grams
- Punctuation counts
- Stylometry:
  - ▶ Type/token ratio
  - ▶ Average sentence/word length
  - ▶ Number of words in paragraph
  - ▶ Number of hapax legomena
- ...

# Example: Authorship Attribution

## Classification models

- K-nearest neighbour
- Random forest classifier
- Logistic regression
- ...

## Logistic Regression

# Logistic regression

- Our goal is, given a document represented with a feature vector  $\mathbf{x}$  and classes  $c \in \mathcal{Y}$ , to learn a classifier discriminating the right class for  $\mathbf{x}$ :

$$\hat{p}(y = c|\mathbf{x})$$

- Let us start with a binary classifier and two classes, thus  $\mathcal{Y} = \{0, 1\}$ .
- We need to estimate  $\hat{p}(y = 1|\mathbf{x})$ , and  $\hat{p}(y = 0|\mathbf{x}) = 1 - \hat{p}(y = 1|\mathbf{x})$  will follow suit.
- Logistic regression uses two components for this: a **linear model** of the inputs and the **Sigmoid (or logistic) function**. So, it is like the perceptron but with a different classification function.



## Sigmoid (or logistic) function

- Let us consider the set of features  $x_1, x_2, \dots, x_d$  we used to represent our input document  $\mathbf{x}$ . We add  $x_0 = 1$  to model the intercept, and create a linear model with them:

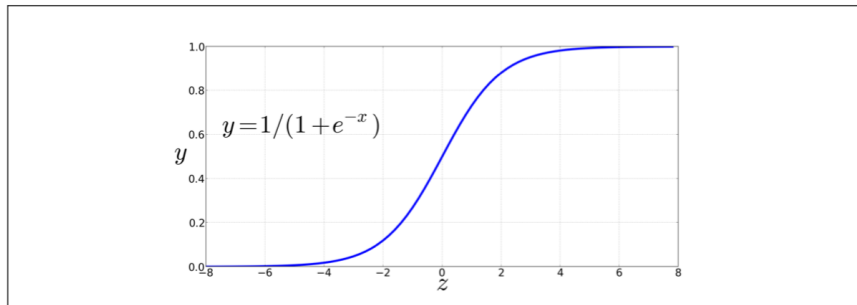
$$z = \sum_{j=0}^d w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

- To create a probability distribution, we pass  $z$  through the Sigmoid  $\sigma(z)$ :

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

- The Sigmoid squeezes  $z$  within 0 and 1 and is always positive.

# Sigmoid (or logistic) function



**Figure 5.1** The sigmoid function  $y = \frac{1}{1+e^{-x}}$  takes a real value and maps it to the range  $[0, 1]$ . Because it is nearly linear around 0 but has a sharp slope toward the ends, it tends to squash outlier values toward 0 or 1.

*Credit: M&J, Ch. 5.*

# Practicalities: Data splitting

	training	development	testing
size	80%	10%	10%
purpose	training models	model selection; hyperparameter tuning	evaluation; never look at it until the very end

*Credit: David Bamman (UC Berkeley).*

# Accuracy and baselines

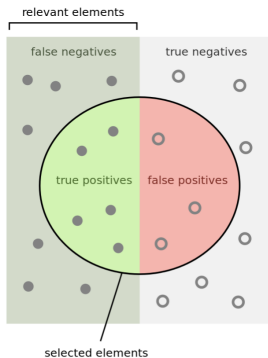
- **Accuracy** is the fraction of correctly predicted data points over the total. It can be calculated on any dataset split: train, development and test. Very good starting point.
- **Baseline**: important to have one. It can be a random classifier (i.e., flip a coin for a binary classifier), or a fast and reasonable model (e.g., logistic regression with TF-IDF features).

# Precision and recall

Given a binary classifier:

- **True positive:** a data point correctly predicted to be 1.
- **True negative:** a data point correctly predicted to be 0.
- **False positive:** a data point incorrectly predicted to be 1.
- **False negative:** a data point incorrectly predicted to be 0.

# Precision and recall



How many selected  
items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant  
items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

*Credit: Wikipedia.*

## F-measure and accuracy reloaded

- F-measure (harmonic mean of precision and recall):

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

- Accuracy:

$$A = \frac{tp + tn}{tp + tn + fp + fn}$$

# Parameters and hyperparameters

Parameters whose values are *learned*

Feature	$\beta$
the	0.01
and	0.03
bravest	1.4
love	3.1
loved	1.2
genius	0.5
<i>BIAS</i>	-0.1

Hyperparameters whose values are *chosen*

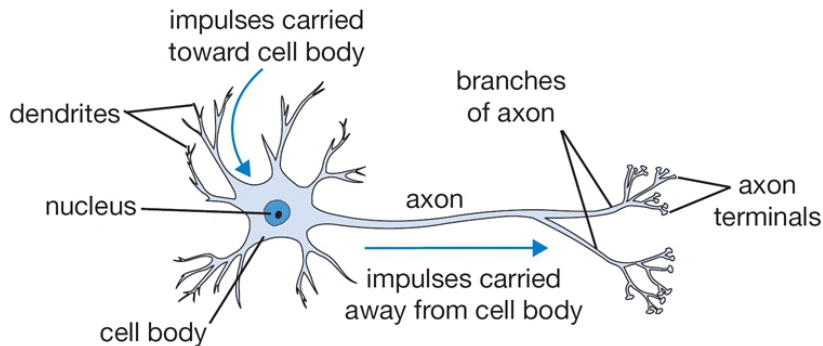
Hyperparameter	value
minimum word frequency	5
max vocab size	10000
lowercase	TRUE
regularization strength	1.0

*Credit: David Bamman (UC Berkeley).*



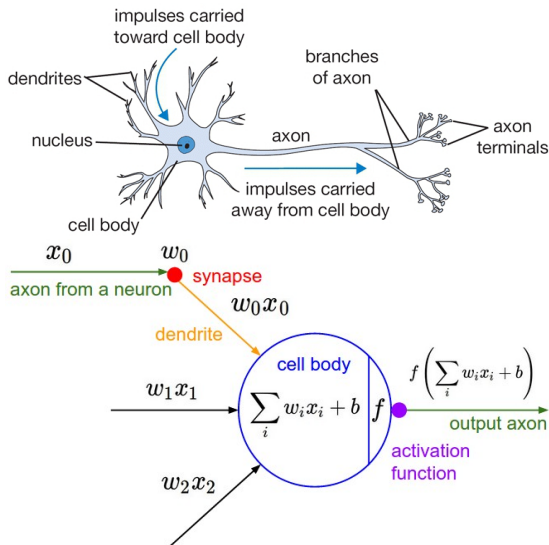
# Neural Networks

# A single neuron



*Credit: Andrej Karpathy via Stanford's CS231N.*

# A single neuron



*Credit: Andrej Karpathy via Stanford's CS231N.*

# Logistic regression as a neural network

Following the notation in the previous slide, we have:

- $\mathbf{x} = \langle x_0, x_1, x_2, \dots, x_d \rangle$  is our input representation.
- We aggregate the features  $\mathbf{x}$  into a linear combination using weights  $\mathbf{w}$ . We also include the bias term  $b$  into the matrix by adding an appropriate dimension fixed at 1 to  $\mathbf{x}$ , so that we can use matrix notation:

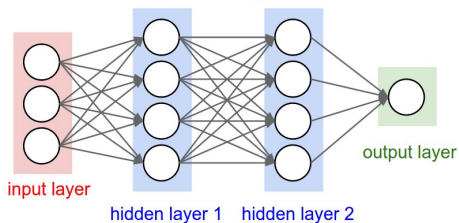
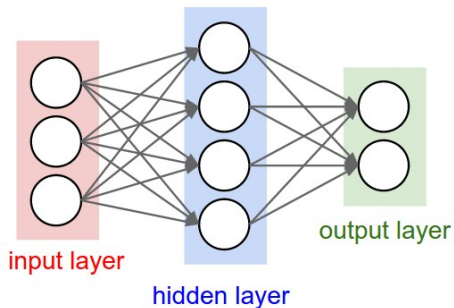
$$z = \sum_{i=0}^d w_i x_i = \mathbf{w} \cdot \mathbf{x}$$

- We pass  $z$  through the sigmoid function to map it to range  $[0,1]$ :

$$f = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- **Linear models are a single neuron.** *Question: what is the activation function for linear regression?*

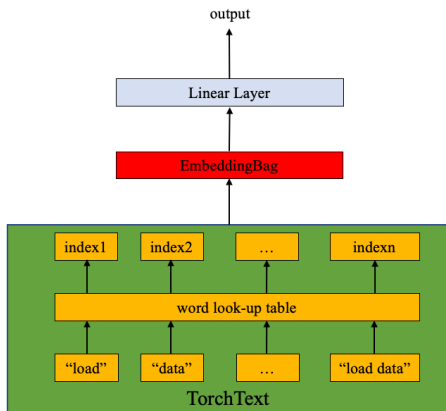
## From single layer to multi-layer



*Credit: Andrej Karpathy via Stanford's CS231N.*

## Using embeddings as features

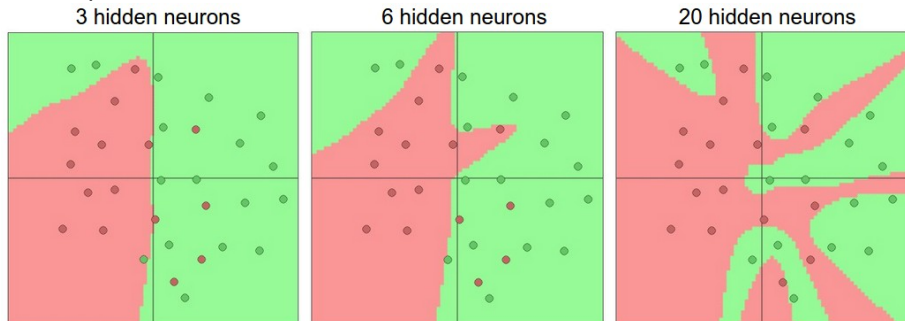
Neural networks are **modular**: we can piece them together into advanced architectures. For example, we can use embeddings to represent our input (either training them or using pre-trained ones). *More on this in the lab.*



*Credit: TorchText.*

# Why do we need non-linearities?

Multiple layers and **non-linear functions** (such as the sigmoid) allow us to fit complex decision boundaries.



*Credit: Andrej Karpathy via Stanford's CS231N.*

# How do we train neural networks?

- Key idea: use a smart way to apply SGD, called **backpropagation**.
- Backpropagation combines using the chain rule to calculate local derivatives (called gradients) with the re-use of pre-computed operations to speed the computation up.
- *More on this in the external materials for the course.*



# Neural networks practicalities

Training neural networks entails a lot more than stacking up layers. Several topics require practical and theoretical knowledge beyond this course:

- Weight initialization
- Regularization (e.g., via dropout)
- Which non-linearities to use
- Which loss functions to use
- How to monitor and adjust the learning process (e.g., optimizers and learning rates) to avoid dying neurons and overfitting

## Neural Language Models

# Recap

With language models, we want to compute:

- The probability of a sequence of words:  $P(w) = P(w_1, w_2, \dots, w_n)$ .
- The probability of a new word given what came before it:  
 $P(w_n | w_1, w_2, \dots, w_{n-1})$ .
- We have also seen n-gram language models which make use of the Markov assumption. For example, a trigram language model would predict the probability of a sequence of words by conditioning on the two previous words at each step:

$$P(w) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots P(w_n|w_{n-2}, w_{n-1})$$

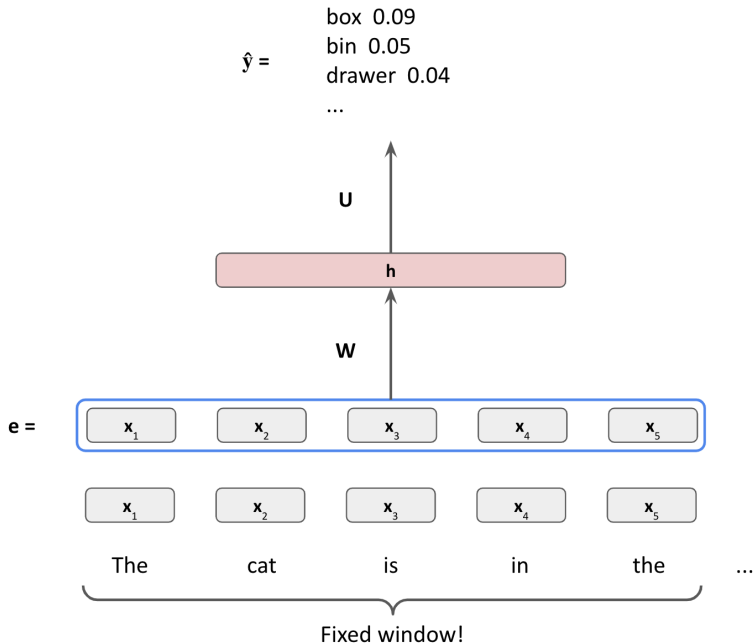
- n-gram language models have some issues, including sparsity and limited use of context.

# A first neural language model

We can start by using word embeddings as features for a **fixed-window neural language model**:

- Given a fixed-window sequence of words represented using their embeddings  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$ ;
- we are interested in estimating the probability of the next word  $\hat{p}(\mathbf{x}_{t+1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$ .
- We can use the previous words' embeddings as features, concatenate them and feed them to a hidden layer with a non-linearity, and concluding by estimating probabilities with a softmax.

# A first neural language model



# A first neural language model

Where:

- Each  $\mathbf{x}$  is a  $1 \times d$  embedding vector of dimensionality  $d$ .
- $\mathbf{e} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t]$  is the concatenation of the embeddings of the words in the fixed-window preceding  $w_{t+1}$ . Therefore  $\mathbf{e}$  has dimensionality  $1 \times dk$ , where  $k$  is the number of words in the fixed-window.
- $\mathbf{h} = f(\mathbf{W}\mathbf{e}^T)$  is the hidden layer, with  $f$  an appropriate activation function and  $\mathbf{W}$  a weight matrix.  $\mathbf{W}$  has dimensionality  $dk \times h$  and  $\mathbf{h}$  has dimensionality  $1 \times h$ . The dimensionality of the embeddings and the hidden layer are hyperparameters of the model.
- Finally,  $\hat{y} = \text{softmax}(\mathbf{U}\mathbf{h}^T)$  is the predicted next word probability as estimated using a softmax function. Here  $\mathbf{U}$  is another weight matrix of dimensionality  $|V| \times h$ , so that with the softmax we predict a probability distribution over the vocabulary  $V$ .

# A first neural language model

This model is interesting, yet unfortunately:

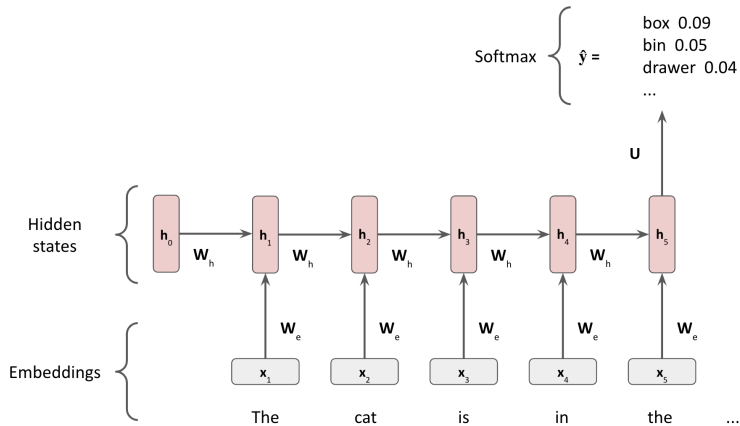
- it does not solve the use of a broader context with flexible word windows;
- larger windows would mean more parameters so scale is also an issue;
- and it does not make an efficient use of parameters and shared weights.

Can we do better? Yes, with **Recurrent Neural Networks**.

# Recurrent Neural Networks



# Recurrent architecture



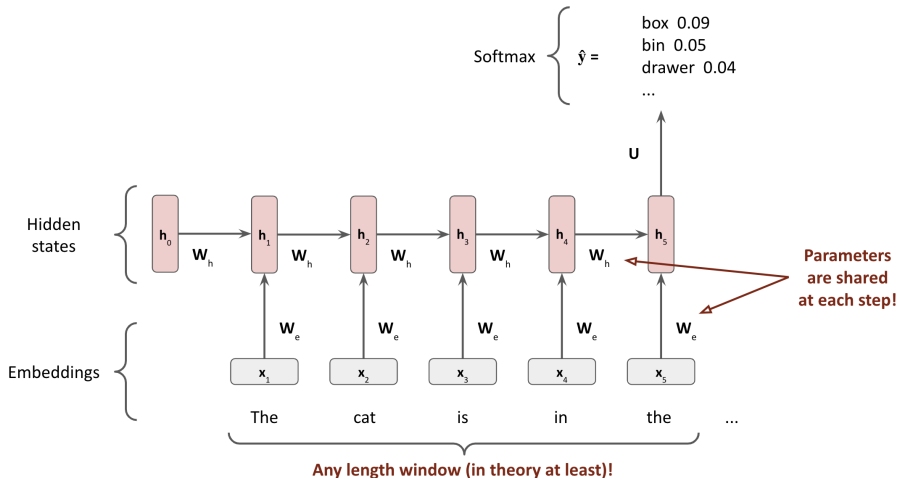
# Recurrent architecture

Where:

- Each  $\mathbf{x}$  is a  $1 \times d$  embedding vector of dimensionality  $d$ .
- $\mathbf{h}_t = f(\mathbf{W}_h \mathbf{h}_{t-1}^T + \mathbf{W}_e \mathbf{x}_t^T)$  is the hidden layer, with  $f$  an appropriate activation function,  $\mathbf{W}_e$  and  $\mathbf{W}_h$  weight matrices.  $\mathbf{W}_h$  has dimensionality  $h \times h$ ,  $\mathbf{W}_e$  has dimensionality  $h \times d$ , and  $\mathbf{h}$  has dimensionality  $1 \times h$ .  $\mathbf{h}_0$  is the initial hidden layer. The dimensionality of the embeddings and the hidden layer are hyperparameters of the model.
- Finally,  $\hat{y} = \text{softmax}(\mathbf{U} \mathbf{h}^T)$  is the predicted next word probability as estimated using a softmax function. Here  $\mathbf{U}$  is another weight matrix of dimensionality  $|V| \times h$ , so that with the softmax we predict a probability distribution over the vocabulary  $V$ .

(Notebook 8.1 Part II: Model)

# Recurrent architecture



# Why RNNs?

RNNs are a big step forward re. our previous concerns:

- Can process inputs on any length and use previous context of any length (in theory);
- model size does not depend on window size ( $W$  matrices remain of the same dimension);
- weights are shared across time steps.

Nevertheless, RNNs can be slow and won't really remember information from many steps back.

# Training RNNs

How can we train RNNs?

- Predict the next word at each step, and calculate the loss accordingly.
- The loss at step  $t$  is the usual **cross-entropy** (now on multiple classes), calculated for the word to be predicted at  $t + 1$ :

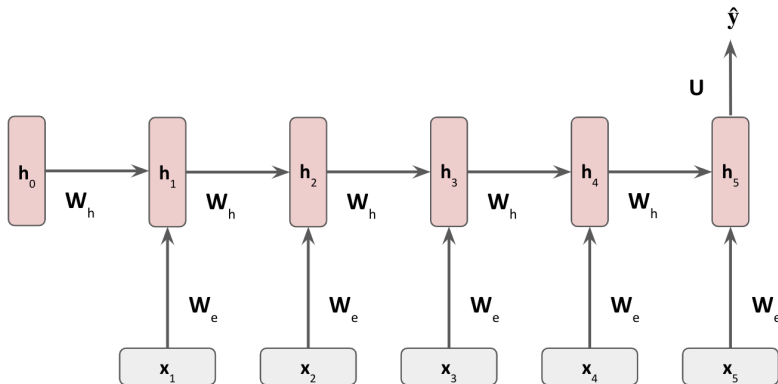
$$\mathcal{L}_t(\mathbf{W}) = - \sum_{w \in V} y_{w_{t+1}} \log(\hat{y}_{w_{t+1}}) = -\log(\hat{y}_{w_{t+1}})$$

- The overall loss is the average of the sum of the losses, calculated at each step:

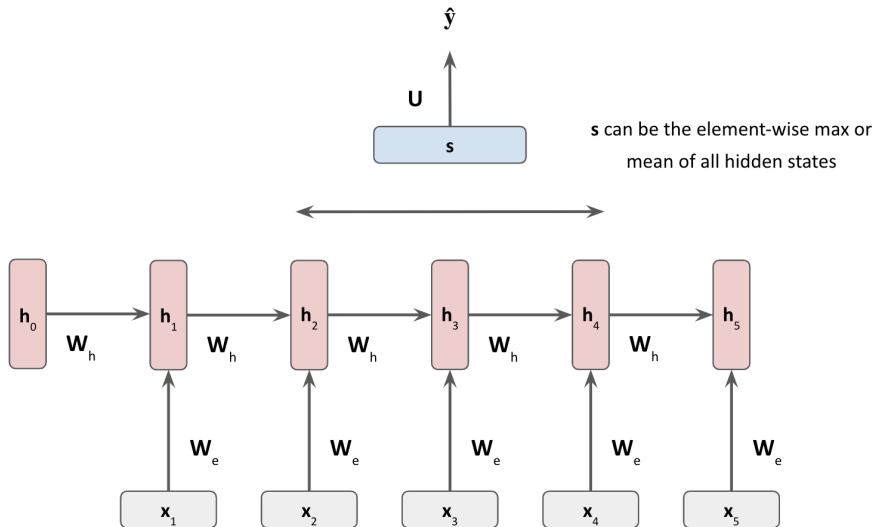
$$\mathcal{L}(\mathbf{W}) = \frac{1}{T} \sum_{t=1}^T \mathcal{L}_t(\mathbf{W}) = \frac{1}{T} \sum_{t=1}^T -\log(\hat{y}_{w_{t+1}})$$

- Optimization can be done via SGD (backpropagation). Since the parameters  $\mathbf{W}$  are used repeatedly, this is sometimes called **backpropagation through time**.

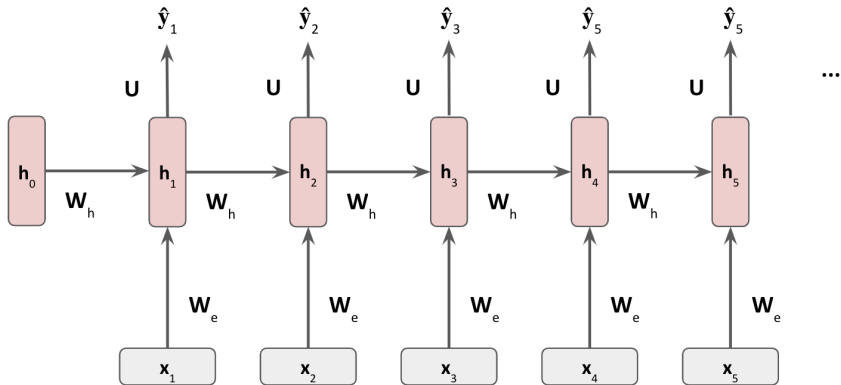
## RNN flavors: Many to one



## RNN flavors: Many to one



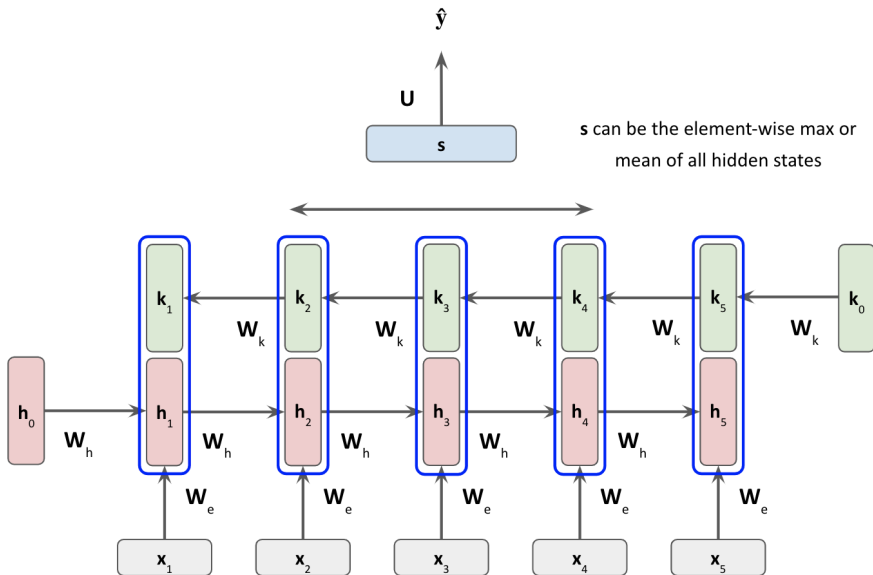
## RNN flavors: Many to many



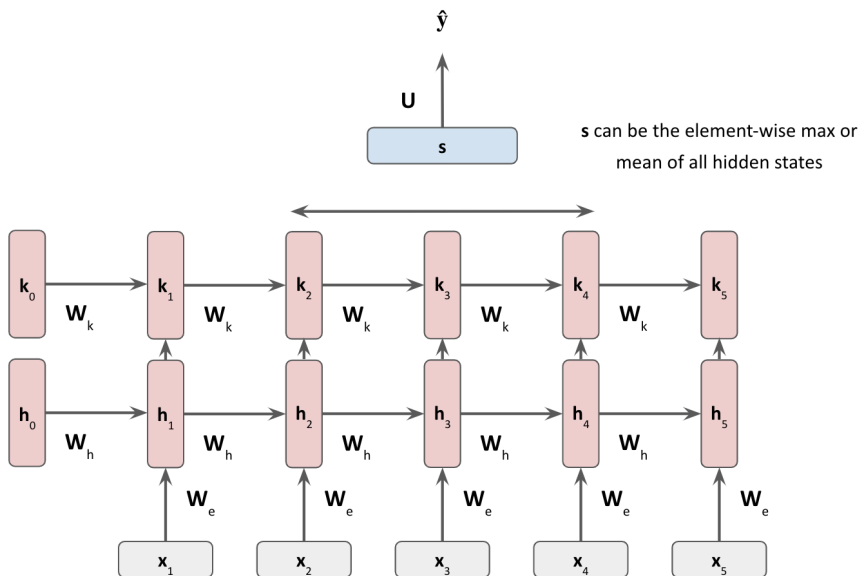


## **Towards Transformer Models**

# Bi-RNNs

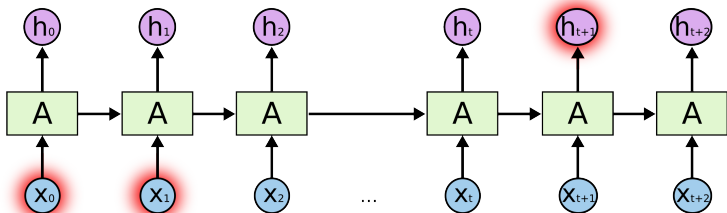


# Multilayer RNNs



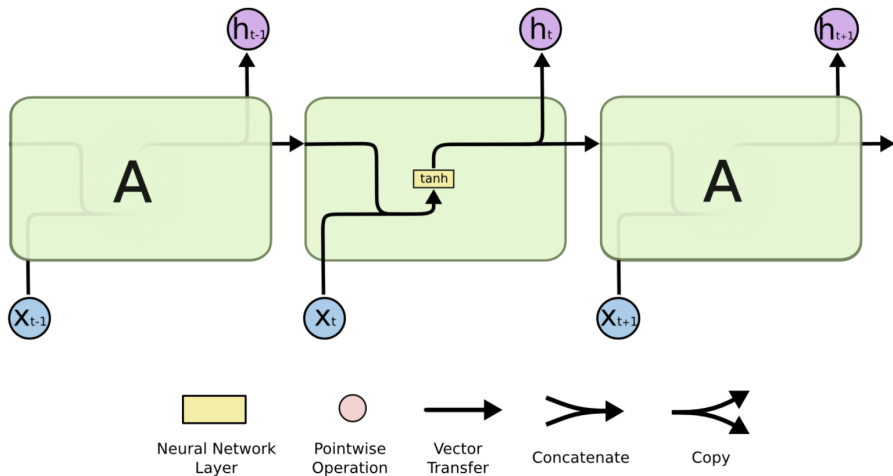
# Vanishing Gradients

- RNNs have a crucial issue: **vanishing and exploding gradients**.
- Both occur when we backpropagate through time with multiplying several times by  $W$ . If  $W$ 's parameters are small, gradients can vanish to zero. If they are large, they can explode.
- This is an issue as it does not allow to model far away context (vanishing) or to properly converge (exploding).
- Solutions:
  - 1 Exploding: **gradient clipping**.
  - 2 Vanishing: **Long Short-Term Memory** networks.



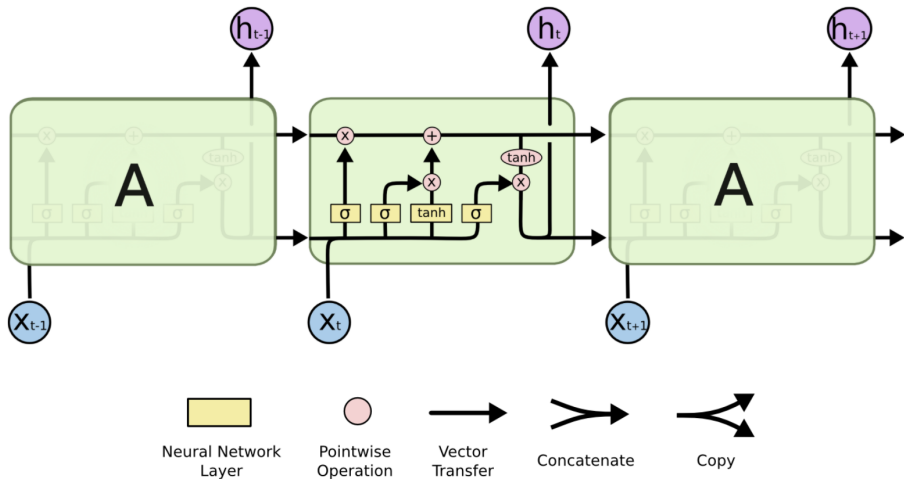
Credit: <https://colah.github.io/posts/2015-08-Understanding-LSTMs>

## A different view on RNNs



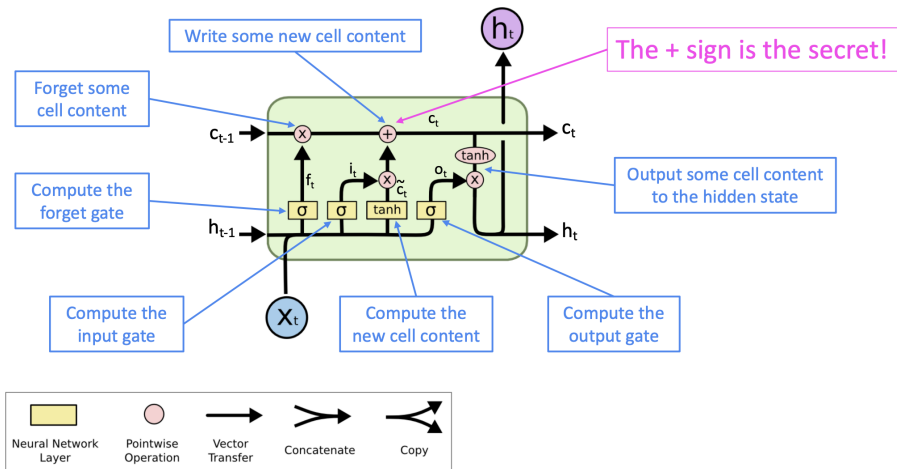
Credit: <https://colah.github.io/posts/2015-08-Understanding-LSTMs>

# LSTMs



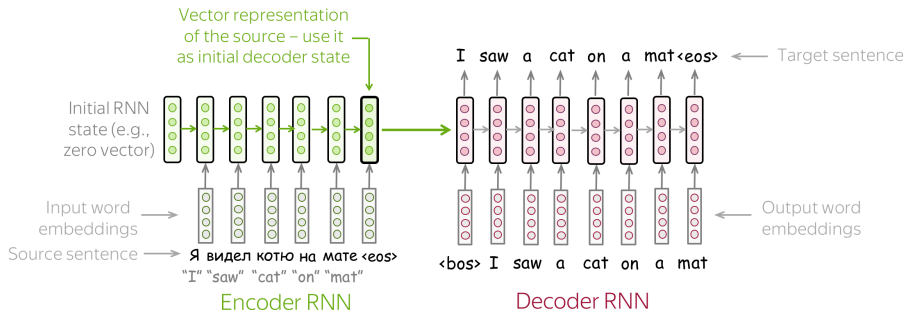
Credit: <https://colah.github.io/posts/2015-08-Understanding-LSTMs>

# LSTMs



*Credit: Stanford CS224N*

# Seq2Seq



Credit: Lena Voita [https://lena-voita.github.io/nlp\\_course.html](https://lena-voita.github.io/nlp_course.html)

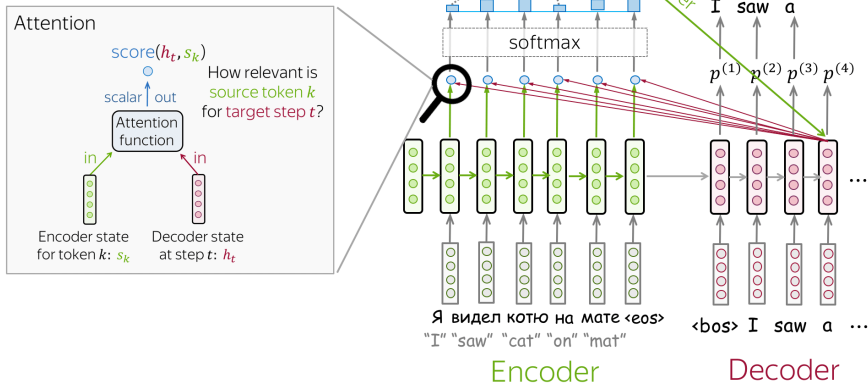


# Attention

Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

A model can learn to “pay attention” to the most relevant source tokens for each step



Credit: Lena Voita [https://lena-voita.github.io/nlp\\_course.html](https://lena-voita.github.io/nlp_course.html)

## What's next

- The most recent advances in neural networks for NLP come from the shift from recurrent architectures to using **attention-based architectures**.
- For your reading assignments, you have explored **transformers** (which combine attention with other ideas from neural networks literature) and will explore **BERT** (which is a neural language model making use of transformers).
- While there is much more to it, in this way you will have a window into contemporary models for NLP.

The next part of the course turns to other topics instead: Web scraping and APIs, recommender systems, corpus annotation, sentiment analysis and clustering with topic modelling, ethics.

## Extras

# Logistic regression

- Applied to our binary classification task, we have that:

$$\hat{p}(y = 1|\mathbf{x}) = \sigma(z_{\mathbf{x}}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$
$$\hat{p}(y = 0|\mathbf{x}) = 1 - \sigma(z_{\mathbf{x}}) = \frac{e^{-\mathbf{w} \cdot \mathbf{x}}}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Then, we just need to use a **decision boundary** to assign the class given the estimated probabilities:

$$\hat{y} = \begin{cases} 1 & \text{if } \hat{p}(y = 1|\mathbf{x}) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

- So, we have defined our data and task, and have a model.  
What do we miss?

## Logistic regression: Cross-entropy

- We need a loss function. Let us use MLE to find one.
  - ▶ Maximize the conditional probability of observing the data given a distribution
- We have that  $p(y|\mathbf{x})$  follows a Bernoulli distribution given that we only have two discrete outcomes  $(0, 1)$ , hence:

$$p(y|\mathbf{x}) = \hat{y}^y(1 - \hat{y})^{1-y}$$

- As usual, let us move to log space and add a minus to switch to a minimization problem (note we work with a single data point  $(\mathbf{x}, y)$  for now):

$$\begin{aligned} -\log p(y|\mathbf{x}) &= -\log [\hat{y}^y(1 - \hat{y})^{1-y}] \\ &= -[y\log \hat{y} + (1 - y)\log(1 - \hat{y})] \end{aligned}$$

- Let us now plug-in the Sigmoid and call it the loss:

$$\mathcal{L}_{\mathbf{x}}(\mathbf{w}) = -[y\log \sigma(\mathbf{w}\mathbf{x}) + (1 - y)\log(1 - \sigma(\mathbf{w}\mathbf{x}))]$$

## Logistic regression: Cross-entropy

- Let us now plug-in the Sigmoid and call it the loss:

$$\mathcal{L}_x(\mathbf{w}) = -[y \log \sigma(\mathbf{w}\mathbf{x}) + (1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))]$$

- The loss on the whole dataset is going to be (note we are already in log space thus we can sum):

$$\mathcal{L}(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log \sigma(\mathbf{w}\mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}\mathbf{x}_i))]$$

- Equivalent to calculating cross-entropy for the Bernoulli distribution
- To this we can, as usual, attach regularization:

$$\mathcal{L}_{L_2}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

## Logistic regression: Optimization via SGD

- The last missing bit is how to find good parameters  $\mathbf{w}$ : we can use SGD.
  - ▶ There is no analytical solution, unlike linear regression
- It turns out that the derivative for one data point  $\mathbf{x}$  is (w.o. regularization):

$$\frac{\partial \mathcal{L}_{\mathbf{x}}(\mathbf{w})}{\partial \mathbf{w}_j} = [\sigma(\mathbf{w}\mathbf{x}) - y] \mathbf{x}_j$$

- For multiple data points, we just sum (w.o. regularization), and with this we are good to go for SGD:

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N [\sigma(\mathbf{w}\mathbf{x}_i) - y_i] \mathbf{x}_{ij}$$

- *Full derivation as an extra, below.*

# Full derivation for logistic regression

- First, we need some notable derivatives:

$$\frac{\partial \log(x)}{\partial x} = \frac{1}{x}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \rightarrow \text{chain rule}$$



# Full derivation for logistic regression

- Then:

$$\begin{aligned}\frac{\partial \mathcal{L}_{\mathbf{x}}(\mathbf{w})}{\partial w_j} &= -\partial [y \log \sigma(\mathbf{w}\mathbf{x}) + (1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))] \\ &= -[\partial y \log \sigma(\mathbf{w}\mathbf{x}) + \partial(1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))] \\ &= -\frac{y}{\sigma(\mathbf{w}\mathbf{x})} \partial \sigma(\mathbf{w}\mathbf{x}) - \frac{1 - y}{1 - \sigma(\mathbf{w}\mathbf{x})} \partial(1 - \sigma(\mathbf{w}\mathbf{x})) \rightarrow \text{chain rule} \\ &= -\left[ \frac{y}{\sigma(\mathbf{w}\mathbf{x})} - \frac{1 - y}{1 - \sigma(\mathbf{w}\mathbf{x})} \right] \partial \sigma(\mathbf{w}\mathbf{x}) \rightarrow \text{re-arrange}\end{aligned}$$

- *Exercise: plug-in the derivative of the Sigmoid and re-arrange yourself to reach:*

$$\dots = [\sigma(\mathbf{w}\mathbf{x} - y)] x_j$$

# Full derivation for logistic regression

- In case you were wondering:

$$\begin{aligned}\frac{\partial \sigma(x)}{\partial x} &= \partial \frac{1}{1 + e^{-x}} \\ &= \partial [1 + e^{-x}]^{-1} \\ &= \frac{e^{-x}}{1 + e^{-x}} \frac{1}{1 + e^{-x}} \\ &= \frac{(1 + e^{-x}) - 1}{1 + e^{-x}} \sigma(x) \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

- *Exercise, derive:*

$$\frac{\partial \log \sigma(x)}{\partial x} = \sigma(-x)$$

# Why MSE and cross-entropy?

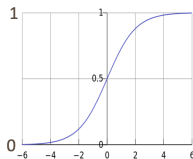
- It turns out that, given some standard assumptions on our models, using those two losses corresponds to doing Maximum Likelihood Estimation. See <https://www.expunctis.com/2019/01/27/Loss-functions.html>.
- If you are curious about the information theory underpinning cross-entropy, read this: <http://colah.github.io/posts/2015-09-Visual-Information>.

# NN activation functions

Several non-linear activation functions have been proposed. A good default options is ReLU.

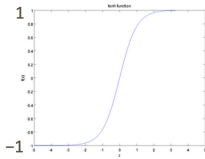
logistic ("sigmoid")

$$f(z) = \frac{1}{1 + \exp(-z)}.$$



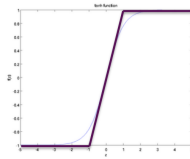
tanh

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}},$$



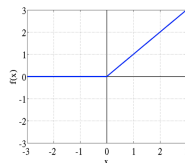
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$



ReLU (Rectified Linear Unit)

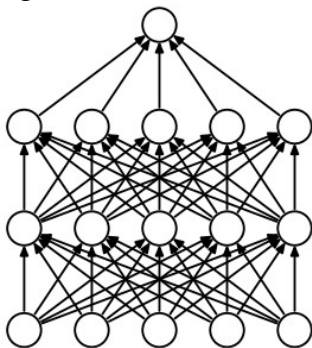
$$\text{rect}(z) = \max(z, 0)$$



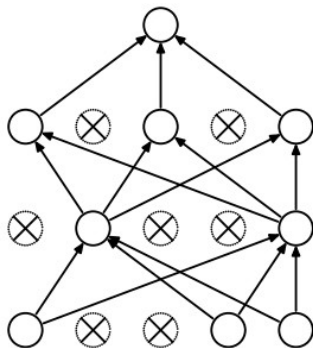
*Credit: Stanford CS224N.*

# NN regularization via dropout

Dropout's idea is to mask a random set of neuron connections at training time, in order to compel the network to learn redundant paths and avoid overfitting.



(a) Standard Neural Net



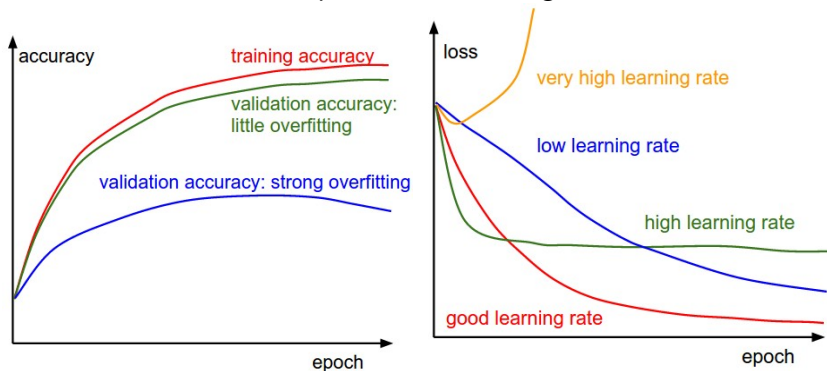
(b) After applying dropout.

*Credit: Srivastava et al.*

*<https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>.*

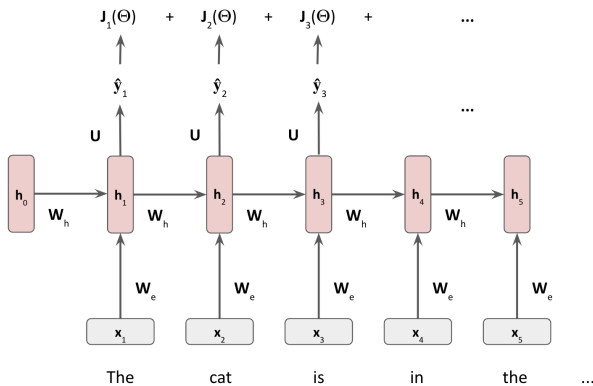
# NN under/overfitting and learning rates

Two illustrations on how to spot correct learning behaviour.



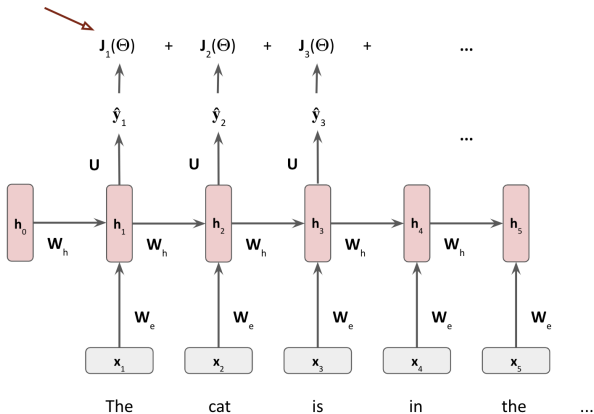
*Credit: Andrej Karpathy via Stanford's CS231N.*

# Training RNNs



# Training RNNs

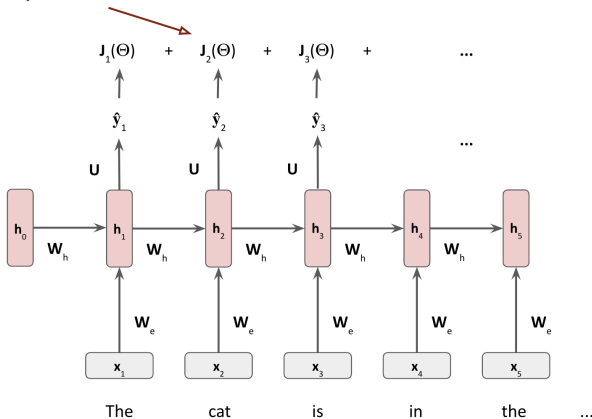
Negative log  
probability for "cat"





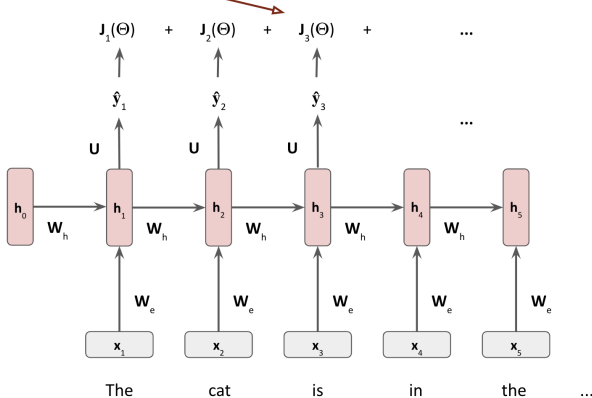
# Training RNNs

Negative log  
probability for "is"



# Training RNNs

Negative log  
probability for "in"



## Recall perplexity?

A second look at perplexity:

- We defined it as the inverse probability of the corpus, normalized by the number of words. For a corpus composed of  $n$  words:

$$PP = P(w_1, w_2, \dots, w_n)^{-\frac{1}{n}} = \prod_{i=1}^n \left( \frac{1}{P(w_i | w_1, \dots, w_{i-1})} \right)^{\frac{1}{n}}$$

- It is actually equal to the exponential of the cross-entropy loss:

$$PP = \prod_{i=1}^n \left( \frac{1}{\hat{y}_{w_{i+1}}} \right)^{\frac{1}{n}} = \exp \left( \frac{1}{n} \sum_{i=1}^n -\log(\hat{y}_{w_{i+1}}) \right) = \exp(\mathcal{L})$$

- So *lower perplexity == lower loss == higher data likelihood*.