

IN5550: Neural Methods in
Natural Language Processing
Sub-lecture 5.1
Distributional and distributed

Andrey Kutuzov

University of Oslo

21 February 2023





- 1 Distributional and Distributed
 - Distributional hypothesis
 - Representing words with vectors



Bag-of-words problems

Simple Bag-of-Words approaches (without **representation learning**) do not take into account **semantic relationships between linguistic entities**.



Bag-of-words problems

Simple Bag-of-Words approaches (without **representation learning**) do not take into account **semantic relationships between linguistic entities**.

No way to detect semantic similarity between documents which **do not share words**:

- ▶ *The war was devastating for the region.*
- ▶ *This military conflict left the country in ruins.*



Bag-of-words problems

Simple Bag-of-Words approaches (without **representation learning**) do not take into account **semantic relationships between linguistic entities**.

No way to detect semantic similarity between documents which **do not share words**:

- ▶ *The war was devastating for the region.*
- ▶ *This military conflict left the country in ruins.*

It means we need more sophisticated **semantically-aware methods**.



Bag-of-words problems

Simple Bag-of-Words approaches (without **representation learning**) do not take into account **semantic relationships between linguistic entities**.

No way to detect semantic similarity between documents which **do not share words**:

- ▶ *The war was devastating for the region.*
- ▶ *This military conflict left the country in ruins.*

It means we need more sophisticated **semantically-aware methods**.

Like **distributional word embeddings**.



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*
- ▶ Neural language models **learn vector representations for words as a byproduct** of their training process.



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*
- ▶ Neural language models **learn vector representations for words as a byproduct** of their training process.
- ▶ These **representations are similar for semantically similar words.**



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*
- ▶ Neural language models **learn vector representations for words as a byproduct** of their training process.
- ▶ These **representations are similar for semantically similar words**.
- ▶ Good word embeddings from an auxiliary task:
 - ▶ Language models (LMs) are **trained on raw texts**, no manual annotation needed.



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*
- ▶ Neural language models **learn vector representations for words as a byproduct** of their training process.
- ▶ These **representations are similar for semantically similar words**.
- ▶ Good word embeddings from an auxiliary task:
 - ▶ Language models (LMs) are **trained on raw texts**, no manual annotation needed.
 - ▶ One can train an LM on the texts collected from the whole Internet.
 - ▶ Internet Archive
 - ▶ CommonCrawl
 - ▶ etc



Distant memory from the last lecture

- ▶ *'Generalizations: similar words get similar representations in the embedding layer'*
- ▶ Neural language models **learn vector representations for words as a byproduct** of their training process.
- ▶ These **representations are similar for semantically similar words**.
- ▶ Good word embeddings from an auxiliary task:
 - ▶ Language models (LMs) are **trained on raw texts**, no manual annotation needed.
 - ▶ One can train an LM on the texts collected from the whole Internet.
 - ▶ Internet Archive
 - ▶ CommonCrawl
 - ▶ etc

How come that we can get good word embeddings without any manually annotated data?



All of this week' sub-lectures in one slide

- ▶ **Vector space models of meaning** based on distributional information are not something new [Turney et al., 2010].



All of this week' sub-lectures in one slide

- ▶ **Vector space models of meaning** based on distributional information are not something new [Turney et al., 2010].
- ▶ But around 2011-2013, such representations trained using **machine learning** became extremely popular in NLP.



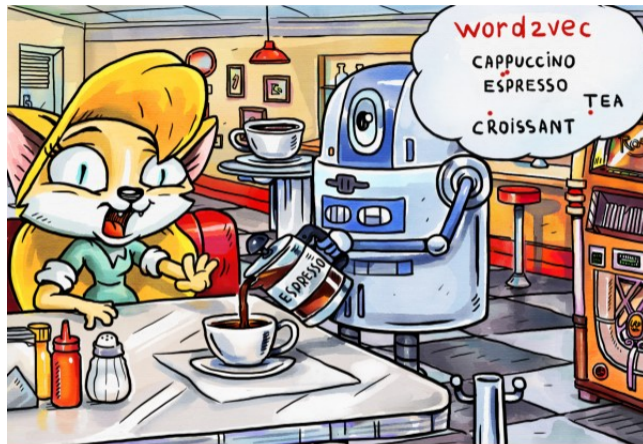
All of this week' sub-lectures in one slide

- ▶ **Vector space models of meaning** based on distributional information are not something new [Turney et al., 2010].
- ▶ But around 2011-2013, such representations trained using **machine learning** became extremely popular in NLP.
- ▶ Commonly used in research and large-scale industry projects (web search, opinion mining, tracing events, plagiarism detection, document collections management, etc.)



All of this week' sub-lectures in one slide

- ▶ **Vector space models of meaning** based on distributional information are not something new [Turney et al., 2010].
- ▶ But around 2011-2013, such representations trained using **machine learning** became extremely popular in NLP.
- ▶ Commonly used in research and large-scale industry projects (web search, opinion mining, tracing events, plagiarism detection, document collections management, etc.)
- ▶ All this is based on their ability to efficiently predict **semantic similarity** between linguistic entities (in particular, words).
- ▶ Semantic information is **distributed** across word vectors, making them non-interpretable.



- Espresso? But I ordered a cappuccino!
- Don't worry, the cosine distance between them is so small that they are almost the same thing.



OK, why does it work at all?



OK, why does it work at all?

Tiers of linguistic analysis



OK, why does it work at all?

Tiers of linguistic analysis

Computational approaches model various **tiers of language**:

- ▶ **graphematics** – how words are spelled,



OK, why does it work at all?

Tiers of linguistic analysis

Computational approaches model various **tiers of language**:

- ▶ **graphematics** – how words are spelled,
- ▶ **phonetics** – how words are pronounced,



OK, why does it work at all?

Tiers of linguistic analysis

Computational approaches model various **tiers of language**:

- ▶ **graphematics** – how words are spelled,
- ▶ **phonetics** – how words are pronounced,
- ▶ **morphology** – how words inflect,



OK, why does it work at all?

Tiers of linguistic analysis

Computational approaches model various **tiers of language**:

- ▶ **graphematics** – how words are spelled,
- ▶ **phonetics** – how words are pronounced,
- ▶ **morphology** – how words inflect,
- ▶ **syntax** – how words interact in sentences,



OK, why does it work at all?

Tiers of linguistic analysis

Computational approaches model various **tiers of language**:

- ▶ **graphematics** – how words are spelled,
- ▶ **phonetics** – how words are pronounced,
- ▶ **morphology** – how words inflect,
- ▶ **syntax** – how words interact in sentences,
- ▶ **pragmatics** – how sentences serve communicative purposes of human beings.



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '**judge**':



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '*judge*':

1. consists of 3 phonemes [j e j];



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '*judge*':

1. consists of 3 phonemes [j e j];
2. is a **singular noun in the nominative case**;



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '*judge*':

1. consists of 3 phonemes [j e j];
2. is a **singular noun in the nominative case**;
3. is a **nominal subject** dependent of the word 'sits' in the syntactic tree of our sentence.



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '*judge*':

1. consists of 3 phonemes [j e j];
2. is a **singular noun in the nominative case**;
3. is a **nominal subject** dependent of the word 'sits' in the syntactic tree of our sentence.

Such representations describe many important features of the word '*judge*'.

But not meaning (semantics).



To **model** means to capture important features of some phenomenon. For example, in a phrase '*The judge sits in the court*', the word '*judge*':

1. consists of 3 phonemes [j e j];
2. is a **singular noun in the nominative case**;
3. is a **nominal subject** dependent of the word 'sits' in the syntactic tree of our sentence.

Such representations describe many important features of the word '*judge*'.
But not meaning (semantics).

Question

Are these representations **discrete** or **continuous**?

How to represent meaning?

How to represent meaning?

- ▶ **Semantics** is difficult to represent formally.

How to represent meaning?

- ▶ **Semantics** is difficult to represent formally.
- ▶ Words which are **similar in their meaning** should possess **mathematically similar representations** independent of their surface forms.

How to represent meaning?

- ▶ **Semantics** is difficult to represent formally.
- ▶ Words which are **similar in their meaning** should possess **mathematically similar representations** independent of their surface forms.
- ▶ '*Judge*' must be similar to '*court*' but not to '*kludge*'...
- ▶ ...even though their surface form suggests the opposite.

How to represent meaning?

- ▶ **Semantics** is difficult to represent formally.
- ▶ Words which are **similar in their meaning** should possess **mathematically similar representations** independent of their surface forms.
- ▶ '*Judge*' must be similar to '*court*' but not to '*kludge*'...
- ▶ ...even though their surface form suggests the opposite.
- ▶ Why so?



Arbitrariness of a linguistic sign



Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]



Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



► **lantern**

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



- ▶ **lantern**
- ▶ **lykt**

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



- ▶ **lantern**
- ▶ **lykt**
- ▶ **ЛАННА**

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



- ▶ lantern
- ▶ lykt
- ▶ ЛАМПА
- ▶ lucerna

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



- ▶ lantern
- ▶ lykt
- ▶ лампа
- ▶ lucerna
- ▶ гэрэл

Arbitrariness of a linguistic sign

Unlike in **road signs**, there's no direct link between form and meaning in words [Saussure, 1916]
The concept of '*Lantern*' can be expressed by any sequence of letters or sounds in different languages:



- ▶ lantern
- ▶ lykt
- ▶ лампа
- ▶ lucerna
- ▶ гэрэл
- ▶ ...



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources

Methods of computationally representing semantics in natural languages fall into 2 large groups:



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources

Methods of computationally representing semantics in natural languages fall into 2 large groups:

1. **Manually building semantic networks or ontologies** (knowledge-based approach). Works top-down: from abstractions to real texts.



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources

Methods of computationally representing semantics in natural languages fall into 2 large groups:

1. **Manually building semantic networks or ontologies** (knowledge-based approach). Works top-down: from abstractions to real texts. For example, **WordNet** [Miller, 1995].



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources

Methods of computationally representing semantics in natural languages fall into 2 large groups:

1. **Manually building semantic networks or ontologies** (knowledge-based approach). Works top-down: from abstractions to real texts. For example, **WordNet** [Miller, 1995].
2. **Extracting semantics from usage patterns in text corpora** (distributional approach). Works bottom-up: from real texts to abstractions.



How do we humans know that '*lantern*' and '*lamp*' have similar meaning? What is **meaning**, after all?

And how can our ML models get this information?

Possible data sources

Methods of computationally representing semantics in natural languages fall into 2 large groups:

1. **Manually building semantic networks or ontologies** (knowledge-based approach). Works top-down: from abstractions to real texts. For example, **WordNet** [Miller, 1995].
2. **Extracting semantics from usage patterns in text corpora** (distributional approach). Works bottom-up: from real texts to abstractions.

The **second** approach is behind '**word embeddings**' (and most modern NLP).



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



- ▶ **Words with similar typical contexts have similar meaning.**



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



- ▶ **Words with similar typical contexts have similar meaning.**
- ▶ First formulated by:
 - ▶ philosopher **Ludwig Wittgenstein** (1930s);
 - ▶ linguists **Zelig Harris** [Harris, 1954] and **John Firth**.



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



- ▶ **Words with similar typical contexts have similar meaning.**
- ▶ First formulated by:
 - ▶ philosopher **Ludwig Wittgenstein** (1930s);
 - ▶ linguists **Zelig Harris** [Harris, 1954] and **John Firth**.
- ▶ *'You shall know a word by the company it keeps'*
[Firth, 1957]



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



- ▶ **Words with similar typical contexts have similar meaning.**
- ▶ First formulated by:
 - ▶ philosopher **Ludwig Wittgenstein** (1930s);
 - ▶ linguists **Zelig Harris** [Harris, 1954] and **John Firth**.
- ▶ *'You shall know a word by the company it keeps'*
[Firth, 1957]
- ▶ More details in [Brunila and LaViolette, 2022].



Hypothesis: meaning is actually a sum of contexts.

Distributional differences will always be enough to explain **semantic differences**:



- ▶ **Words with similar typical contexts have similar meaning.**
- ▶ First formulated by:
 - ▶ philosopher **Ludwig Wittgenstein** (1930s);
 - ▶ linguists **Zelig Harris** [Harris, 1954] and **John Firth**.
- ▶ *'You shall know a word by the company it keeps'* [Firth, 1957]
- ▶ More details in [Brunila and LaViolette, 2022].
- ▶ **Distributional semantics models** (DSMs) get information from lexical co-occurrences in large natural corpora.



Contexts for 'tea':

establishments, besides two livery stables, a en things, because their methods of family to , as, indeed, A waiter comes in with the let me always remain here.' "I prefer weak hell. I should think you had drunk enough . Not a bit. Come in and have some responsibility. And greatly as we enjoyed our your naturally liking me. (She is and had] Tell them I shan't be home to , that was Mr. McComas will not come to , or asked you to have a cup of woman can hardly know one places Gilbey's of my - my hopes.' BROADBENT. He'll want : THE MANAGER. Can I take any order? Some to Tramp.} Will you drink a sup of are trying to sleep." the evening after your GUINNESS. I'll go get you some fresh sional men, artists, and even with laborers the Lutches and Mrs. Rance the attendance at she came over to the great house to . The Baroness found it amusing to go to tea; she dressed as if for dinner. The would be dead in two years, as the	tea They never boasted of Robert Acton, nor indulg tea at once. pose, which you carry so well tea. He places the tray on the table. Jasper tea!" cried Daisy, and she went off with the tea in Chin a. life; it is a failure, tea. Stay to dinner. every year. Don't persist tea Crusoe island. Then there's the religious diff tea in the evening. Afraid though as he was tea, will you, LADY BRITOMART. I must get the tea, ma'am: he has gone to call upon tea. It's not human. ugly woman must have tea on the table before him]. The lady that tea. Let us have some. BURGE-LUBIN [_resolutely ge tea? would THE SHE-ANCIENT. Speak, Arjillax: you w tea with myself and the the happiest person in tea. "Better still-then there you are!" And Streth tea, ducky. [She takes up the its burden, is tea services out and made the people who had tea just in the right place on the west tea. She had let the proposal that she should tea; she dressed as if for dinner. The tea- tea-table offered an anomalous and picturesque rep tea-table. Be serious, Felix. You forget that I
--	---

Distributional hypothesis



Contexts for 'tea':

establishments, besides two livery stables, a
en things, because their methods of family to
, as, indeed, A waiter comes in with the
let me always remain here.' "I prefer weak
hell. I should think you had drunk enough
. Not a bit. Come in and have some
responsibility. And greatly as we enjoyed our
your naturally liking me. (She is and had
) Tell them I shan't be home to
, that was Mr. McComas will not come to
, or asked you to have a cup of
woman can hardly know one places Gilbey's
of my - my hopes.' BROADBENT. He'll want
: THE MANAGER. Can I take any order? Some
to Tramp.} Will you drink a sup of
are trying to sleep." the evening after your
GUINNESS. I'll go get you some fresh
sional men, artists, and even with laborers
the Lutches and Mrs. Rance the attendance at
she came over to the great house to
. The Baroness found it amusing to go to
tea; she dressed as if for dinner. The
would be dead in two years, as the

tea They never boasted of Robert Acton, nor indulg
tea at once. pose, which you carry so well
tea. He places the tray on the table. Jasper
tea!" cried Daisy, and she went off with the
tea in Chin a. life; it is a failure,
tea. Stay to dinner. every year. Don't persist
tea Crusoe island. Then there's the religious diff
tea in the evening. Afraid though as he was
tea, will you, LADY BRITOMART. I must get the
tea, ma'am: he has gone to call upon
tea. It's not human. ugly woman must have
tea on the table before him]. The lady that
tea. Let us have some. BURGE-LUBIN [_resolutely ge
tea? would THE SHE-ANCIENT. Speak, Arjillax: you w
tea with myself and the the happiest person in
tea. "Better still--then there you are!" And Streth
tea, ducky. [She takes up the its burden, is
tea services out and made the people who had
tea just in the right place on the west
tea. She had let the proposal that she should
tea; she dressed as if for dinner. The tea-
tea-table offered an anomalous and picturesque rep
tea-table. Be serious, Felix. You forget that I

Contexts for 'coffee':

prompt his an incident in my life as
ek her out all courteously, PETKOFF (over his
could be done, too,' he remarked, sipping his
, of us. I should like a cup of
UKA (innocently). Perhaps you would like some
her in public because he has fallen head
manners for he was novels, broken backed,
stretches her hand across the table for the
ittle sitting-room, and cigarettes, after the
had just given me a pannikin of hot
a heavy roll coming; tried to save my
I'll have a claret cup instead of
t of trouble travelling. And then, with fresh
. Your word had such weight with me!" fresh
wont press you. "Try a weed with your
coffee for breakfast. Of course, hes too _Two fig
coffee and cigaret). I don't believe in going
coffee. 'Bury him in some sort,' I explained. 'One
coffee. MICHAEL. If you'd come in better hours,
coffee, sir? DISCOVERY ANTICIPATED BY DIVINATION s
coffee-colored heathens and pestilential white agi
coffee stained, torn and "'This was the last time
coffee pot.) welcome, an expression which drops in
coffee, had been permitted by the ladies, and in
coffee...Slapped it down there, on my chest--bange
coffee, burnt my fingers...and fell out of my
coffee. Put some first night that we've come
coffee, a clean cup, and a brandy bottle on
coffee? He gave his friend a glance as to
coffee. Local tobacco. The black coffee you get at



- ▶ Your neural classifiers in Obligatory 1 implicitly learned **vector representations** for words (embeddings).
- ▶ In practice, representing word meaning with vectors was first popularized in **psychology** by [Osgood et al., 1964]...
- ▶ ...then developed by many others.



- ▶ Your neural classifiers in Obligatory 1 implicitly learned **vector representations** for words (embeddings).
 - ▶ In practice, representing word meaning with vectors was first popularized in **psychology** by [Osgood et al., 1964]...
 - ▶ ...then developed by many others.
- ▶ Word vectors can be created manually...
 - ▶ ...but in most cases, corpus-driven **distributional** methods are much more efficient.

Representing words with vectors



Componential analysis: manual creation of word vectors

TABLE 2: THE HYPOTHESED COMPONENTIAL ANALYSIS OF KINSHIP TERMS

Kinship terms	[MALE]	[ASCEND]	[DESCEND]	[LINEAL]
<i>Father</i>	+	+	-	+
<i>Mother</i>	-	+	-	+
<i>Uncle</i>	+	+	-	-
<i>Aunt</i>	-	+	-	-
<i>Brother</i>	+	-	-	+
<i>Sister</i>	-	-	-	+
<i>Son</i>	+	-	+	+
<i>Daughter</i>	-	-	+	+
<i>Nephew</i>	+	-	+	-
<i>Niece</i>	-	-	+	-
<i>Cousin</i>	+/-	-	-	-

[Widyastuti, 2010]

We **will not** do this. We will use **distributional** vector models (*next sub-lecture 5.2*).

References I



Brunila, M. and LaViolette, J. (2022).

What company do words keep? revisiting the distributional semantics of J.R. Firth & Zellig Harris.

In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4403–4417, Seattle, United States. Association for Computational Linguistics.



Firth, J. (1957).

A synopsis of linguistic theory, 1930-1955.

Blackwell.







Harris, Z. S. (1954).

Distributional structure.

Word, 10(2-3):146–162.

References II

-  Miller, G. (1995).
Wordnet: a lexical database for English.
Communications of the ACM, 38(11):39–41.
-  Osgood, C. E., Suci, G. J., and Tannenbaum, P. H. (1964).
The measurement of meaning.
University of Illinois Press.
-  Saussure, F. d. (1916).
Course in general linguistics.
Duckworth.
-  Turney, P., Pantel, P., et al. (2010).
From frequency to meaning: Vector space models of semantics.
Journal of artificial intelligence research, 37(1):141–188.



Widyastuti, S. (2010).

Componential analysis of meaning: Theory and applications.

Journal of English and Education, 4(1):116–128.