

THE UNIVERSITY OF YORK

WORD EMBEDDINGS

A DISSERTATION SUBMITTED TO  
THE DEPARTMENT OF MATHEMATICS  
IN CANDIDACY FOR THE DEGREE OF  
MATHEMATICS BSC

BY  
JOEL STROUTS

YORK, UK

2021

## ABSTRACT

In this report we attempt to provide a complete introduction to the core theories of word embeddings (numerical vector representations of words as they appear in text) sufficient to enable further reading of more advanced topics. This introduction covers briefly the linguistic context, theoretical development, most impactful implementations, and practical considerations of application of word embedding methods. In the first section we discuss the nature of the words, in the second we describe the cross disciplinary events which lead to development the key mathematical techniques for producing embeddings and explain their operation, and finally we discuss practical considerations of these algorithm's applications. As the primary relevance of these methods within research is their practical application we take a broader contextual approach opposed to a purely mathematical one and focus on the developments and methods most influential in the sphere of application. We assume familiarity with core mathematics at an undergraduate level as prerequisite for our descriptions of the selected implementation methods but otherwise assume no prior knowledge of the subject.

# TABLE OF CONTENTS

ABSTRACT . . . . .	ii
INTRODUCTION . . . . .	1
1 WORDS . . . . .	2
1.1 Units of Meaning? . . . . .	2
1.2 Features of Written Language . . . . .	4
1.3 Objects of Representation . . . . .	5
2 WORD REPRESENTATIONS . . . . .	7
2.1 Preliminaries . . . . .	7
2.2 Count Models . . . . .	11
2.3 Count Models . . . . .	12
2.3.1 History . . . . .	12
2.3.2 Term-Document & Term-Term Matrices . . . . .	13
2.3.3 Term Weighting . . . . .	14
2.3.4 Smoothing . . . . .	14
2.3.5 Singular Value Decomposition . . . . .	14
2.4 Predict Models . . . . .	14
2.4.1 Early Developments . . . . .	14
2.4.2 <code>word2vec</code> : Skip-Gram with Negative Sampling . . . . .	14
2.4.3 <code>word2vec</code> : Continuous Bag of Words . . . . .	14
2.5 Measuring Similarity . . . . .	14
2.6 Comparison . . . . .	14
2.7 Other Approaches . . . . .	14
3 APPLICATION . . . . .	15
3.1 Preprocessing . . . . .	15
3.2 Hyper-parameters . . . . .	15
3.3 Bias . . . . .	15
FURTHER READING . . . . .	16

## INTRODUCTION

Word embeddings are numerical representations of words with the property that qualitative semantic observations like “the word *car* is more similar to the word *bicycle* than to the word *think*” are encoded by the differences in each word’s representation. Words are represented as points in a high dimensional ‘semantic space’ with relative positions corresponding to the ‘relatedness’ of the given words. In this way, the statement above could be restated as follows:

$$d(\text{CAR}, \text{BICYCLE}) > d(\text{CAR}, \text{THINK})$$

where  $d$  is a distance function, and  $ABC$  is the embedding of the word  $abc$ .

All methods for generating word embeddings apply some form of the \*distributional hypothesis\*: that assertion that words which appear in similar contexts have similar meanings. This hypothesis was famously summarised by J.R. Firth in the statement;

*“You will know a word by the company it keeps”*

The first techniques for producing embeddings were developed in the field of information retrieval (IR) and are based on word-context frequency matrices. These approaches are now referred to as count models, in contrast to more recent methods for producing word embeddings using neural networks which are referred to as predict models.

In the first chapter we discuss the object of representation: words, then in the second chapter we discuss the methods of representation: first count models and then predict models, and finally in the last chapter we discuss practical considerations before concluding with suggested further reading.

# CHAPTER 1

## WORDS

To make use of either word embedding algorithms and their outputs for some practical means an understanding of related linguistic concepts is not necessarily needed, however, we consider a grounding in these terms illuminating context for a more reasoned application and understanding of these methods. In particular, understanding these concepts is key to understanding some of the shortcomings of the models discussed in this report and the modifications made by later models to address those.

### 1.1 Units of Meaning?

*“To many people the most obvious feature of a language is that it consists of words” -Halliday, *Lexicology and Corpus Linguistics* 2004*

What is a word? The ‘Word’ workshop, held at the international research centre for linguistic typology 12-16 August 2000 consisted of 16 presentations discussing the answer to this question. Ten of those were published in the book *word: a cross linguistic typology*. The difficulty of the question is summarised in the book’s conclusion as follows:

*“The problem of the word has worried general linguists for the best part of a century. In investigating any language, one can hardly fail to make divisions between units that are word-like by at least some of the relevant criteria. But these units may be simple or both long and complex, and other criteria may establish other units. It is therefore natural to ask if ‘words’ are universal, or what properties might define them as such.”*

Consider *so called* filler-words like “uh” or “um”, or compound forms like “rock ‘n’ roll” and “New York”. These are examples of language elements that are sometimes treated as words and sometimes not; it is not clear where the boundary should be drawn. Even agreed upon boundaries change with time; the word “tomorrow”,

listed in Johnson’s 1755 dictionary of the English language as two distinct components “to”, and “morrow” yet these parts have merged and are considered a single indivisible word. Halliday asks “How do we decide about sequences like lunchtime (lunch-time, lunch time), dinner-time, breakfast time? How many words in isn’t, pick-me-up, CD?”.

The understanding of ‘word’ formed within the English language cannot simply be applied beyond that scope either, (Word 2002) notes “The idea of ‘word’ as a unit of language was developed for the familiar languages of Europe”. Languages written without spaces, like Chinese, Japanese, and Thai present one particular example of the challenges in applying the ‘word’ concept beyond this scope: a single sentence may be divided up into different word-like units depending on the approach taken, (chen et al 2017) gives the following example:

The sentence “Yao Ming reaches the finals” (姚明进入总决赛)

May be divided up to yield either three words or five using the two segmentations methods; ‘Chinese Treebank’ segmentation, or ‘Peking University’ segmentation, respectively:

姚明      进入      总决赛

YaoMing reaches finals

姚    明      进入      总            决赛

Yao Ming reaches overall finals

Different languages present different difficulties -In the case of Arabic, expressions formed by adding consecutive suffixes to a single root word (like establish → establishment → establishmentarianism) are a more essential part of the language,

to the extent that complete sentences can be expressed as single word-units. As a result, ‘words’ in the sense of ‘things occurring between spaces’ describe a unit of meaning larger and more complex than the term describes in English, meaning representations of Arabic words defined this way describe something quite different to the English analogue.

The term **lexical item** describes ‘units’ of language more broadly; in Chinese they may be individual characters or groups with a conventional meaning. in Arabic they may be the set of roots and suffixes In English they may be the words, inflections, and multi-word phrases used as single parts.

## 1.2 Features of Written Language

While the linguistic & philosophical discussion of ‘word’ informs the theory of computational linguistics, in our context the question of “what is a word?” is realised as a question of text-preprocessing. It is assumed that the input to the embedding process is a written text and it is the task of text preprocessing to extract ‘words’ from that text. The implementation used is more motivated by the application requirements & task performance than linguistic theory.

Embedding algorithms are not discriminatory about what you classify as ‘a word’ in their input (in fact they can be applied to forms of structured data other than just language, as we will discuss in the final section of this report) and the decision about how to make these divisions will depend on the task at hand. If the intended application of the embeddings is within the field of chemistry, then a compound name like (NH4)2SO4 should absolutely be treated as a word, while in other contexts it is possible that such an irregular form may be discounted without negative consequences. The input material impacts these decisions too; text in newspapers has a narrower scope of common word-like forms than a medium like tweets, and text preprocessing decisions should reflect this.

The process of extracting word-units from a text is called **tokenization** and the resulting units are called **word-tokens**. Word tokens with the same form (like *red*

and *red*) are said to be instances of same word-class (or word-type). For example, the sentence;

*The limits of my language means the limits of my world* - Wittgenstein,  
1922

Contains 11 word-tokens, and 7 word-classes (assuming the sentence is tokenized by identifying word boundaries at both ends and each blank space).

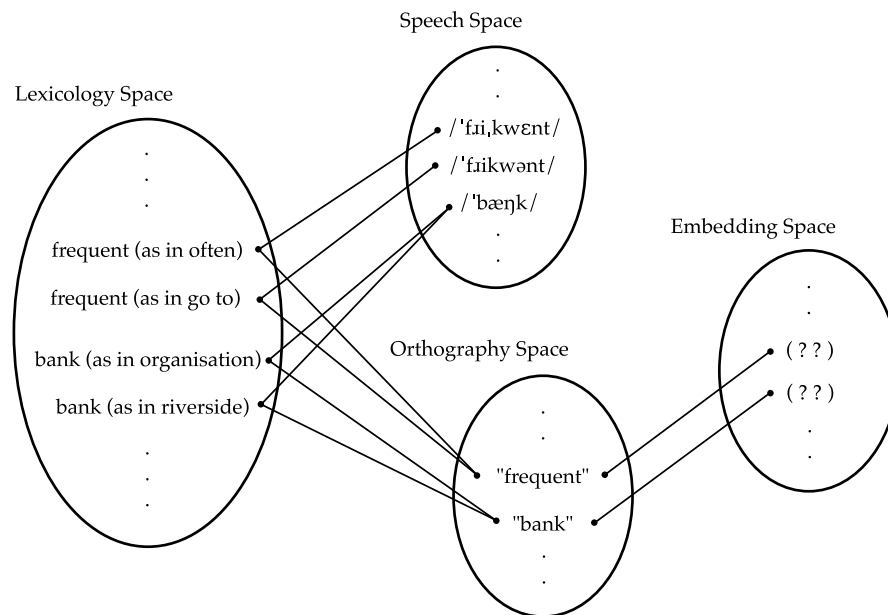
A detailed discussion of tokenization & text preprocessing in general is beyond the scope of this report but some practical considerations can be found in ??

### 1.3 Objects of Representation

**Lexicology** is the study of words, and **semantics** is the study of meaning. The ‘embedding’ in ‘word embeddings’ derives from the mathematical sense ‘to embed in a space’ (for example, to embed a graph in a coordinate space). Word embedding algorithms are methods for embedding *lexicological objects* in *semantic space*, i.e., words a space where position relates to meaning somehow). This is achieved by statistical analysis of large quantities of written language (by extraction of word-tokens).

The conventions of written language, called **orthography**, make word-tokens an imperfect representation of ‘words’ in the lexicological sense. In particular, one problem is that distinct words do not always have distinct realisations in written language. Orthographic forms with this property are called **polysemus** (poly as in *many* and semus as in *meaning*, from the same root as semantics). Alternatively, the distinct words sharing this same form are referred to as **homonyms** (having the same name). The words *bank* as in money holding organisation and *bank* as in riverside are often used to illustrate this property. Indeed, some words, like *frequent* as in often, and *frequent* as in attend are examples of words distinct in spoken language despite being identical in written language; such words are said to be **homographs** (same orthography, different names).



Figure 1.1: The *problem of polysemy* illustrated

Embedding algorithms which generate their representations without taking this *problem of polysemy* into account produce representations which conflate the multiple meanings of a word-form into a single point in the embedding space.

## CHAPTER 2

### WORD REPRESENTATIONS

#### 2.1 Preliminaries

We begin by defining terms already encountered in mathematical terms.

##### **Definition 2.1.1 – Orthography**

An orthography is a finite set of symbols.

*Example 2.1.1.* The sets:

$$\mathcal{O}_L = \{c \mid c \text{ is a symbol that can be copied and pasted the pdf of this report}\} \quad (2.1)$$

$$\mathcal{O}_P = \{c \mid c \text{ is the inscription of a regular polygon with eight sides or less}\} \quad (2.2)$$

$$\mathcal{O}_K = \{c \mid c \text{ is a typable character marked on Figure 2.1}\} \quad (2.3)$$

are all examples of orthographies.

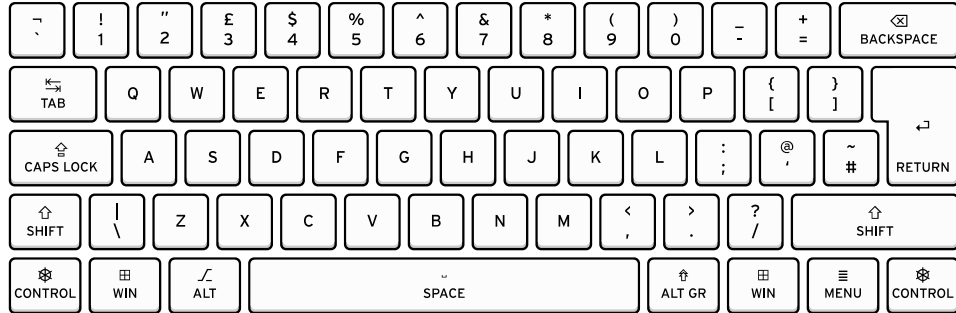


Figure 2.1: A standard uk keyboard layout

Note, the term ‘orthography’ in linguistics refers to the conventions of written language in general and is distinct from the use of ‘an orthography’ we use here.

##### **Definition 2.1.2 – Document**

A document  $d$  is a finite sequence  $c_0, c_1, \dots, c_n$  ( $c \in \mathcal{O}$ ), where  $\mathcal{O}$  is an orthography.

*Example 2.1.2 (This Report).* The sequence of symbols obtained by copy pasting this entire report into a text editor (ordered by starting with the cursor at the top of the file and moving the cursor with the arrow keys one position to the right at a time -including newlines) is a document using the orthography 2.1.

*Example 2.1.3 (Shapes Document).* The sequence of shapes:  $\triangle\triangle\square\triangle\triangle\circ$  is a document using the orthography 2.2.

*Example 2.1.4 (Plato Documents).* Any of the works of Plato available at the Project Gutenberg website, copied and pasted from their plain text formats and using the same procedure described above, are documents using the orthography 2.3.

### **Definition 2.1.3 – Word-Token**

A word token  $w$  has the same attributes as a document; it is a finite sequence  $c_0, c_1, \dots, c_n$  ( $c \in \mathcal{O}$ ), where  $\mathcal{O}$  is an orthography. We only refer to sequences of symbols as word tokens when they are the output of a tokenizer, however.

### **Definition 2.1.4 – Tokenizer**

A tokenizer  $T$  is a function that maps a document to a sequence of word tokens;  $T(d) = w_1, w_2, \dots, w_n$ . The output is referred to as the 'tokenized' form of the document. We use the shorthand  $W_d$  to refer to the tokenized document  $T(d)$  when it is clear which tokenization has been used.<sup>1</sup>

*Example 2.1.5 (This Report Tokenized).* The document from example 2.1.2 if tokenized by treating everything occurring between two spaces (or between a space and the start or end of a line) yields a sequence of word tokens beginning “THE” “UNIVERSITY” “OF”.

*Example 2.1.6 (Shapes Document Tokenized).* The document from example 2.1.3

if tokenized using the function:

$T(d)$ : Group repeated symbols in the sequence  $d$  into word-tokens  
(in the order of occurrence).

produces the output: “ $\triangle \triangle$ ”, “ $\square$ ”, “ $\triangle \triangle$ ”, “ $\diamond$ ”

*Example 2.1.7 (Simple Text Tokenizer).* Given a document  $d$  using the orthography  $\mathcal{O}_K$  from example 2.1.1 apply the following algorithm:

```

input : A list of characters: document
output: A list of word tokens: tokens
1 tokens  $\leftarrow$  list();
2 next_token  $\leftarrow$  list();
3 for character in document do
4   if character is alphanumeric then
5     next_token  $\leftarrow$  append(next_token, character);
6   else if character is a space or a newline then
7     if next_token is not empty then
8       next_token  $\leftarrow$  convert_to_lowercase(next_token);
9       tokens  $\leftarrow$  append(tokens, next_token);
10      next_token  $\leftarrow$  list();
11    end
12  end
13 end

```

**Algorithm 1:** Simple Text Tokenizer

### Definition 2.1.5 – *Corpus*

A corpus  $C$  is a finite set of documents  $\{d_1, d_2, \dots, d_n\}$ .

*Example 2.1.8 (Plato Corpus).* All the of Plato works referred to in example 2.1.4, if converted to documents using the same procedure described there, considered together form a corpus of documents. We will refer to this corpus as  $C_p$ .

### Definition 2.1.6 – *Vocabulary*

The vocabulary of a tokenized document  $W_d$  is the set of unique word tokens in

the tokenization of that document;  $V(W_d) = \{w \mid w \in \text{the sequence } W_d\}$ . The vocabulary of a corpus is the union of the vocabularies of each document in that corpus;  $V(C) = V(d_1) \cup V(d_2) \cup \dots \cup V(d_n)$ .

*Example 2.1.9 (Shapes Document Vocabulary).* Using the same tokenization as before, the document from example 2.1.3 has the vocabulary:

$$\begin{aligned} V(\text{"}\triangle \triangle\text{"}, \text{"}\square\text{"}, \text{"}\triangle \triangle\text{"}, \text{"}\diamond\text{"}) \\ = \{\text{"}\triangle \triangle\text{"}, \text{"}\square\text{"}, \text{"}\diamond\text{"}\} \end{aligned}$$

### Definition 2.1.7 – # Operator

The # symbol is used to denote the number of elements of the object it precedes. Before sequences it counts the number of elements in the sequence:  $\#W_d = n$  (for  $W_d = w_1, \dots, w_n$ )<sup>2</sup>. Before sets it is a shorthand for the size of that set:  $\#V(d) = |V(d)|$  = the number of words in the vocabulary  $V(d)$ . Before a corpus it counts the number of documents in that corpus.

*Example 2.1.10 (Shapes Document Sizes).* The word count of the tokenized shapes document

$$W_d = \text{"}\triangle \triangle\text{"}, \text{"}\square\text{"}, \text{"}\triangle \triangle\text{"}, \text{"}\diamond\text{"}$$

is  $\#W_d = 4$  and the vocabulary size is  $\#V(W_d) = 3$ .

*Example 2.1.11 (Plato Document Sizes).* Word token counts and vocabulary sizes for a selection Plato documents from 2.1.4, tokenized using the simple method:<sup>3</sup> have

*Example 2.1.12 (Plato Corpus Size).* The complete Plato corpus (2.1.8) consists of 25 documents and, each tokenized the same way as before, has the word count

---

<sup>3</sup> the code used to create these examples and others examples on the same corpus is available in the appendix

Attribute	Document ( $d$ )					
	laws	republic	cratylus	meno	apology	crito
$\#W_d$	140,406	118,283	23,883	12,716	11,392	5,332
$\#V(W_d)$	8100	8,105	2,963	1,493	1,789	1,030

Table 2.1: Number of Word Tokens & Vocabulary Sizes of a selection of Documents in the Plato Corpus (example 2.1.8)

$\sum_{d \in C_P} \#W_d = 670,936$  and the vocabulary size  $\#V(C_P) = 19,432$ .

## 2.2 Count Models

The embedding algorithms discussed in this report represent the confluence of two key ideas: first that the meaning of a word corresponds with the contexts in which it occurs, and second that the meaning of a word can be represented by a vector of numerical values.

This first idea is called the **distributional hypothesis**, was “widespread in linguistic theory in the 1950s” Jurafsky (2021), with key contributions to its conception contributed by Joos (1950) *Description of Language Design*, Harris (1954) *Distributional Structure*, and Firth (1957) *A synopsis of linguistic theory 1930-1955*.

To illustrate the distributional hypothesis, consider for example the possible meanings of “bank” discussed above in relation to polysemy: *bank* as in the organisation, and *bank* as in riverside. The first sense of the word is identified with neighboring word contexts like: account, manager, desk, etc. whereas the second sense is distinguished by context words like: river, grassy, marsh.

The way these distributional relationships are inferred by analysis of a **corpus** (a collection of reference texts) and synthesised in vector form differ according to the embedding algorithm used. The ‘count models’ developed in the 1960s and refined thereafter were the first to apply the distributional hypothesis to create word vectors, and remained the state of the art until a series of papers published by Milokov et al (2013) presented efficient methods for generating word vectors using

neural networks, along with a standard software implementation `word2vec` which then became ubiquitous.

We will discuss count models first as they lay the contextual foundation for the innovations of the later predict models.

## 2.3 Count Models

### 2.3.1 History

The adoption of distributional perspectives within linguistic spheres in the 1950s coincided with an emerging interest in mechanical and computational methods for document indexing and information retrieval (IR); this purpose was one of the first applications considered for modern computers. The question of quantifying document ‘aboutness’ and relatedness suggested, as viewed through the right mathematical lens, a notion of a ‘document space’ populated by ‘document vectors’ where relatedness could be quantified by a distance function on the space. From this the notion word vectors naturally followed. We summarise some key events in this lineage below.

In 1948 two chemists and leading advocates of punched cards, James W. Perry and G. Malcolm Dyson, petitioned IBM head, Thomas Watson, Sr to invest in research & development on “mechanical solutions for scientists’ specific information problems”, who selected Hans Peter Luhn for the task [cite: Rieder, B 2020].

This appointment resulted in “a remarkable series of papers” [cite: salton 1987] published by Luhn between 1957 & 1959, in particular: *A Statistical Approach to Mechanized Encoding and Searching of Literary Information*, Luhn (1957). The concepts discussed in this paper were prescient of future developments in the mechanical processing of documents but details of computation are omitted.

Maron & Kuhns (1960) provided details of a potential approach to the computation of such a system in the influential *On relevance, Probabilistic Indexing and Information Retrieval*. published around the same time. Crucially, the paper suggests the use of conditional probabilities to compute ‘closeness of index terms’. This

was an application of distributional thinking to quantify word similarity, but it did not specify a vector-space conception of the problem.

This final piece was supplied by Paul Switzer in *Vector Images in Document Retrieval* (1964), which references Maron & Kuhn’s paper and elaborates by suggesting the construction of a ‘term-term association matrix’ based on co-occurrence frequencies producing an ‘index term vector space’ in which documents are then located. It is the first paper to suggest that document similarity could be determined by projecting index terms and documents into a vector space.

Details of the resulting model and the most widely adopted refinements are the focus of this section.

### 2.3.2 Term-Document & Term-Term Matrices

#### **Definition 2.3.1 – Document Word Occurences**

The number of times a word-token  $w$  occurs in a tokenized document  $W_d$  is given by  $\text{count}(w, W_d)$ .

#### **Definition 2.3.2 – Term Document Matrix**

A term-document matrix  $X$ , (for corpus  $C$  and tokenizer  $T$ ) has columns corresponding to each corpus document, and rows corresponding to each word in the corpus vocabulary (hence having dimensions  $\#V(C) \times \#C$ ). The matrix entry



*2.3.3 Term Weighting*

*2.3.4 Smoothing*

*2.3.5 Singular Value Decomposition*

## **2.4 Predict Models**

*2.4.1 Early Developments*

*2.4.2 word2vec: Skip-Gram with Negative Sampling*

*2.4.3 word2vec: Continuous Bag of Words*

## **2.5 Measuring Similarity**

## **2.6 Comparison**

## **2.7 Other Approaches**

## **CHAPTER 3**

### **APPLICATION**

#### **3.1 Preprocessing**

#### **3.2 Hyper-parameters**

#### **3.3 Bias**

## FURTHER READING

