



Natural Language Processing



# Language Models and Word Embeddings

Natural Language Processing

Some slide content based on textbooks:

*Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition* by Daniel Jurafsky and James H. Martin

# Lecture content:

ALICE was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversational dialogue: "and what is the use of a book," thought Alice, "without pictures or conversation?" This was the question she was considering in her heart when suddenly she said, for the hot day made her very nervous and stupid,) whether the pleasure of eating the daisy-chain would be worth the trouble of getting up and picking the daisies, when all at once a white rabbit with pink eyes ran close by her; there was nothing so very remarkable in that, nor did Alice think it so very much out of the way to hear the Rabbit say to itself, "Oh dear! Oh dear! I shall be too late!" (when she thought over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually took a watch out of its waistcoat-pocket, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket or a watch to take out of it, and

Minor image source: [http://fr.wikipedia.org/wiki/Fichier:White\\_rabbit\\_in\\_Alice\\_in\\_Wonderland.jpg](http://fr.wikipedia.org/wiki/Fichier:White_rabbit_in_Alice_in_Wonderland.jpg)

- What is language modelling?
- Markov models
- Evaluating language models
- Problems with Markov models
- Revision: Neural Networks
- Word embeddings
- Properties of word embeddings
- Applications of word embeddings
- Extensions

# What is Language Modeling?

And why should I care about it?



Source: <https://pixabay.com/photos/handsome-female-girl-people-school-16811/>

# A language what?

**What is the probability of a sequence of words?**

- just how likely you are to hear somebody say it
- some sequences more likely than others, for example:  
 $P(\text{"why does that kid have a bike?"}) > P(\text{"why does that kid have a moustache?"}) > P(\text{"a moustache why that kid does have?"})$

A statistical **language model** is a **probability distribution** over **sequences of words**

- If we have a distribution over sequences of words:

$$P(w_1, w_2, w_3, \dots, w_n) = \phi([w_1, w_2, w_3, \dots, w_n])$$

- we can **condition** the next word on previous words:

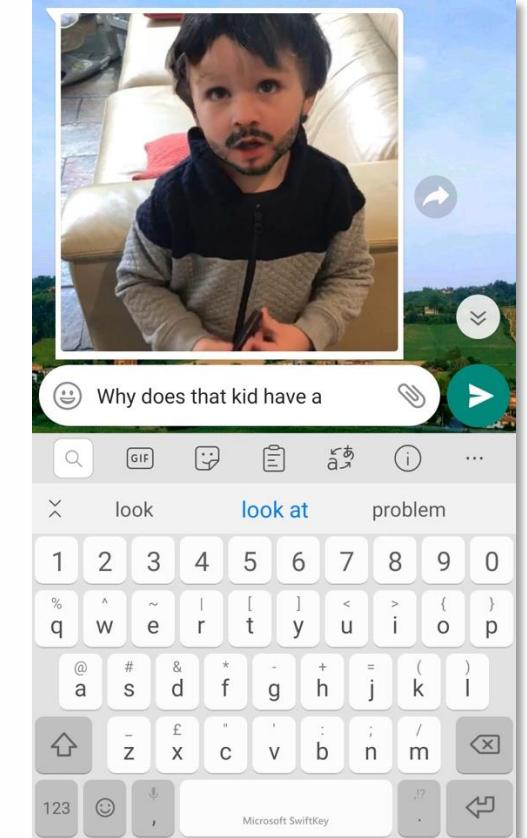
$$P(w_n | w_1, w_2, \dots, w_{n-1}) = P(w_1, w_2, \dots, w_{n-1}, w_n) / P(w_1, w_2, \dots, w_{n-1})$$

- and **sample new sequences** from it

$$w^* \sim P(w | w_1, w_2, \dots, w_{n-1})$$

In other words, **language models** are general-purpose **text generators**

- just like continually selecting the predicted text on your mobile phone



# Buzz

Way back in 2019 there was a lot of buzz about language models when **GPT-2** was released

- OpenAI didn't initially release model parameters due to security concerns
- but could a **language model** really be that dangerous?

The Guardian website screenshot showing a story titled "New AI fake text generator may be too dangerous to release, say creators". The story discusses OpenAI's decision not to release their AI text generator publicly due to fears of misuse. The headline is "New AI fake text generator may be too dangerous to release, say creators". Below the headline, it says "The Elon Musk-backed nonprofit company OpenAI declines to release research publicly for fear of misuse". There is a black and white photo of a man with a cigarette in his mouth.

Forbes website screenshot showing a story titled "OpenAI's Realistic Text-Generating AI Triggers Ethics Concerns". The author is William Falcon, a contributor from Science. The story discusses the ethical concerns raised by GPT-2. Below the article, there is a long quote from Epicurus about pleasure and pain.

CNN Business website screenshot showing a story titled "This AI is so good at writing that its creators won't let you use it". The author is Rachel Metz. The story discusses the creators' decision not to release GPT-2's parameters.

The screenshot shows a BBC News article from June 2017. The headline is "Prices for fake news campaigns revealed". Below the headline is a graphic titled "Tips for spotting false news" with a subtext: "It's possible to spot false news. As we work to limit the spread, check out a few ways to identify whether a story is genuine." A callout box highlights a quote: "Mounting a year-long fake news campaign can cost about \$400,000 (£315,000), suggests a report." The article discusses costs for setting up fake social media profiles, writing stories, and spreading them via followers. The BBC navigation bar at the top includes Home, News, Sport, More, and a search icon.

# Impact ...

In 2019 lots of discussion about fake news campaigns

- using Language Models to create fake news could be a **lucrative business model** ...
- and in geopolitics *the automated pen may be mightier than the sword*

The screenshot shows an NPR article titled "Did Fake News On Facebook Help Elect Trump? Here's What We Know". The author is Danielle Kurtzleben. The article features a photo of Mark Zuckerberg. The background of the page has a red border around the title area.

# How do Language Models work?

# How do Language Models work?

## Language Models

- discover **statistical regularities** in text
- and use them to predict the next word

## Example of a statistical regularity?

- next token is often a **repeat** of a previous word from same text

If we can predict with some accuracy next word:

- iterating forward allows us to predict entire sentences

So how can we discover and use such statistical regularities?

## CHAPTER I. Down the Rabbit-Hole

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversations?'

So she was considering in her own mind (as well as she could, for the hot day made her feel very sleepy and stupid), whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a White Rabbit with pink eyes ran close by her.



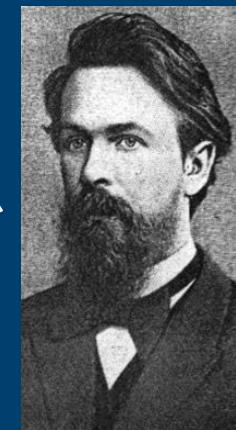
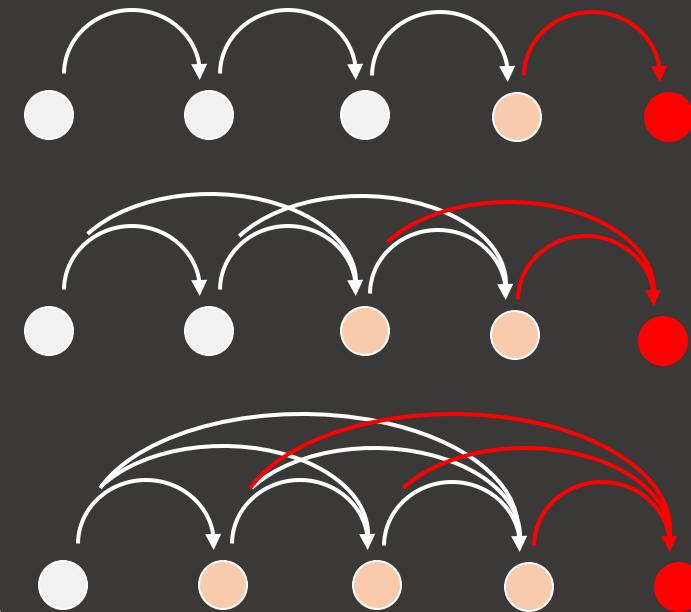
Source: <http://www.illustrationsbyjohn Tenniel.com/Black-and-white-illustration-from-The-Nursery-Alice-in-Wonderland-07.jpg>

# Markov Models

Hey Andrey, you look like a smart guy. So tell me, ...

- if natural language utterances can be of **arbitrary length**,
- and we only have a **finite set of parameters** to model them
- how can we define a probability distribution over them?

Why not predict the next word based on **fixed number** of **previous words**? \*\*



Andrey Markov

famous mathematician

Image source: [https://en.wikipedia.org/wiki/Andrey\\_Markov](https://en.wikipedia.org/wiki/Andrey_Markov)

# Markov models

Simplest models count **n-grams** in large corpus

- n-gram = sequence of  $n$  words

Query	Google Hits
play station	16100000
play sport	2680000
play gym	2430000
...	
play <b>croquet</b>	221000

$$P(\text{croquet} \mid \text{play}) = \frac{N(\text{play croquet})}{N(\text{play})}$$

'For the Duchess. An invitation from the Queen to play **croquet**.' The Frog-Footman repeated, in the same solemn tone, only changing the order of the words a little, 'From the Queen. An invitation for the Duchess to play ... ???'

- longer n-grams give **better predictions**

Query	Google Hits
to play sport	403000
to play <b>croquet</b>	119000
...	
to play gym	1790
to play station	178

$$P(\text{croquet} \mid \text{to play}) = \frac{N(\text{to play croquet})}{N(\text{to play})}$$

Query	Google Hits
Duchess to play <b>croquet</b>	11000
Duchess to play station	0
Duchess to play sport	0
Duchess to play gym	0
...	

$$P(\text{croquet} \mid \text{Duchess to play}) = \frac{N(\text{Duchess to play croquet})}{N(\text{Duchess to play})}$$

# Smoothing, Back-off & Interpolation

**Smoothing:** add small constant to all counts before estimating probabilities

$$P(w_n | w_{n-1}, w_{n-2}) = [count(w_{n-2}w_{n-1}w_n) + \alpha] / [count(w_{n-2}w_{n-1}) + V^*\alpha]$$

- here  $\alpha$  is *pseudocount* (often  $\alpha = 1$ ), which **prevents zero probability** for previously unseen n-grams
- and  $V$  is the vocab size

**Backoff:** don't invent value for unseen  $n$ -gram just but *backoff* to  $(n-1)$ -gram

- so use trigram count if you have it, otherwise revert to bigram model, else unigram
- e.g., when estimating  $P(z|x,y)$ : if  $count(xyz)=0$  use  $P(z|y)$  instead
- and if  $count(yz)$  is also zero, use  $P(z)$  instead

**Interpolation:** interpolate higher and lower order models

- e.g. trigram, bigram, unigram estimates:

$$\hat{P}(w_n | w_{n-1}w_{n-2}) = \lambda_1 P(w_n | w_{n-1}w_{n-2}) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n) \quad \sum_i \lambda_i = 1$$

- how to set the lambdas?
- choose values that maximize probability on held-out **development** data

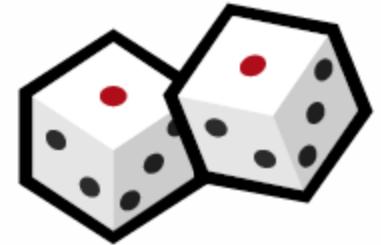
# Generating text from a Markov model

# Choices for generating text

Can use estimate of **next word probability** in various **ways** to generate text:

- **Greedy**  
choose most probable term:  $w^* = \operatorname{argmax}_t P(w_n=t / w_{n-k}, \dots, w_{n-1})$
- **Random sampling**  
sample term according to probabilities:  $w^* \sim P(w_n / w_{n-k}, \dots, w_{n-1})$
- **Top- $k$  sampling**  
limit sampling to  $k$  most likely terms:  $w^* \sim P(w_n / w_{n-k}, \dots, w_{n-1}) \mathbf{1}(w_n \in \text{top-}k)$
- **Temperature sampling**  
limit sampling to likely terms by raising probabilities to power  $w^* \sim P(w_n / w_{n-k}, \dots, w_{n-1})^{1/T}$   
where  $T$  is the temperature (higher temperature = more uniform sampling)
- **Beam search**  
search forward one step at a time for most likely sequence  $(w_n, w_{n+1}, w_{n+2}, w_{n+3}, w_{n+4}, \dots)$   
while limiting the search space to a maximum set of  $k$  candidate sequences

Greedy techniques always produce exactly same text, while sampling produces different text each time



# Examples of generated text

Trained on a corpus of **Wall Street Journal** articles

- using unigrams  $P(\text{word})$ , text is just random:

*unigram:* Months the my and issue of year foreign new exchange's september  
were recession exchange new endorsed a acquire to six executives

- using bigrams  $P(\text{word}/\text{word}_1)$ , text is a lot better:

*bigram:* Last December through the way to preserve the Hudson corporation  
N. B. E. C. Taylor would seem to complete the major central planners one  
point five percent of U. S. E. has already old M. X. corporation of living on  
information such as more frequently fishing to keep her

- and trigrams  $P(\text{word}/\text{word}_1, \text{word}_2)$ , text is almost believable:

*trigram:* They also point to ninety nine point six billion dollars from two  
hundred four oh six three percent of the rates of interest stores as Mexico and  
Brazil on market conditions

# Examples of generated text

- Generated text using a **corpus of Shakespeare's plays** to estimate different Markov models:

Unigram

- To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have
- Every enter now severally so, let
- Hill he late speaks; or! a more to leg less first you enter
- Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like

Bigram

- What means, sir. I confess she? then all sorts, he is trim, captain.
- Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.
- What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?
- Enter Menenius, if it so many good direction found'st thou art a strong upon command of fear not a liberal largess given away, Falstaff! Exeunt

Trigram

- Sweet prince, Falstaff shall die. Harry of Monmouth's grave.
- This shall forbid it should be branded, if renown made it empty.
- Indeed the duke; and had a very good friend.
- Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

Quadrigram

- King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;
- Will you not tell me who I am?
- It cannot be but so.
- Indeed the short and the long. Marry, 'tis a noble Lepidus.

# Bit like Finnegans Wake ...

Output from low-order Ngram LM

- reads a bit like Finnegans Wake from James Joyce
- non-sensical but potentially grammatical text

## CHAPTER ONE

riverrun, past Eve and Adam's, from swerve of shore to bend of bay, brings us by a  
commodius vicus of recirculation back to Howth Castle and Environs.  
Sir Tristram, violer d'amores, fr'over the short sea, had passen-core rearived from  
North Armorica on this side the scraggy isthmus of Europe Minor to wielderfight  
themselfe to Laurens County's gorgios while they went doublin their mumper all  
the time: nor avoice from afire bellowsed mishe mishe to tauftauf thuartpeartick not  
yet, though venissoon after, had a kidscad buttended a bland old isaac: not yet,  
though all's fair in vanessy, were sosie sesthers wroth with twone nathandjoe. Rot a  
peck of pa's malt had Jhem or Shen brewed by arclight and rory end to the  
regginbow was to be seen ringsome on the aquaface.

The fall (bababadalgharaghakamminarronnkonnbronntonne-  
ronntuonnthunntravarhounawnskawntooohohoordenenthur — nuk!) of a once  
wallstrait oldparr is retaled early in bed and later on life down through all christian  
minstrelsy. The great fall of the offwall entailed at such short notice the pftjschute  
of Finnegans, erse solid man, that the humptyhillhead of humself promptly sends an  
unquiring one well to the west in quest of his tumptytumtoes: and their  
upturnpikepointandplace is at the knock out in the park where oranges have been  
laid to rust upon the green since dev-linsfirst loved livvy.

What clashes here of wills gen wonts, oystrygods gaggin fishy-gods! Brékkék  
Kékkek Kékkek! Kóax Kóax Kóax! Ualu Ualu Ualu! Quaouauh! Where  
the Baddelaries partisans are still out to mathmaster Malachus Micgranes and the  
Verdons cata-pelting the camibalistics out of the Whoytebowe... Assiegates and 1

How can we evaluate  
a Language Model?



# Types of LM evaluation

How do we know if one model is better than another?

## **Extrinsic** evaluation

- use model in a downstream task (e.g. as a spelling corrector) and evaluate performance on that task

## **Intrinsic** evaluation

- train parameters of model on training set and test model performance on held-out dataset
- use likelihood that model would produce the new data to evaluate how well it is doing

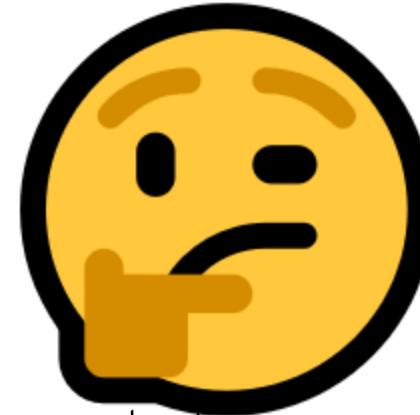
# What is Perplexity?

Dictionary  
Definitions from [Oxford Languages](#) · Learn more

 **perplexity**  
*noun*

1. inability to deal with or understand something.  
"she paused in perplexity"

Similar: confusion, bewilderment, puzzlement, bafflement, incomprehension



Perplexity of language model quantifies level of **surprise/confusion** at seeing new text

- measures how **unlikely** the **observed data** is under the model

Calculating perplexity:

- compute probability of observed sequence  $P(w_1, w_2, \dots, w_n)$  under model
  - for a bigram model have:  $P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2/w_1)P(w_3/w_2)\dots P(w_n/w_{n-1})$
- normalize probability for length of text sequence
  - i.e. compute “average” per-word probability  $P(w_1, w_2, \dots, w_n)^{1/n}$
- invert probability to compute uncertainty:

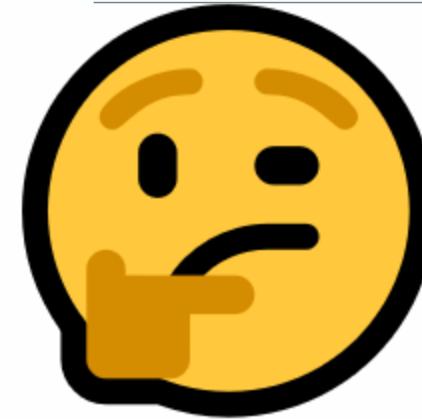
$$PP(\mathbf{w}) = 1/P(w_1, w_2, \dots, w_n)^{1/n} = P(w_1, w_2, \dots, w_n)^{-1/n} = \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

- minimizing perplexity same as maximizing probability
- so **lower perplexity = better model**

# Perplexity, nLL and CE

How does Perplexity relate to measures we usually use to train/evaluate predictive models?

- Spoiler: all pretty much the same concept!
- **Perplexity**
  - inverse of geometric mean of probability of correctly predicting next word
- **Negative log Likelihood (nLL)**
  - negative logarithm of probability of sequence
  - dividing by sequence length gives per-word nLL
  - so perplexity is just  $2^{nLL}$
- **CrossEntropy (CE)**
  - expected (empirical average) log surprise under model
  - number of bits needed to quantify surprise



$$\begin{aligned} PP(W) &= 2^{-\frac{1}{N} \log_2 P(w_1, w_2, \dots, w_N)} \\ &= (2^{\log_2 P(w_1, w_2, \dots, w_N)})^{-\frac{1}{N}} \\ &= P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}} \end{aligned}$$

$$\begin{aligned} H(W) &= -\frac{1}{N} \log_2 P(W) \\ &= -\frac{1}{N} \log_2 P(w_1, w_2, \dots, w_N) \end{aligned}$$

$$PP(W) = 2^{H(W)} = 2^{-\frac{1}{N} \log_2 P(w_1, w_2, \dots, w_N)}$$

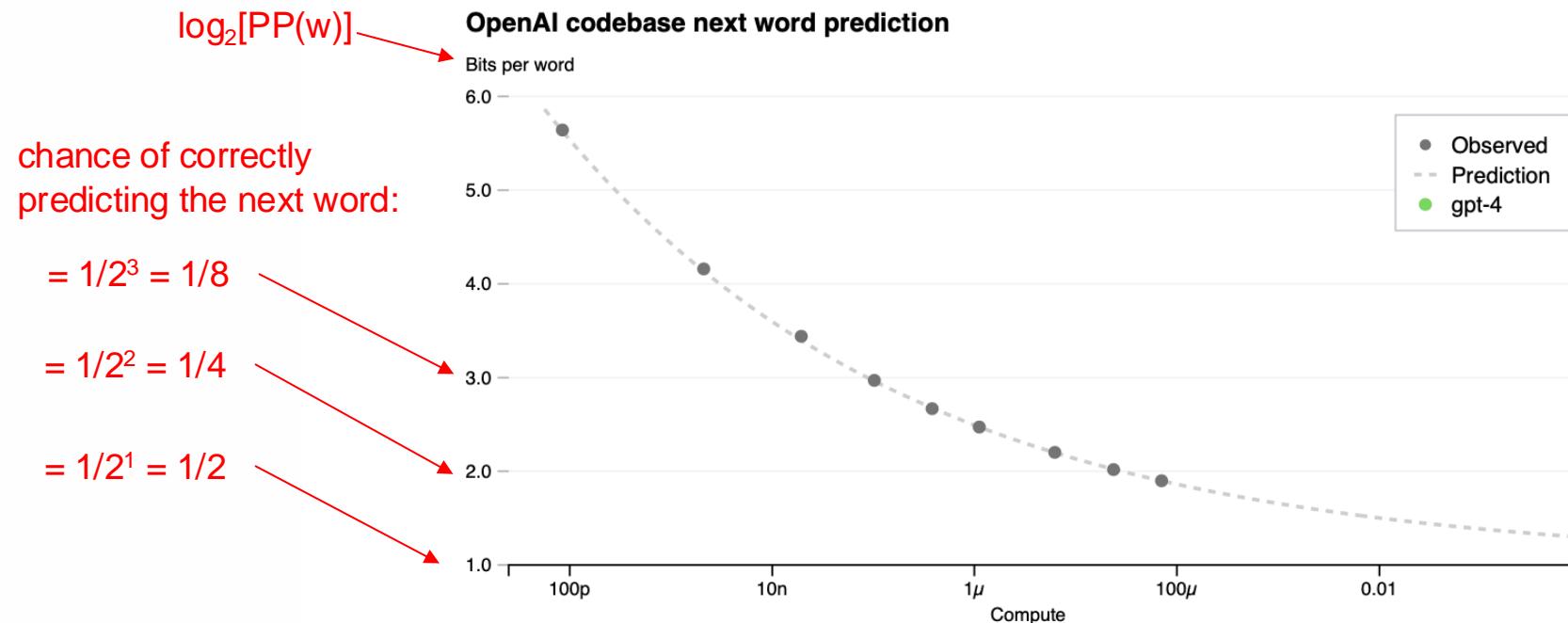
$$\begin{aligned} H(p, q) &= - \sum p(x_i) \log_2 q(x_i) \\ H(p, q) &\approx -\frac{1}{N} \log_2 q(W) \end{aligned}$$

Formulae source <https://towardsdatascience.com/perplexity-in-language-models-87a196019a94>

# Aside: Perplexity of GPT-4

What is the perplexity of GPT-4?

- according to the GPT-4 Technical Report (<https://cdn.openai.com/papers/gpt-4.pdf>)



**Figure 1.** Performance of GPT-4 and smaller models. The metric is final loss on a dataset derived from our internal codebase. This is a convenient, large dataset of code tokens which is not contained in the training set. We chose to look at loss because it tends to be less noisy than other measures across different amounts of training compute. A power law fit to the smaller models (excluding GPT-4) is shown as the dotted line; this fit accurately predicts GPT-4's final loss. The x-axis is training compute normalized so that GPT-4 is 1.

# Problem with Ngram Language Models

# Problem with Markov models

Query	Google Hits
play station	16100000
play sport.	2680000
play gym	2430000
...	
play <b>croquet</b>	221000

Query	Google Hits
to play sport	403000
to play <b>croquet</b>	119000
...	
to play gym	1790
to play station	178

Query	Google Hits
Duchess to play <b>croquet</b>	11000
Duchess to play station	0
Duchess to play sport	0
Duchess to play gym	0
...	

$$P(\text{croquet} | \text{play}) = \frac{N(\text{play croquet})}{N(\text{play})}$$

$$P(\text{croquet} | \text{to play}) = \frac{N(\text{to play croquet})}{N(\text{to play})}$$

$$P(\text{croquet} | \text{Duchess to play}) = \frac{N(\text{Duchess to play croquet})}{N(\text{Duchess to play})}$$

'For the Duchess. An invitation from the Queen to play **croquet**.' The Frog-Footman repeated, in the same solemn tone, only changing the order of the words a little, 'From the Queen. An invitation for the Duchess to play ... ??'

## Problem:

- as  $n$  gets larger, chance of finding sequence in corpus drops exponentially
- meaning there is **never enough training data**
- continuously backing off to shorter n-grams greatly **limits power of the model**

# Markov models just don't scale

In order to generate reasonable language

- need to model **very long distance dependencies**

Memory and data requirements:

- scale exponentially in length of observable dependency
- so **Markov models just don't scale**
  - nonetheless, they were still state-of-the-art not that long ago

Need instead methods that can both

- **generalise** from limited data
- handle **longer dependencies**

'For the Duchess. An invitation from the Queen to play **croquet**. The Frog-Footman repeated, in the same solemn tone, only changing the order of the words a little, 'From the Queen. An invitation for the Duchess to play ... ???'

# Alternative: Word Embeddings

# Word Embeddings - what are they?

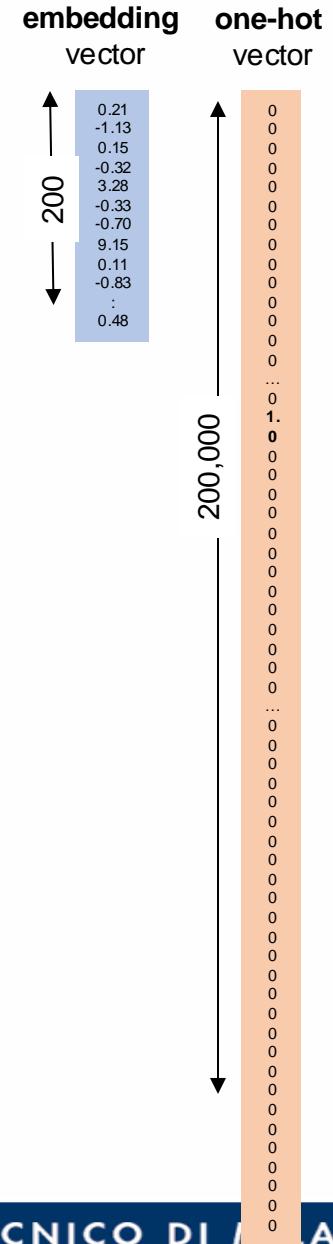
Word embeddings appeared around 2013 and improved performance on just about every NLP task

What are they?

- **dense vectors** representing **words** in a **high dimensional space**
- typically have between 100 & 1000 dimensions
- very low dimensional compared to **one-hot encoding** of terms, since typical document collections have vocabularies of 100k to 1m tokens

Note:

- Just like one-hot encodings word embeddings can be **aggregated** to represent sentences and documents



# Motivating word embeddings

Your task is to fill in the blank in the sentence:

'Sure Sally, let's have a skype call at 3pm \_\_\_\_\_ the 3<sup>rd</sup> of June.'

What word could fit here?

- prepositions: **on, by, before, around, near, ...**
- days of the week: **Monday, Tuesday, Wednesday, ...**
- timezones: **GMT, CET, EST, AEST, ...**

Note: very **few words fit** the context

- those that do come in **groups** that are **semantically related** to one other

Embeddings are produced by **supervised machine learning models**

- models trained to **predict missing word** based on **surrounding context**
- context may include only previous words (causal models)
- or also future words (non-causal models)

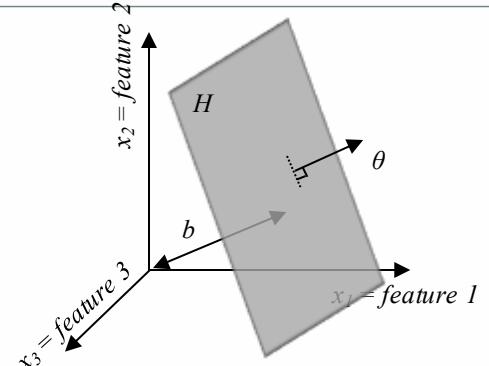
# Supervised Learning Problem

Predicting missing word:

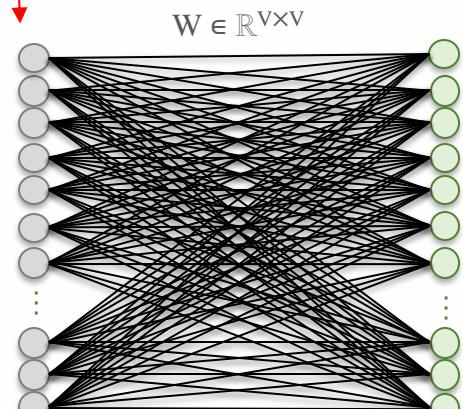
- **Features:** words in current context:
  - “Sure”, “Sally”, “let’s”, “have”, “a”, “skype”, “call”, “at”, “3pm”
- **Target:** missing word from sequence
  - multi-class problem (estimate probability for every word in vocab)

Issue: even **linear classifier** requires **very large number of parameters!**

- multi-class linear classifier (e.g. Logistic Regression ) to predict missing word
- with bag-of-words feature vector (so ignoring word order)
- requires parameter vector (of size of vocabulary) for each vocab term
- so overall number of parameters will be **quadratic** in the size of the vocab
- if vocab=100 thousand, then we would have 10 billion parameters!!
- which **used to be a lot** before deep learning came along 😱 😂



bag-of-words representation  
of words in context window



nodes to predict each possible  
missing word (applying softmax  
produces probability distribution over  
vocabulary)

# Revision: What are Neural Networks

# Brains and Neurons

Brains composed of large number of **neurons**

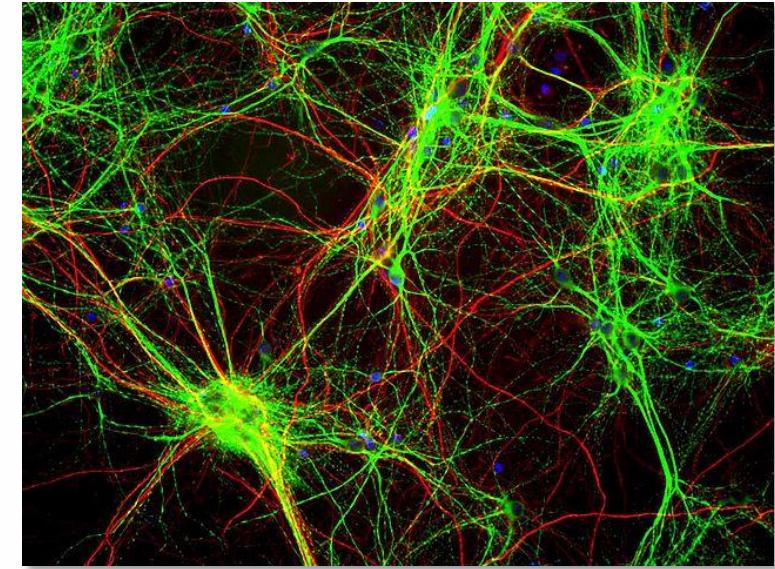
- human brain has tens of billions
- **massive graph:** each node connects to 10,000+ others

Neuron is **simple processing unit**:

- receives electrochemical stimulus from other neurons
- triggers output signal if sum of inputs exceeds threshold

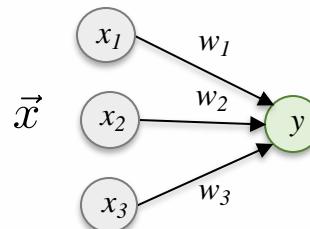
Brain learns to perform complicated tasks

- by connecting neurons in certain patterns



Rat brain under microscope (by Gerry Shaw), source:  
[https://commons.wikimedia.org/wiki/File:Culture\\_of\\_rat\\_brain\\_cells\\_stained\\_with\\_antibody\\_in\\_MAP2\\_\(green\)\\_Neurofilament\\_\(red\)\\_and\\_DNA\\_\(blue\).jpg](https://commons.wikimedia.org/wiki/File:Culture_of_rat_brain_cells_stained_with_antibody_in_MAP2_(green)_Neurofilament_(red)_and_DNA_(blue).jpg)

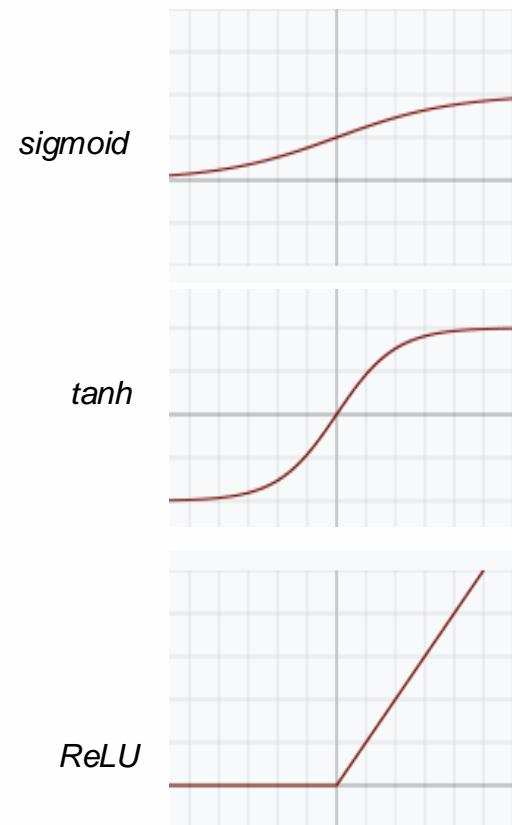
# Artificial Neural Networks



$$y = \tanh(\vec{w} \cdot \vec{x} + w_0)$$

## Artificial Neural Networks (ANNs)

- network of artificial (**simple model** of) neurons
- long history:
  - 1940s: understanding how neurons in nervous system work
  - 1970s: von-Neumann architectures outpaced NNs
  - 1980s: renewed interest, but computation prohibitive
  - 2010s: massive resurgence in popularity
- model neurons **as simple step functions**
  - weight inputs and threshold to produce output value
  - using non-linear **activation function**
    - typical examples: *sigmoid*, *tanh*, *rectified-linear unit (ReLU)*
    - see: [https://en.wikipedia.org/wiki/Activation\\_function](https://en.wikipedia.org/wiki/Activation_function)



# Neural Networks: configuration

Create different types of models by varying:

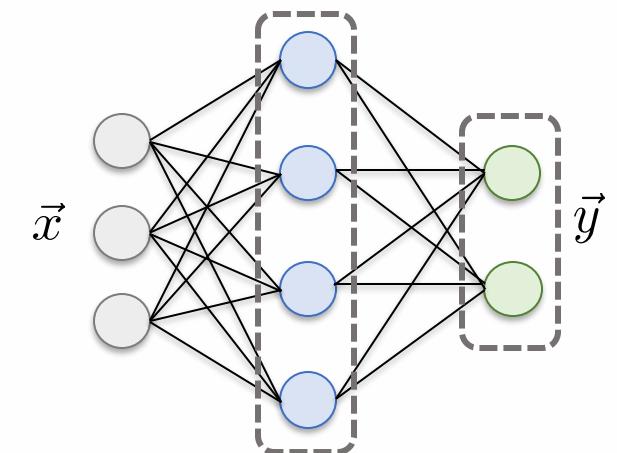
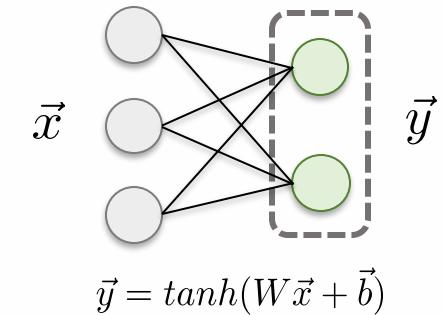
- topology (shape) of network
- activation function used
- loss function optimized (e.g. binary cross-entropy)

NNs made up of neurons arranged in *layers*

- each layer receives inputs from previous layer
- and if activated, emits outputs to next layer
- **internal** (neither input nor output) **layers** are called **hidden layers**
  - allow model to find non-linear decision boundary

Parameters:

- connection of two neurons has associated **weight**
- neuron has associated **bias** (threshold)

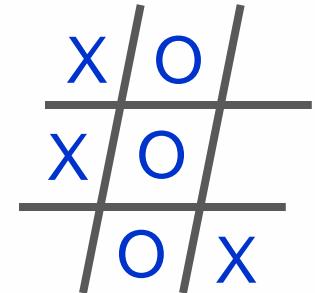


$$\text{CrossEntropy}(\theta) = -\frac{1}{n} \sum_i \log P(y_i|x_i; \theta)$$

# Example non-linear classification problem

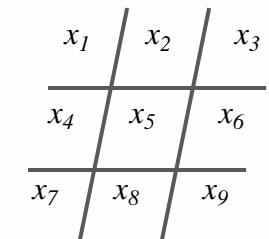
Consider game of **tic-tac-toe**

- aim is to classify boards a **X**-wins, **O**-wins or **nobody**-wins
  - 9 dimensional feature vector:
$$\begin{aligned}x &= (\textcolor{blue}{X}, \textcolor{blue}{O}, \_, \textcolor{blue}{X}, \textcolor{blue}{O}, \_, \_, \textcolor{blue}{O}, \textcolor{blue}{X}) \\&= (1, -1, 0, 1, -1, 0, 0, -1, 1)\end{aligned}$$



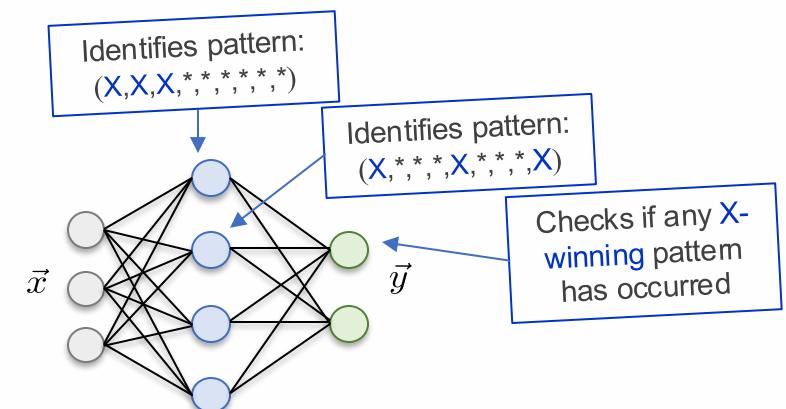
Example of non-linear classification problem

- i.e. problem where no linear decision boundary can be found to separate the classes



Simple two layer neural network can

- first layer of neurons learn to identify simple winning pattern
  - E.g. 3 vertical crosses or 3 diagonal crosses
- second layer identifies the winners
  - i.e. performs disjunction



# Neural Networks: training

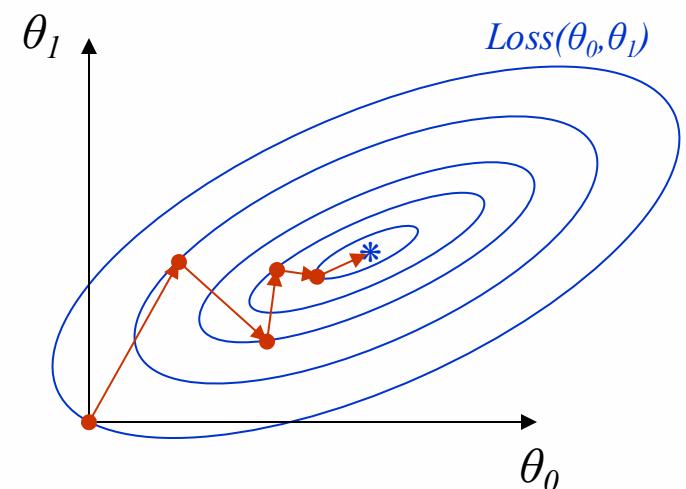
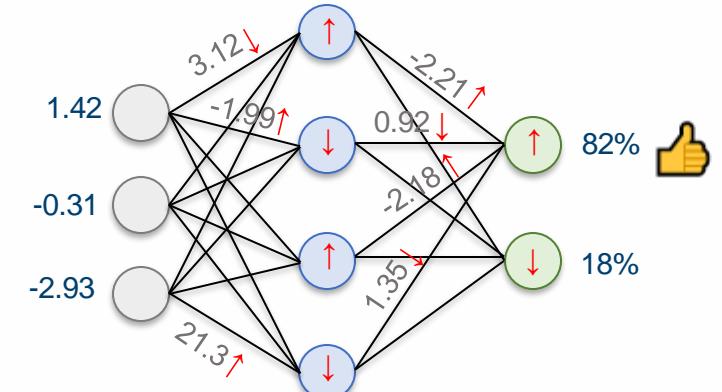
Parameters learnt by process of *backpropagation*

- initialize all **weights/biases** randomly
- for each sample in training set, perform:
  - 1) **forward pass**: feed input through network to calculate **prediction**
  - 2) **backward pass**: update parameters:
    - reward connections producing correct prediction
    - and penalize those contributing to wrong result

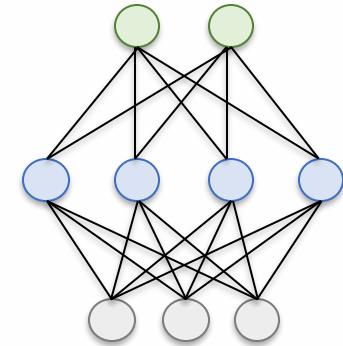
Just a **gradient descent** routine

- where **nested function** is being optimised
- and **chain rule** from calculus used to compute derivative of loss function with respect to each layer of parameters

$$\nabla Loss = \nabla \text{MSE}(a, b) = \left( \frac{\partial \text{MSE}(a, b)}{\partial a}, \frac{\partial \text{MSE}(a, b)}{\partial b} \right)$$



# Neural Networks: properties



Neural networks are extremely powerful learners:

- network with single hidden layer **can learn any function**, providing layer is wide enough
- lots of variation in network structures

Many hyperparameters to tune for NN, including:

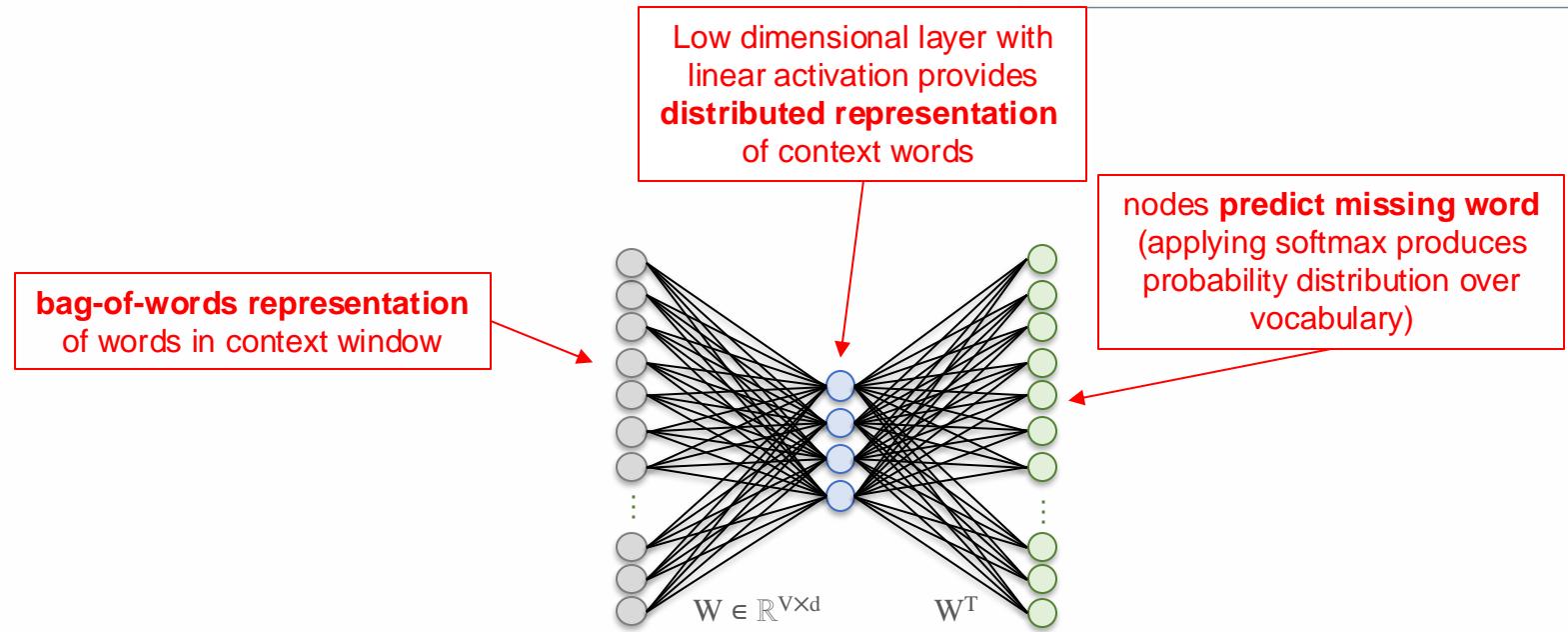
- network architecture,
- type of activation function

Lots of hardware & software optimizations available for training NNs

- more *computationally expensive* than other learners, e.g. decision trees, linear classifiers

# Word Embeddings: Word2Vec, GloVe & related models

# Word2Vec



**Word2Vec** developed in [2013](#) by Mikolov et al.

- following early work by [Bengio et al.](#) in 2003
- later GloVe in [2014](#) by Penington et al.

Word2Vec solved the parameter space issue by using:

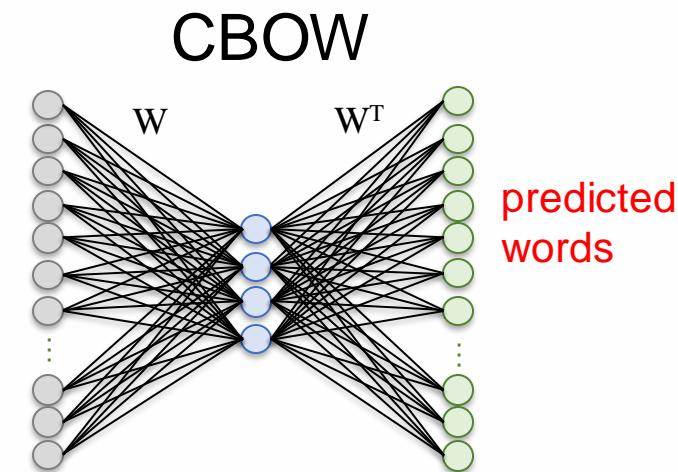
1. bag-of-words representation
2. neural network with single (**linear**) hidden layer
3. training model in “discriminative” fashion by inventing negative examples

# Word2Vec (cont.)

Two versions of Word2Vec:

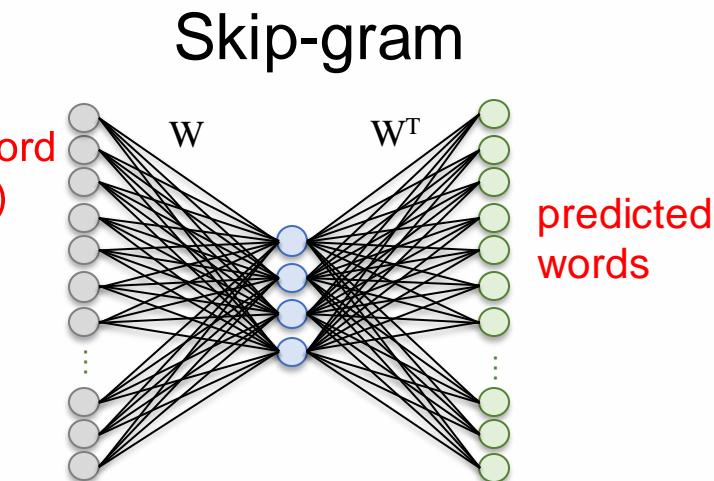
- Continuous Bag of Words (cbow) model
  - trained to predict observed word based on all surrounding **context words**
  - context consists of **all terms occurring** in symmetric window around target term
- Skip-gram (skip) model
  - trained to predict observed word given a single context word

ALL context words  
(bag-of-words representation)



So skip-gram is just 1-to-1 prediction, while cbow is many-to-1 prediction

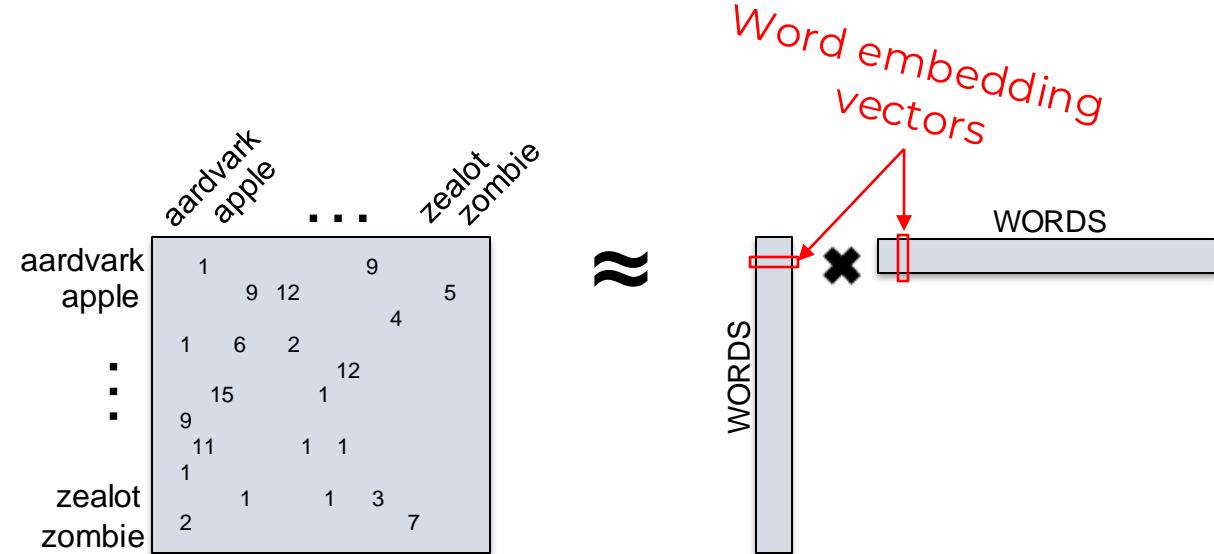
SINGLE context word  
(one-hot-encoding)



# It's all just matrix decomposition ...

Word embeddings can be seen as a form of **matrix decomposition**

- square **count matrix**: vocabulary  $\times$  vocabulary
- contains word **co-occurrences** in text within a **fixed-size context window**
- factorizing **generalises** the information in those windows
- and produces word embedding vectors



Techincal details of  
Word2Vec and GloVe

# Word2Vec – details

Skip-gram model:

- predict target word  $c$  given observed word  $w$  using softmax of dot product of embedding representations:
- directly optimising is computationally intensive because need to sum over all possible target words  $c'$ :

$$p(c|w; \theta) = \frac{e^{v_c \cdot v_w}}{\sum_{c' \in C} e^{v_{c'} \cdot v_w}}$$

- so Mikolov et al. sample some **negative examples** (words that were **NOT observed**)
- and turn problem into **binary classification task** (word pair observed, yes or no?) to avoid nasty summation!
- on GPU optimizing softmax directly is not a problem now

$$p(D = 1|w, c; \theta) = \frac{1}{1 + e^{-v_c \cdot v_w}}$$

Continuous bag of words model (cbow)

- estimated exactly same way except context is **sum of word vectors**

$$\arg \max_{\theta} \sum_{(w,c) \in D} \log \frac{1}{1 + e^{-v_c \cdot v_w}} + \sum_{(w,c) \in D'} \log \left( \frac{1}{1 + e^{v_c \cdot v_w}} \right)$$

Follows explanation in Goldberg et al. <https://arxiv.org/pdf/1402.3722.pdf>

# Alternative Embedding: GloVe

AIM: give probabilistic interpretation to translation in embedding space

- original paper: <https://nlp.stanford.edu/pubs/glove.pdf>
- e.g. translation in space from “steam” to “ice” should increase chance of seeing word “solid”
- so projection of “ice” minus “steam” onto vector “solid” should be function of conditional probabilities

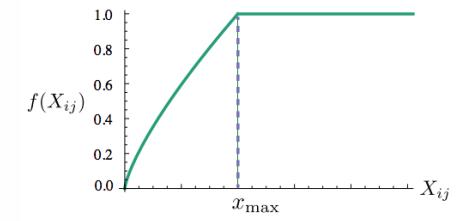
$$(w_i - w_j)^T \tilde{w}_k \approx f\left(\frac{P(k|i)}{P(k|j)}\right)$$

$$P(\text{solid}|\text{ice}) > P(\text{solid}|\text{steam})$$

- solving is equivalent to fitting objective:
  - where  $w_i$  and  $w_k$  are embedding vectors,  $b_i$  and  $b_k$  are biases, and  $X_{ik}$  is their co-occurrence count
- objective approximated by minimising a weighted least squares objective
  - with some tricks needed to make function converge

$$J = \sum_{i,j=1}^V f(X_{ij}) \left( w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2,$$

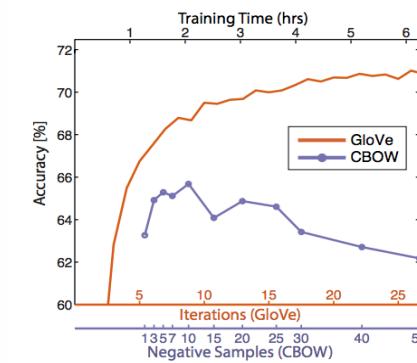
$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik})$$



# Word2Vec or GloVe? which is better?

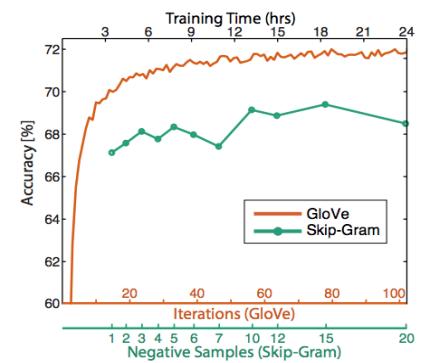
Which is best: Word2Vec (cbow), Word2Vec (skip-gram), or GloVe?

- from GloVe paper:



(a) GloVe vs CBOW

Source: <https://nlp.stanford.edu/pubs/glove.pdf>



(b) GloVe vs Skip-Gram

- according to most sources skip-gram beats cbow
  - except for fastText (see next)
- according to Levy et al. (<http://www.aclweb.org/anthology/Q15-1016>) word2vec is:
  - faster to train
  - has lower memory requirement
  - produces more accurate models

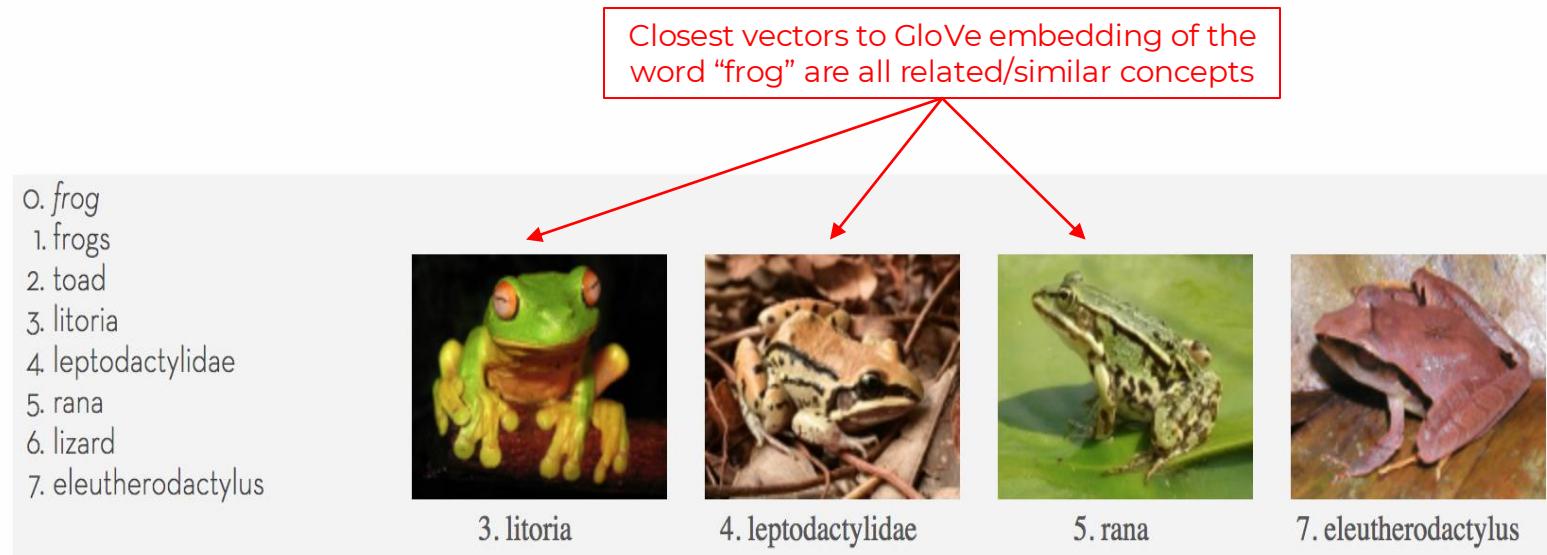
# Properties of Word Embeddings

# Properties of Word Embeddings

Word embeddings have many **interesting** and somewhat surprising **properties**

Semantic clustering:

- neighbours in the embedding space are **semantically related** to one another



Source <https://nlp.stanford.edu/projects/glove/>

# Interpreting Word Embeddings

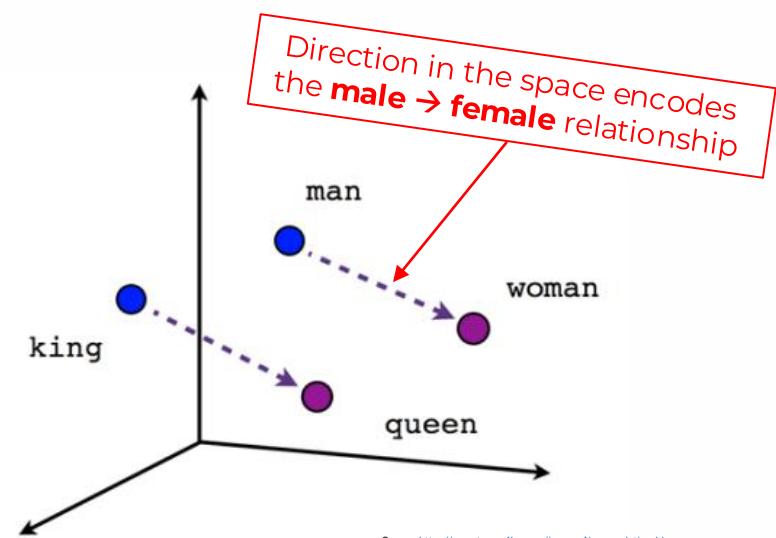
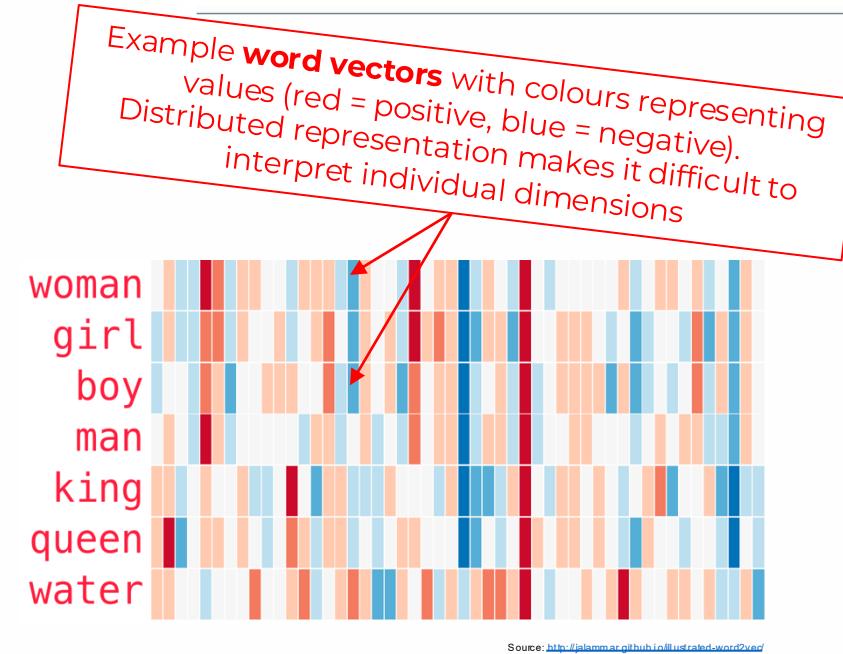
Word embeddings provide **dense distributed representation**

- where individual dimensions aren't usually interpretable

But **translation** in the space is **meaningful**

- certain directions in the space have meaning
- which means that **semantics** is often **additive**
- and **analogies** are encoded in the embedding:

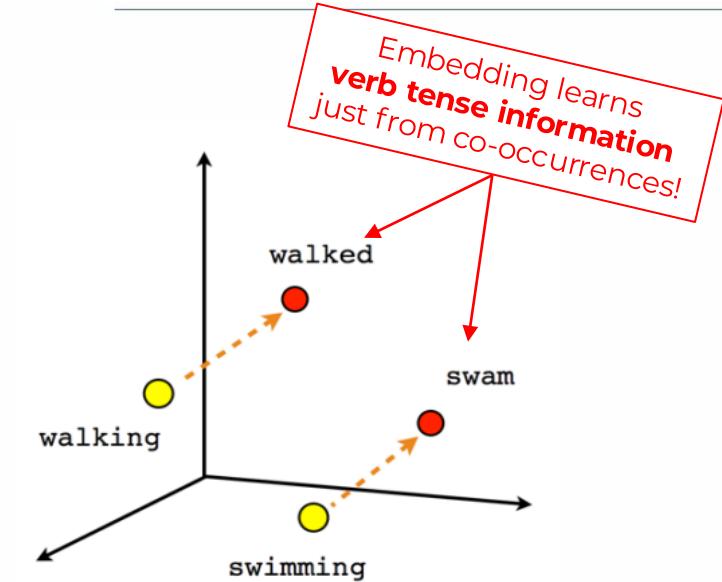
*man* is to *woman* as *king* is to *queen*



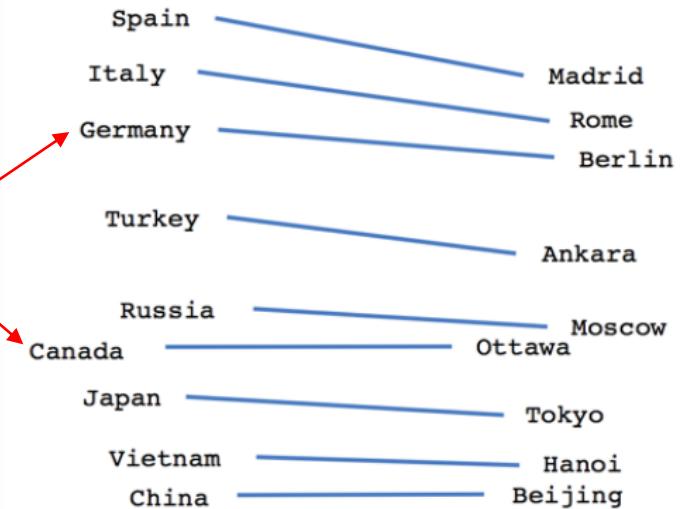
# Knowledge in Word Embeddings

Embeddings discover all sorts of **relationships** between words, for example:

- part-of-speech: *help* → *helpful*, *pain* → *painful*
- type-of relationships: *red*, *green*, *blue* → *colours*
- synonyms: *brave* ≈ *courageous*
- geographic: *Chicago + state* → *Illinois*



and **capital city** → **country** relationships  
again just from context information:  
Madrid + "is the capital city of" → Spain  
Tokyo + "is the capital city of" → Japan  
Hanoi + "is the capital city of" → Vietnam



Source: [https://www.tensorflow.org/images/linear\\_relationships.png](https://www.tensorflow.org/images/linear_relationships.png)

# Uses of Word Embeddings

# Uses of Word Embeddings

## Causal models:

- restrict context words to come before missing word
- language modelling: can be used for predicting next word in a sequence
- with longer dependencies than possible with n-gram models

## Non-causal models:

- can be used as additional **feature vector** for representing words
- improve performance on most tasks
- such as training classifier to detect sentiment
- or translating one language to another
- because they make use of additional domain knowledge (semantics of words)



Christopher Manning (famous NLP researcher)  
[https://en.wikipedia.org/wiki/Christopher\\_D.\\_Manning](https://en.wikipedia.org/wiki/Christopher_D._Manning)

Sriracha sauce: <https://en.wikipedia.org/wiki/Sriracha>

# Word Embeddings help Language Models

Low dimensional representation causes similar terms to share similar descriptions

- allows model to **generalize** from semantically **related examples**
- e.g. **part-of-speech** and **hyponym** (type-of) relationships implicitly encoded in embedding vector in additive manner

the **Duchess** to play ... ??

Examples from corpus with similar contexts:

- the **Queen** to play croquet
- the **Duke** to play chess

the **Duchess** to play ... ??

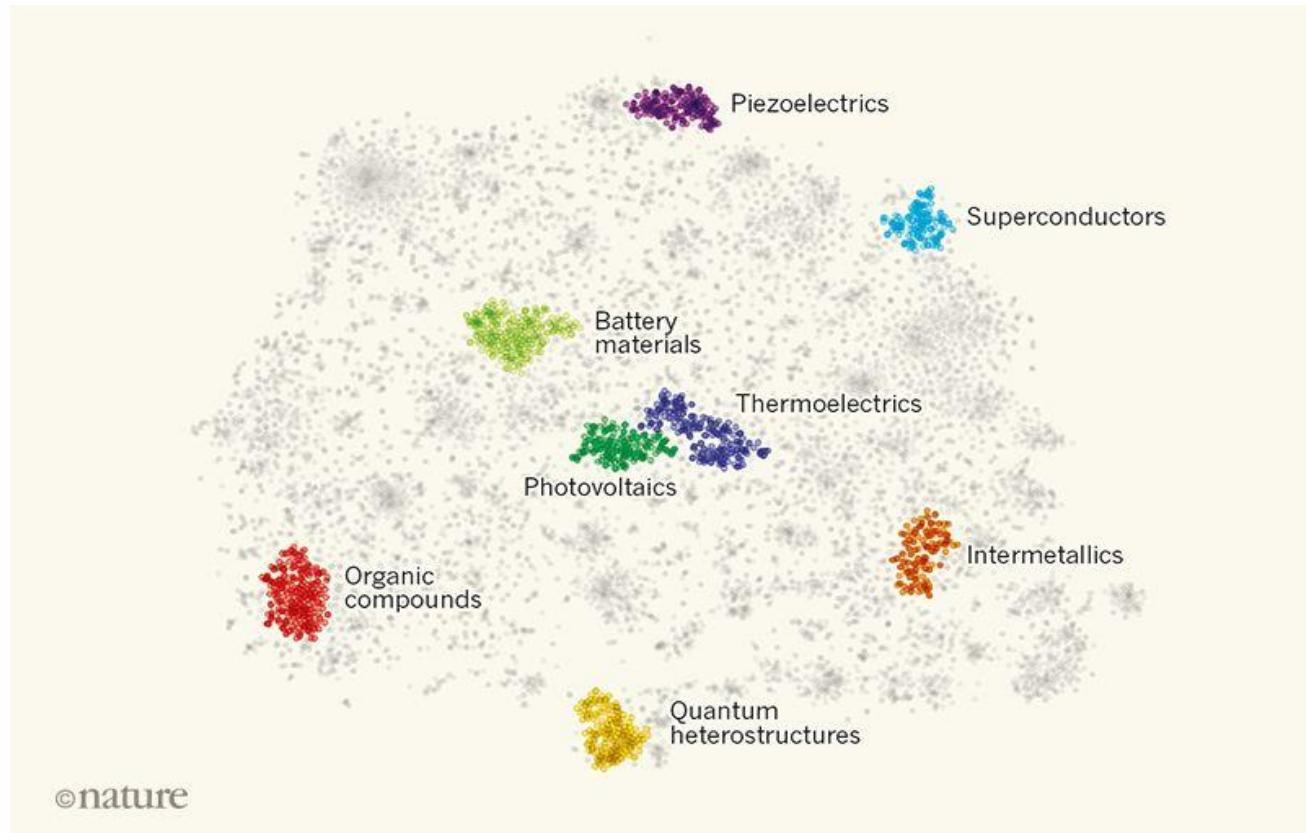
Look for examples with any of these types

- the [**noun+person+female+royal**] to play ...

# Useful for Mining Text

Embeddings place similar concepts close together

- useful for discovering implied (but unknown) properties of them



Visualisation from news article "[Text mining facilitates materials discovery](#)" by Olexandr Isayev ,

**nature**  
Letter | Published: 03 July 2019  
**Unsupervised word embeddings capture latent knowledge from materials science literature**  
Vahe Tshitoyan, John Dagdelen, Leigh Weston, Alexander Isayev, Olga Kononova, Kristin A. Persson, Gerbrand Ceder  
Nature 571, 95–98(2019) | Cite this | 42k Accesses | 13 Citations | PDF version | Subscribe | Search | Login

**Abstract**  
nature > news & views > article

**Text mining facilitates materials discovery**  
Computer algorithms can be used to analyse text to find semantic relationships between words without human input. This method has now been adopted to identify unreported properties of materials in scientific papers.

The total number of materials that can potentially be made – sometimes referred to as materials space – is vast, because there are countless combinations of components and structures from which materials can be fabricated. The accumulation of experimental data that represent pockets of this space has created a foundation for the emerging field of materials informatics, which integrates high-throughput experiments, computations and data-driven methods into a tight feedback loop that enables rational materials design. Writing in *Nature*, Tshitoyan *et al.* report that knowledge of materials science ‘hidden’ in the text of published papers can be mined effectively by computers without any guidance from humans.

RELATED ARTICLES

Read the paper: [Unsupervised word embeddings capture latent knowledge from materials science literature](#)

Machine learning speeds up synthesis of porous materials

Machine-learning techniques used to accurately predict

# Extensions of Word Embeddings

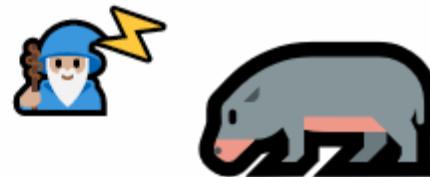
# Sub-word embeddings

Word embeddings work well if vocabulary is fixed

- so cannot deal with **new words** in test set

If we see a new word

- e.g. a made-up word like: **hippopotamification**  
*(act of magically turning someone or into a hippopotamus)*
- we don't have embedding for it!
- and must just ignore it, despite fact that we can guess its meaning of word from letters contained in it ... 😞



**Fasttext** ([2016 Bojanowski et al.](#))

- split words into fixed-length character sequences
- learns embeddings for character n-grams
- combines the embeddings to form words

Advantage: deals nicely with **morphologically** related terms, so:

- “**believe**” and “**believing**” have similar representations
- as do “**rain**” and “**rainfall**”

Embeddings are cool.

=>

— Em  
— Embe  
— mbed  
— bedd  
— eddi  
— ddin  
— ding  
— ings  
— \_are  
— \_coo  
— cool  
— ool.

# Conclusions

# Conclusions

Language models (LM):

- predict the next word in a sequence and can be used to generate text
- Markov models don't scale to long enough n-grams
- Perplexity is used to measure LM performance

Word embeddings

- provide distributed representation of words as dense vectors in high dimensional space
- capture the syntax and semantics of words from their usage
- solve the limited generalization problem of n-gram language models