

COMP34711

Week 2

(Simple)

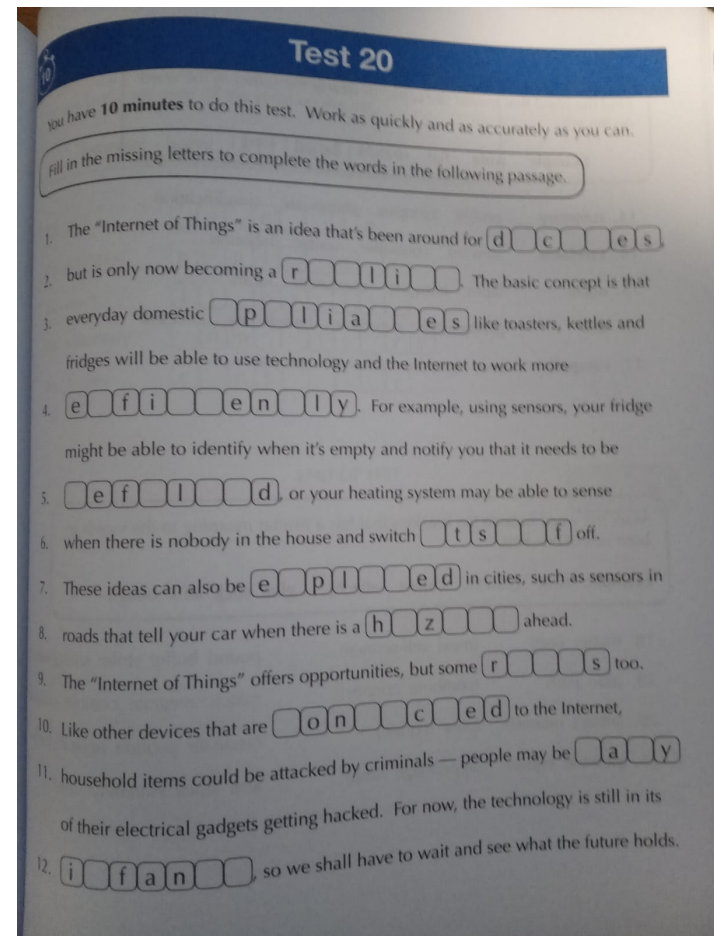
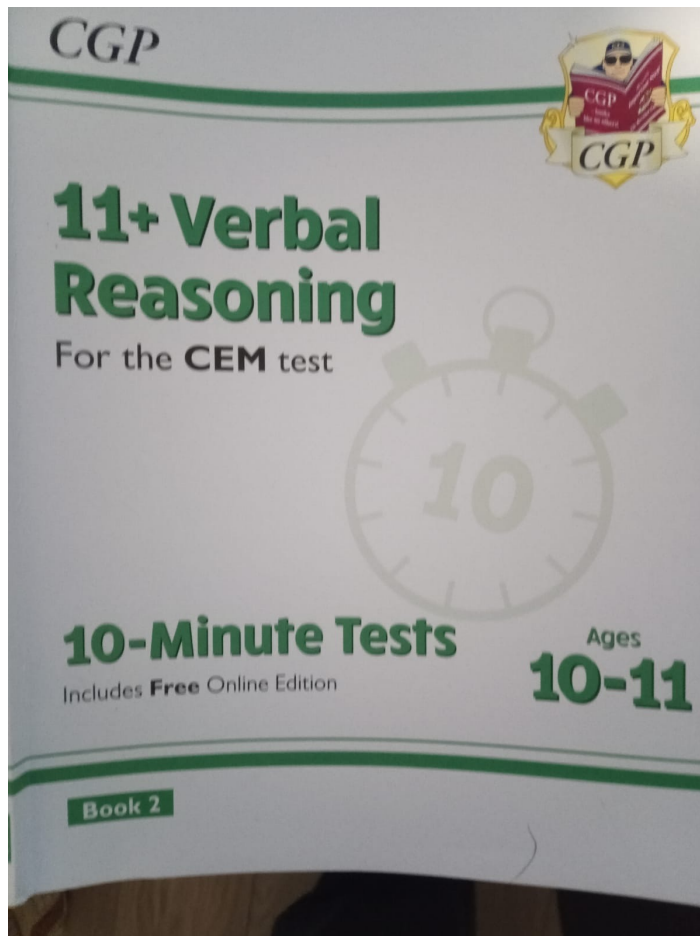
Language Models and Representations

Goran Nenadic

“Learning” language

The “Internet of Things” is an idea that’s been around for _____, but is only now becoming a _____. The basic concept is that everyday domestic _____ like toasters, kettles and fridges will be able to use technology and the Internet to work more _____. For example, using sensors, your fridge might be able to identify when it’s empty and notify you that it needs to be _____, or your heating system may be able to sense when there is nobody in the house and switch _____ off. These ideas can also be _____ in cities, such as sensors in roads that tell your car when there is a _____ ahead. The “Internet of Things” offers opportunities, but some _____ too.

“Learning” language



Language modelling

- An (often probabilistic) language model = a function that assigns a probability over a piece of text so that ‘natural’ pieces have a larger probability

$P(\text{"the the the the the"})$ – pretty small

$P(\text{"the cat in the hat"})$ – pretty big

$P(\text{"the cat in the hat"}) > P(\text{"the hat in the cat"})$

$P(\text{"the cat in the hat"}) \text{ ?? } P(\text{"the hat on the cat"})$

Language models and representations

- *Model* = a “simplified”, abstract representation of something, often in a computational form
 - e.g. tossing a coin
 - e.g. probabilistic model of whether forecast
- How do we represent a piece of text (e.g. word, sentence, paragraph, document, corpus)?
[use ‘document’ as an umbrella term here]

Language models and representations

<https://huggingface.co/roberta-base>

<https://huggingface.co/bert-base-cased>

<https://huggingface.co/bert-base-chinese>

Mask token: <mask>

The goal of life is <mask>.

Compute

Computation time on cpu: cached



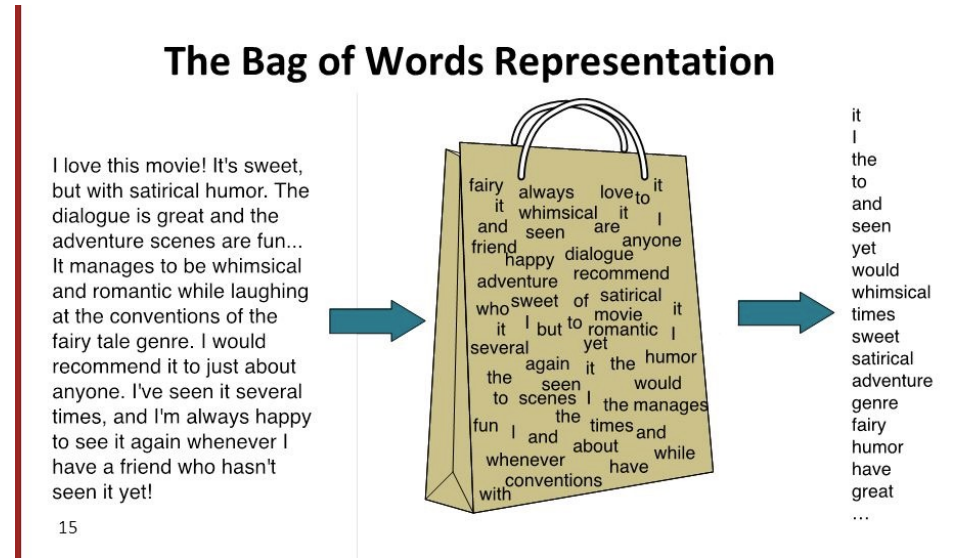
Paris is the <mask> of France.

Compute

Computation time on cpu: cached

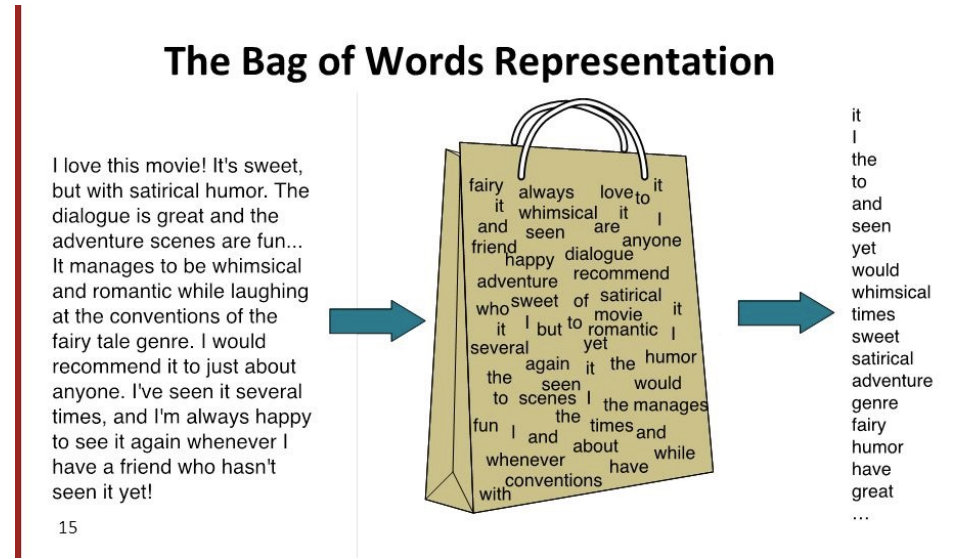


BoW representation



- Simplest model: “bag of words” (BoW)
 - Reduce each document into a bag of words
- Bag is a set with repetitions, but it could be used with binary indicators too.

BoW representation



- Note on terminology:
 - Bag of terms
 - Bag of tokens
 - Later, we will talk about bag of stems, etc.

BoW representation

- Questions:
 - Is the meaning lost without order?
John is quicker than Mary. vs *Mary is quicker than John.*
 - What about negations?
 - Are all words equally important?
 - Is the meaning lost without context? Ambiguities?
 - Would it work for all languages?
 - What about (important) combinations of words?
- BoW is however efficient
 - Used in many NLP systems

BoW representation

- **Use all words? Or skip some?**
 - How about so-called **stop-words** (*the, and, of, it, ...*)?
 - Most frequent ones?
- **Do we want to “rank” them?**
 - Give some weight to more important ones
 - How do you decide on importance/relevance?
- Let's see some linguistic hypotheses and theories

Frequency of words

- Some words are very frequent
 - e.g. “the”, “of”, “to”
- Many words are less frequent
 - e.g. “bazinga”
 - ~50% terms appears once
- Zipf’s Law
 - frequency of any word in a given collection is inversely proportional to its rank in the frequency table.



George Zipf
1902-1950

Zipf's law

Wikipedia abstracts

→ 3.5M En abstracts

$$r \times P_r \cong \text{const} \rightarrow$$

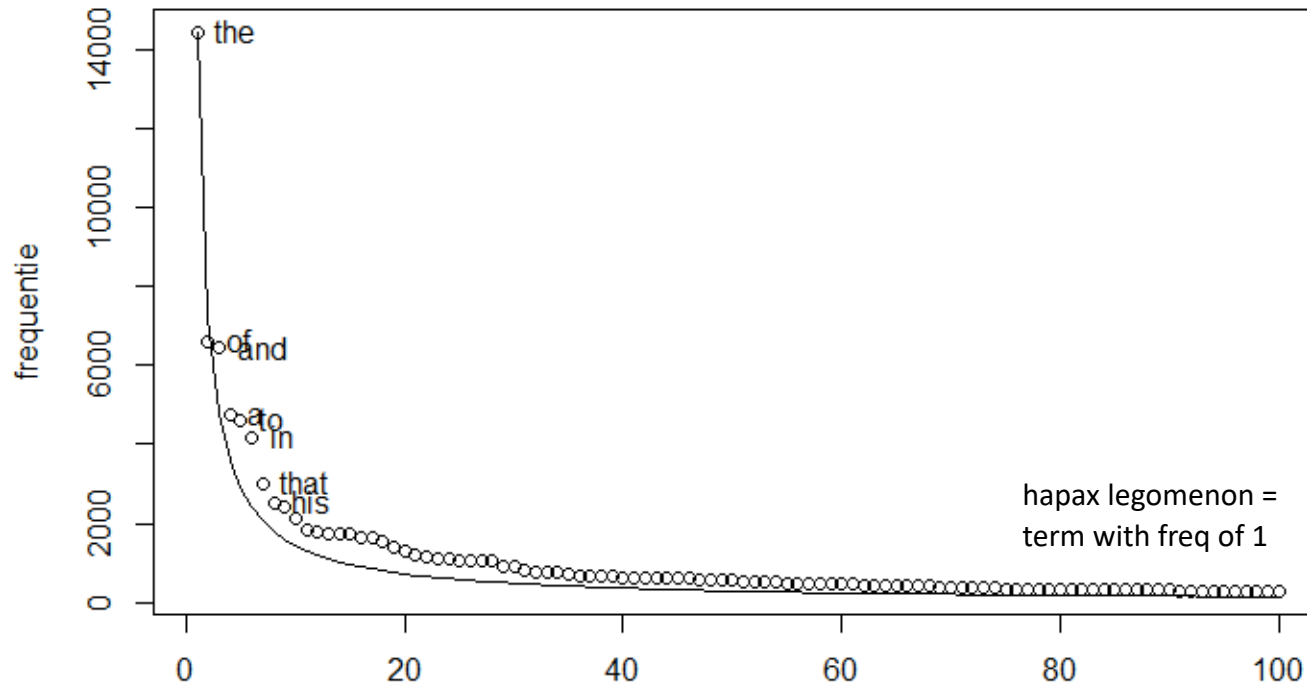
$$r \times \text{freq}_r \cong \text{const}$$

Term	Rank	Frequency	r x freq
the	1	5,134,790	5,134,790
of	2	3,102,474	6,204,948
in	3	2,607,875	7,823,625
a	4	2,492,328	9,969,312
is	5	2,181,502	10,907,510
and	6	1,962,326	11,773,956
was	7	1,159,088	8,113,616
to	8	1,088,396	8,707,168
by	9	766,656	6,899,904
an	10	566,970	5,669,700
it	11	557,492	6,132,412
for	13	493,374	5,970,456
as	14	480,277	6,413,862
on	15	471,544	6,723,878
from	16	412,785	7,073,160

Think about rare words:

- spelling errors
- names
- emails
- codes

Zipf's law

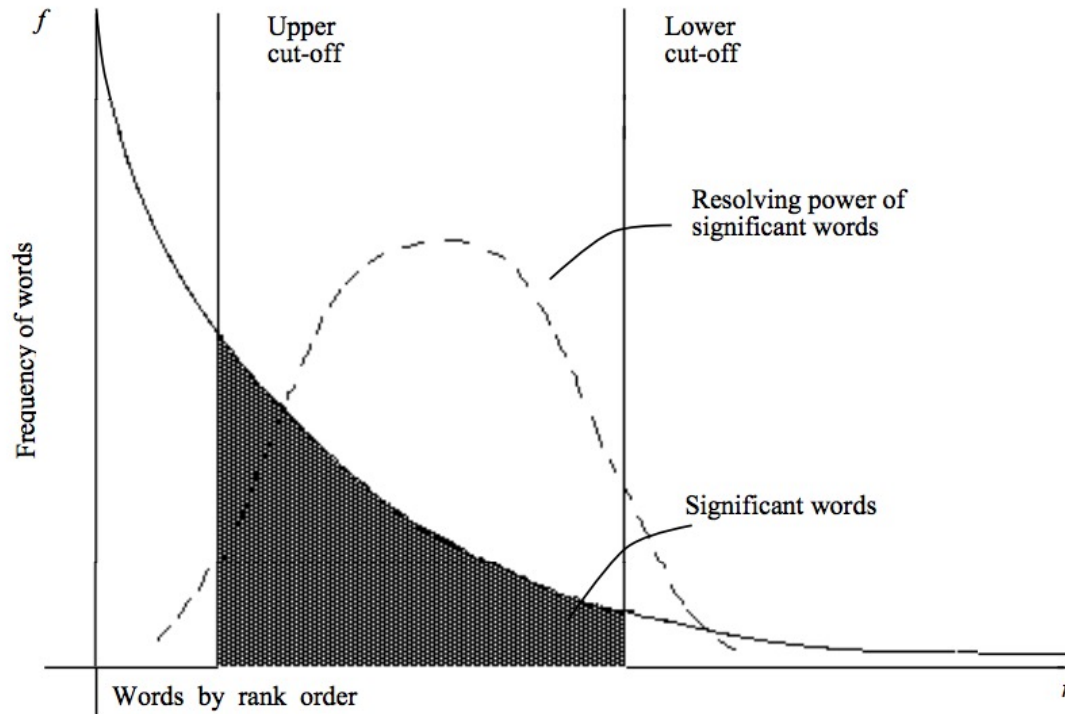


$$r \times Pr \cong const$$

r = rank of term according to frequency

Pr = probability of appearance of term

Luhn's hypothesis



The words exceeding the upper cut-off are considered to be common.

Those below the lower cut-off are rare, and therefore not contributing significantly to the content of the article.



Hans Peter Luhn
1896-1964

<http://www.dcs.gla.ac.uk/Keith/pdf/Chapter2.pdf>

LUHN, H.P., 'The automatic creation of literature abstracts', IBM Journal of Research and Development, 2, 159-165 (1958).

Removing stop words

- Highly frequently occurring words have low distinguishing power for representing documents
 - Top 30 words account for ~30% of mentions
 - Closed class function words (**the, and, of, it, ...**)
- Therefore, such words are often filtered out
- **Stop-word lists**
 - E.g. <https://gist.github.com/sebleier/554280>
 - Different lists for different applications, domains, etc.



[But not every application would benefit from removing stop words]

Vector representation

- Vector representation is a way to implement the BoW model.
- How to represent documents as vectors?
- Vocabulary V = set of terms left after pre-processing (tokenization, stop-word removal, case folding, etc.)
- Each document is represented as a $|V|$ -dimensional vector:

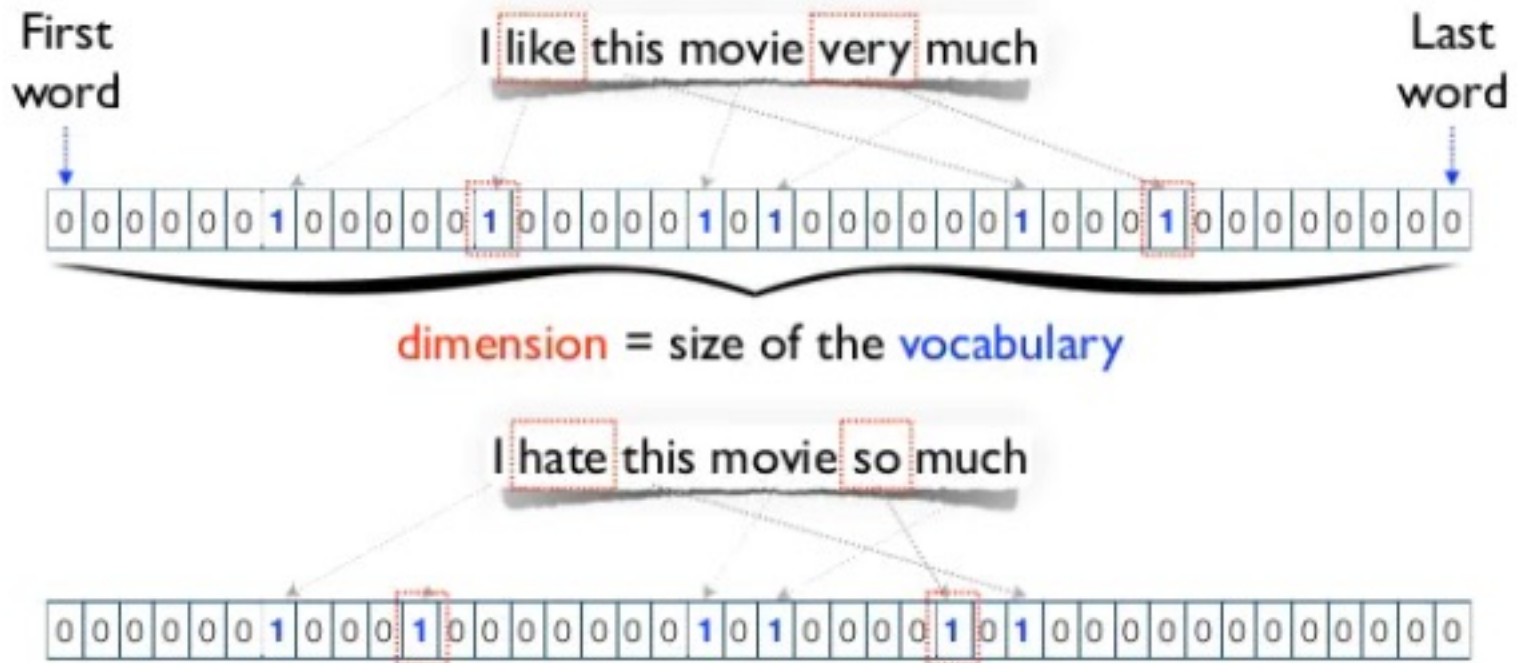
$$d = [w_1, w_2, \dots, w_{|V|}]$$

w_i is a weight of term i in the vocabulary

- Very high-dimensional representation
 - Millions of dimensions
 - These are very **sparse** vectors - most entries are zero

Weights

- Represent terms as an **incidence** (0, 1)



Weights

- Would a **frequency** vector be better?
 - Intuitively: if a term occurs many times in a document, is that word more important?
 - But:
 - Almost all documents will have (many) determiners?
 - Rare terms are more discriminative than frequent terms
 - Documents have different lengths: will clearly affect the counts
- Raw term frequency is not what we want:
 - E.g. term that appears 10 times is more important than a term with 1 occurrence *but not 10 times!*

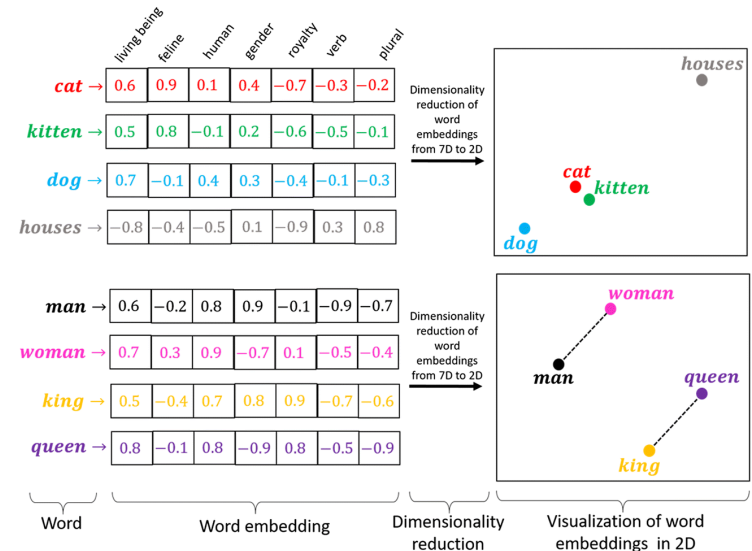
Importance does not increase proportionally with term frequency

Weights

- We will consider different weights when we look into Information Retrieval next week
 - E.g. *inverse document frequency* (in how many documents this word appears)
- Note again: these are very sparse models. Can we get more **dense** models to represent words and thus documents?
 - We will look later at *word embeddings*.

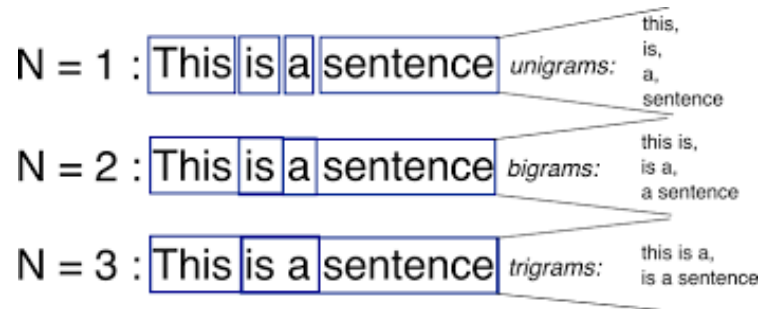
$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents



Combinations of words?

- Why only single words? Maybe a bag of n-grams?
 - Can n-grams represent some context?
- Reminder **n-grams**
 - a contiguous sequence of *n words* from a given sequence of text
 - 1-gram = “unigram”; “bigram” = 2-gram; “trigram” = 3-gram, etc.
 - Frequency of longer n-grams is often quite low



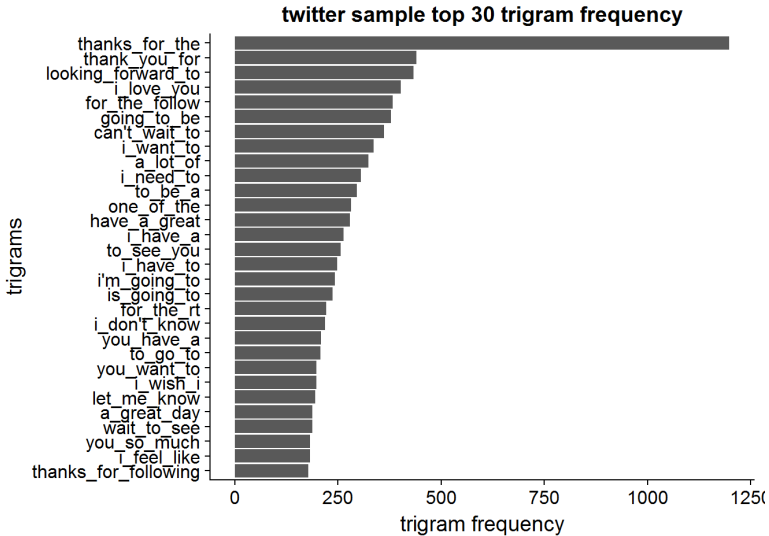
- VSM representation with n-grams?

Combinations of words

3-4-grams

word

Word	↓ Count	Word	↓ Count
1 one of the	303,265 ...	11 to be a	100,285 ...
2 as well as	291,495 ...	12 a number of	98,337 ...
3 part of the	151,966 ...	13 is one of	90,869 ...
4 I do n't	148,894 ...	14 of the most	90,028 ...
5 a lot of	145,674 ...	15 out of the	87,383 ...
6 the United States	144,719 ...	16 the use of	82,268 ...
7 be able to	132,872 ...	17 you want to	81,888 ...
8 in order to	128,631 ...	18 there is a	80,923 ...
9 some of the	119,133 ...	19 the number of	77,974 ...
10 the end of	118,539 ...	20 end of the	77,110 ...



Trigrams	
Frequency	Token
42030	THE U. S.
27260	IN NINETEEN EIGHTY
24165	CENTS A SHARE
18233	NINETEEN EIGHTY SIX
16786	NINETEEN EIGHTY SEVEN
15316	FIVE MILLION DOLLARS
14943	MILLION DOLLARS OR
14517	MILLION DOLLARS IN
12327	IN NEW YORK
11981	A YEAR EARLIER

Probabilistic language models

- Other ways to model text data?
- There is a lot about uncertainty and probability in how we use languages
- For example, is word 'running' more likely to follow word 'rabbit' than word 'flying'
I saw a rabbit running yesterday.
I saw a rabbit flying yesterday.
- **How do we know that?**

Probabilistic language models

$$p(\text{"I saw a rabbit running yesterday"}) = ?$$

> $p(\text{"I saw a rabbit flying yesterday"}) = ?$

$$p(\text{"I} \rightarrow \text{saw} \rightarrow \text{a} \rightarrow \text{girl} \rightarrow \text{running} \rightarrow \text{yesterday}})$$

> $p(\text{"} \neq \text{a} \rightarrow \text{girl} \rightarrow \text{saw} \rightarrow \text{I} \rightarrow \text{running} \rightarrow \text{yesterday}})$

A **language model** (LM) measures the probability of natural language utterances, giving higher scores to those that are more common, more grammatical, more “natural”.

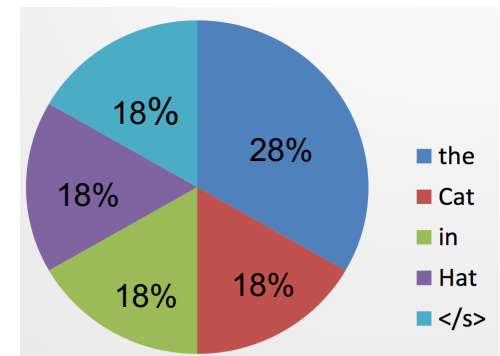
Unigram Language Model

Consider the sequence $S = w_1 w_2 \dots w_n$ as a sequence of independent unigrams (single tokens)

$$P(S) = P(w_1) * P(w_2) * \dots * P(w_n) = \prod_i P(w_i)$$

$$\begin{aligned} & p(\text{"the cat in the hat </s>"}) \\ &= p(\text{"the"}) p(\text{"cat"}) p(\text{"in"}) p(\text{"the"}) p(\text{"hat"}) p(\text{"</s>"}) \\ &= 0.28 \times 0.18 \times 0.18 \times 0.28 \times 0.18 \times 0.18 \end{aligned}$$

$$p(\text{"the hat in the cat </s>"}) = ?$$



Unigram Language Models

- Note: we assumed here that words appeared **independently (and identically distributed, iid)** of each other
 - That can't be true!
 - But it works pretty well (for some applications)
- How do we find these probabilities?
 - Use **corpus frequencies** as relative counts – count how many times they appear in a real-world dataset.
- Parameters of unigram LMs = word probabilities estimated from a corpus
 - How many?

Unigram Language Models

- How to deal with unseen (**out-of-vocabulary, OOV**) words that may appear in a sentence we want to estimate the probability for?
 - E.g. misspellings or new words
- Go with $p(OOV) = 0$? Harsh?
- Or somehow to estimate these probabilities
 - Same for all OOV words?
 - E.g. **add-one smoothing**
 - Other smoothings – later.

$$P(w) = \frac{\#w + 1}{\sum_{w' \in D} (\#w' + 1)}$$

Unigram Language Models

- Key issue: complete loss of order information with the unigram model

$$p(\text{"the the the the the </s>"}) > p(\text{"the cat in the hat </s>"})$$

$$p(\text{"cat the hat in the </s>"}) = p(\text{"the cat in the hat </s>"})$$

- How to add some sequencing into modelling?
- Chain rule

Chain Rule

- Chain rule in probability theory:
 - Calculate the *joint distribution* of a set of random variables using only conditional probabilities.

e.g. *S = the cat in the hat*

$$\begin{aligned}
 p(S) = & p(\text{the}) \cdot p(\text{cat}|\text{the}) \cdot p(\text{in}|\text{the cat}) \cdot \\
 & p(\text{the}|\text{the cat in}) \cdot p(\text{hat}|\text{the cat in the})
 \end{aligned}$$

So, for utterance: $w_1 w_2 w_3 \dots w_L$

$$\begin{aligned}
 & p(w_1, w_2, w_3, \dots, w_L) \\
 = & p(w_1) p(w_2 | w_1) p(w_3 | w_1, w_2) p(w_4 | w_1, w_2, w_3) \cdots p(w_L | w_1, w_2, \dots, w_{L-1}) \\
 = & \prod_{k=1}^L p(w_k | w_1, w_2, \dots, w_{k-1})
 \end{aligned}$$

Chain Rule

- Approximation: take (only) a fixed number of words before the current one

$$p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1})$$

- Unigram model (N=0)

$$p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1}) = \prod_{k=1}^L p(w_k)$$

- Bigram model (N=1)

$$p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1}) = \prod_{k=1}^L p(w_k | w_{k-1})$$

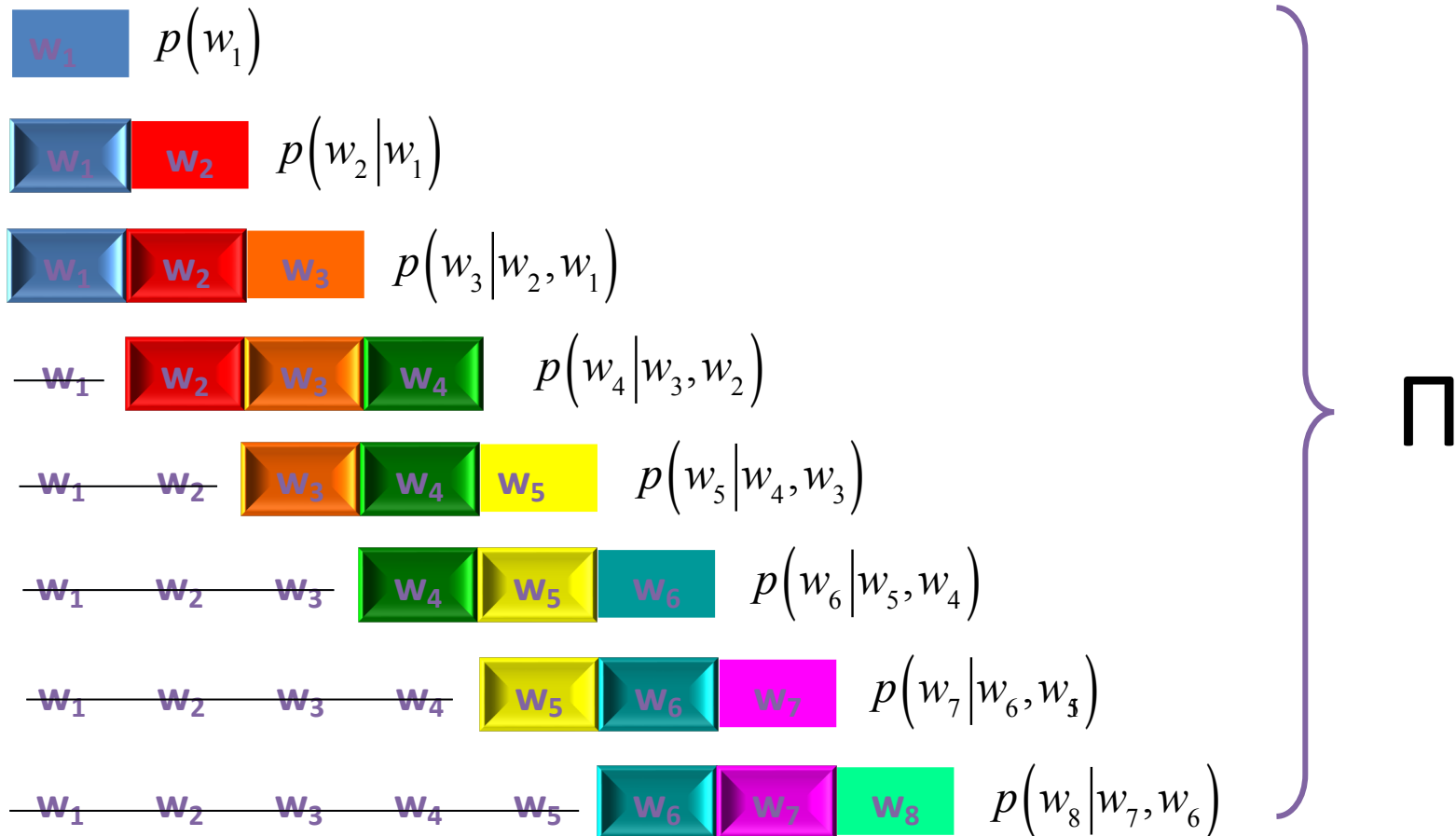
Bigram Model

$$p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1}) = \prod_{k=1}^L p(w_k | w_{k-1})$$

$$\begin{aligned} & p(\text{"the cat in the hat } \langle /s \rangle \text{"}) \\ &= p(\text{"the"}) p(\text{"cat"} | \text{"the"}) p(\text{"in"} | \text{"cat"}) p(\text{"the"} | \text{"in"}) p(\text{"hat"} | \text{"the"}) p(\text{"\langle /s \rangle"} | \text{"hat"}) \end{aligned}$$

Bi-gram model corresponds to a Markov chain. It assumes that a state depends only on its previous state not on the other events that occurred before it.

Trigram Model



N-gram Models

- N-gram model assumes

$$p(w_k | w_{1:k-1}) = p(w_k | w_{k-N+1:k-1})$$

We only care about **N-1** previous “states” (words) of w_k

For example: If we focus on the 10th word in a sequence:

... w_6 w_7 w_8 w_9 **w_{10}** w_{11} w_{12} w_{13} ...

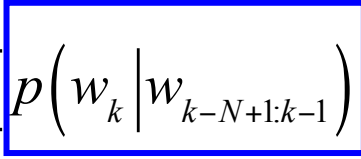
$$N=4: \quad p(w_{10} | w_{1:9}) = p(w_{10} | w_{7:9})$$

$$N=3: \quad p(w_{10} | w_{1:9}) = p(w_{10} | w_{8:9}) \rightarrow \text{trigram model}$$

$$N=2: \quad p(w_{10} | w_{1:9}) = p(w_{10} | w_9) \rightarrow \text{bigram model}$$

$$N=1: \quad p(w_{10} | w_{1:9}) = p(w_{10}) \rightarrow \text{unigram model}$$

N-gram Probabilities

$$\text{N-gram model: } p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1})$$


- How to estimate the *N-gram conditional probability* using a **training corpus**
 - ❑ Based on word counts.
 - ❑ Using a statistical model, e.g., hidden Markov model and Gaussian mixture model.
 - ❑ Using a machine learning model, e.g., neural network.

Count-based Estimation

- Estimate N-gram probability based on counts.

For the bigram case:

$$p(w_k | w_{k-1}) = \frac{\text{count}("w_{k-1}w_k")}{\sum_w \text{count}("w_{k-1}w")} = \frac{\text{count}("w_{k-1}w_k")}{\text{count}("w_{k-1}")}$$

denoted by $\frac{C(w_{k-1:k})}{C(w_{k-1})}$

Generalise to N-gram case

$$p(w_k | w_{k-N+1:k-1}) = \frac{C(w_{k-N+1:k})}{C(w_{k-N+1:k-1})}$$

Count-based Estimation

■ Example:

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P(I | <s>) = 2/3 = 0.67$$

$$P(\text{Sam} | <s>) = 1/3 = 0.33$$

$$P(\text{am} | I) = 2/3 = 0.67$$

$$P(</s> | \text{Sam}) =$$

$$P(\text{Sam} | \text{am}) =$$

$$P(\text{do} | I) =$$

$$P(\text{not} | \text{do}) =$$

$$P(\text{eggs} | \text{green}) =$$

Count-based Estimation

■ Example:

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P(I | <s>) = 2/3 = 0.67$$

$$P(\text{Sam} | <s>) = 1/3 = 0.33$$

$$P(\text{am} | I) = 2/3 = 0.67$$

$$P(</s> | \text{Sam}) = 1/2 = 0.5$$

$$P(\text{Sam} | \text{am}) = 1/2 = 0.5$$

$$P(\text{do} | I) = 1/3 = 0.33$$

$$P(\text{not} | \text{do}) = 1/1 = 1$$

$$P(\text{eggs} | \text{green}) = 1/1 = 1$$

Count-based Estimation

- Get raw bigram frequencies from a larger corpus:

count(want chinese) = 6; count(to spend) = 211

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

From 9220 sentences
 (but can/should get billion-word corpus)

Count-based Estimation

- Normalise with unigram counts:

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

N-gram Probabilities

$$\text{N-gram model: } p(w_{1:L}) = \prod_{k=1}^L p(w_k | w_{k-N+1:k-1})$$

■ Example:

$$p(\text{"I want to eat Chinese food"}) = p(I) * p(\text{want} | I) * \dots$$

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Count-based Estimation

- Examples – similar process for N-grams

$$p(\text{"want"}|\text{"I"}) = \frac{C(\text{"I want"})}{C(\text{"I"})} \text{ (bi-gram example)}$$

$$p(\text{"machine"}|\text{"support vector"}) = \frac{C(\text{"support vector machine"})}{C(\text{"support vector"})} \text{ (tri-gram example)}$$

Homework and further reading

- Do your homework (estimated time ~ 2 hours)
 - Practice tokenisation, **sentence segmentation**, n-gram extraction, n-gram modelling
- Ch3 in: Jurafsky, Martin: *Speech and language processing: an introduction to natural language processing, computational linguistics, and speech recognition*
 - <https://web.stanford.edu/~jurafsky/slp3/>
- Ch3 in: Bird, Klein, Loper: *Natural Language Processing with Python: analyzing text with the Natural Language Toolkit*
 - <http://www.nltk.org/book/>