# Seminar in *Artificial Intelligence*
## Word embedding

Marcin Trebunia, Dominik Rygiel, Konrad Adamczyk

Department of Telecommunications

27.05.2019

**AGH**

## Agenda

**1 Introduction**
- Necessity for encoding text
- Simpler types of encoding

**2 Word embedding details**
- Detail1
- Detail2
- Detail3
- Detail4

**3 Applications**

**4 Problems and limitations**

# What is word embedding?

What a **lovely** day.
What a **nice** day.

- Machine learning models take vectors (arrays of numbers) as input.
- ....

$$
\begin{aligned}
\text{What} &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \\
\text{a} &= \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix} \\
\text{lovely} &= \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \\
\text{nice} &= \begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \\
\text{day} &= \begin{bmatrix} 0 & 0 & 0 & 0 & 1 \end{bmatrix}
\end{aligned}
$$

# One hot encoding ctd.

- Words completely independent of each other
- Inefficient approach: vector is sparse

Example:

- Dictionary of 10,000 words
- One hot encode each word
- Each vector's elements are 99.99% zeros!

$$
\begin{aligned}
\text{What} &= [1] \\
\text{a} &= [2] \\
\text{lovely} &= [3] \\
\text{nice} &= [4] \\
\text{day} &= [5]
\end{aligned}
$$

+ Efficient - dense vector
− Encoding arbitrary - does not catch relationships between words
− Can be challenging for a model to interpret

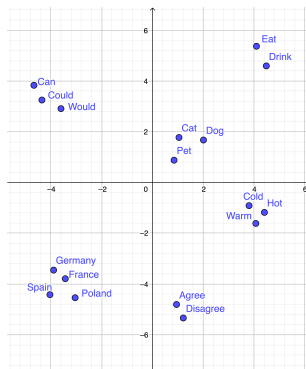$$\text{What} = \begin{bmatrix} 1.2 & -0.1 & 4.3 & 3.2 \end{bmatrix}$$
$$\text{a} = \begin{bmatrix} 0.4 & 2.5 & -0.9 & 0.5 \end{bmatrix}$$
$$\text{lovely} = \begin{bmatrix} 2.1 & 0.3 & 0.1 & 0.4 \end{bmatrix}$$
$$\text{nice} = \begin{bmatrix} 2.0 & 0.4 & 0.3 & 0.5 \end{bmatrix}$$
$$\text{day} = \begin{bmatrix} 3.0 & -0.6 & 3.5 & -0.8 \end{bmatrix}$$
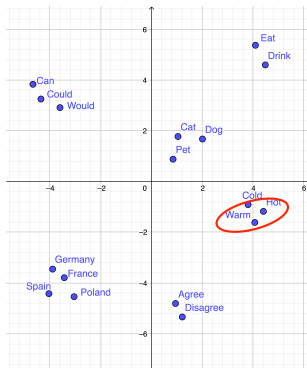
# Word embedding

- Words with similar context occupy close spatial positions
- The cosine of the angle between words' vectors should be close to 1 (angle close to 0)

Caption of the figure

# Word Embedding


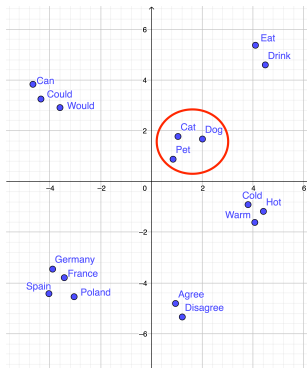
Words are synonyms

# Word Embedding



Words are antonyms

Words are value on a scale

Words are hyponym - hypernym

# Word Embedding



Words appear in similar context

## How can we use it?

- If user search for "Dell notebook battery size" we would like to match it also with "Dell laptop battery capacity"
- If user search for "Cracow Motel" we would like to match it also with "Krakow Hotel"

.....One of the main limitations of word embeddings (word vector space models in general) is that possible meanings of a word are conflated into a single representation (a single vector in the semantic space). Sense embeddings[17] are proposed as a solution to this problem: individual meanings of words are represented as distinct vectors in the space.....

# Thank you for your attention!

# Q & A

**AGH**

🌐 Intro to word embeddings.
https://www.tensorflow.org/alpha/tutorials/text/
word_embeddings.
Accessed: 2019-05-11.

🌐 Intro to word embeddings.
https://towardsdatascience.com/
introduction-to-word-embedding-and-word2vec-652d0c2060f
fbclid=IwAR3c2RpZOmbWC84_
mKFtRI6PwTD7vJRxiquKPp2Y3en3_OfDpBsWjjSinv8.
Accessed: 2019-05-11.