# Seminar in *Artificial Intelligence*
## Word embedding

**Marcin Trebunia, Dominik Rygiel, Konrad Adamczyk**

**Department of Telecommunications**

27.05.2019

**Agenda**

1. **Introduction**
   - Necessity for encoding text
   - Different types of encoding

# What is word embedding?

## Encoding text

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

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

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

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

Source: .

Example:

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

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

# Unique number encoding ctd.

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

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

- 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

Words are synonyms

# Word Embedding



Words are antonyms

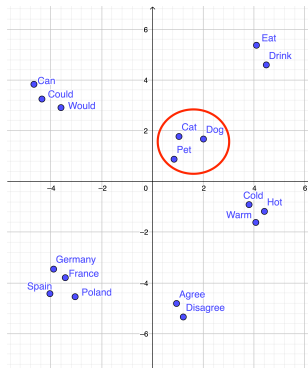# Word Embedding

**Slide with a Figure from a File**



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"

# Title of the Slide

**Subtitle of the Slide: Use Only if Necessary. . .**

- Bulalet point 1.
- Bullet point 2 -e– <span style="color:red">you can emphasise</span> a text.

Text:

Blocks are good for important notions.

Alertblocks are even better to catch the attention.
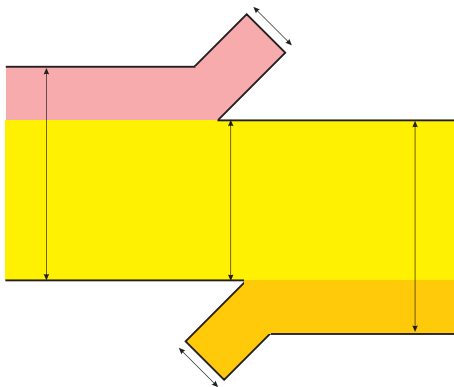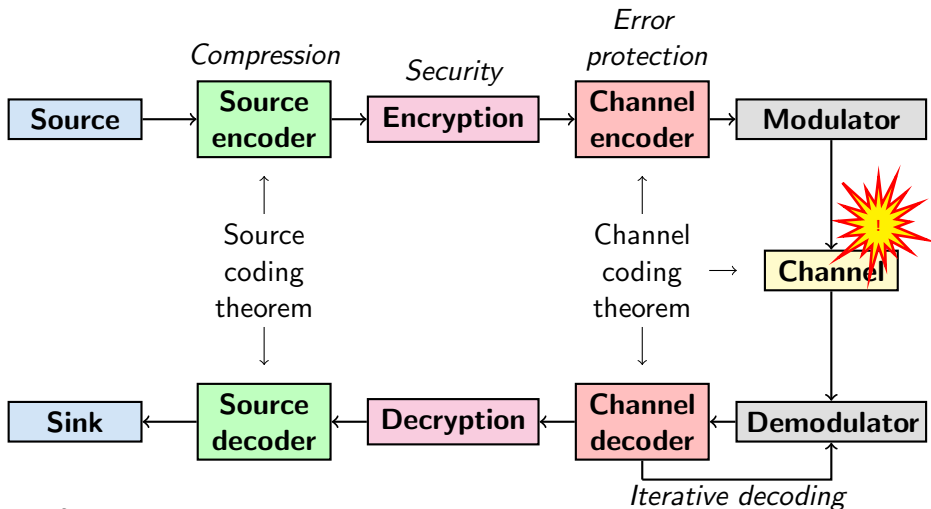
**Title**

You can have the title in the (alert)block.

**Figure:** Caption of the figure

I am sorry that the example figures and the table concern the information theory ☺

**Table:** Caption of the Table

| Data | Entropy |
|---|---|
| Plain text | 4.347 |
| Native executables | 5.099 |
| Packed executables | 6.801 |
| Encrypted executables | 7.175 |

Source: .

# Thank you for your attention!

# Q & A

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

Robert Lyda and James Hamrock.
Using Entropy Analysis to Find Encrypted and Packed
Malware.
*IEEE Security & Privacy*, 5(2):40–45, March/April 2007.

Todd K. Moon.
*Error Correction Coding*.
John Wiley & Sons, Inc., Hoboken, NJ, 2005.