# SkillFactory

# Word embeddings CNN for texts
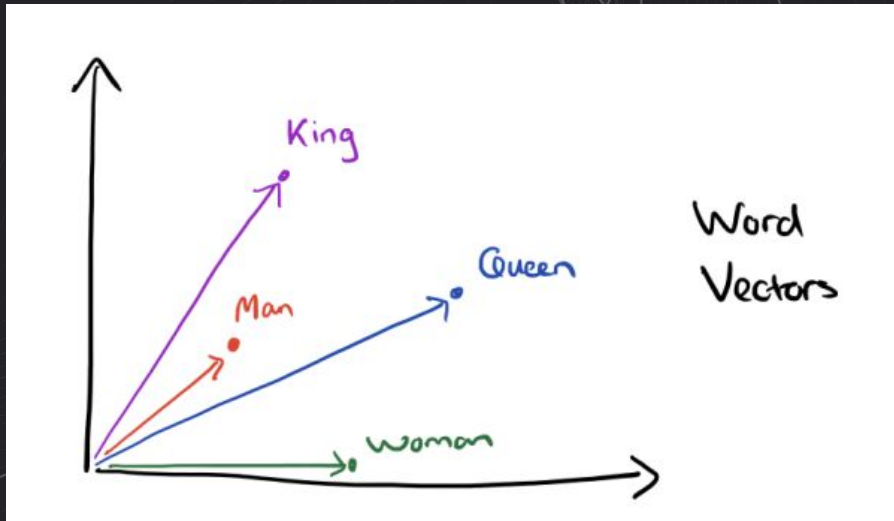
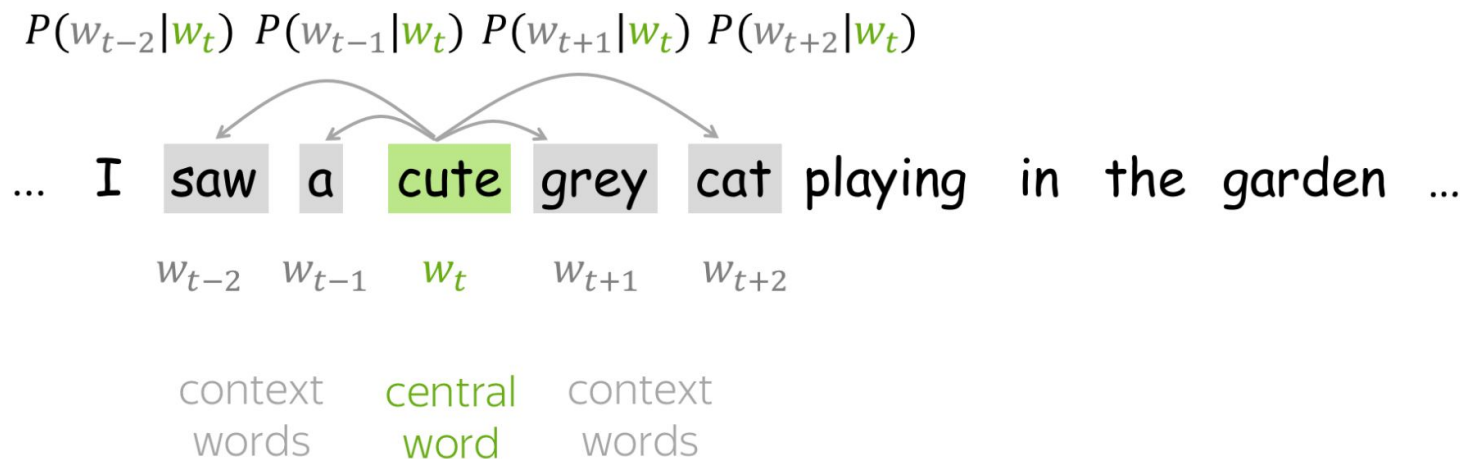**Sidorov Nikita**

MLE (NLP) in Sber

# Word2vec
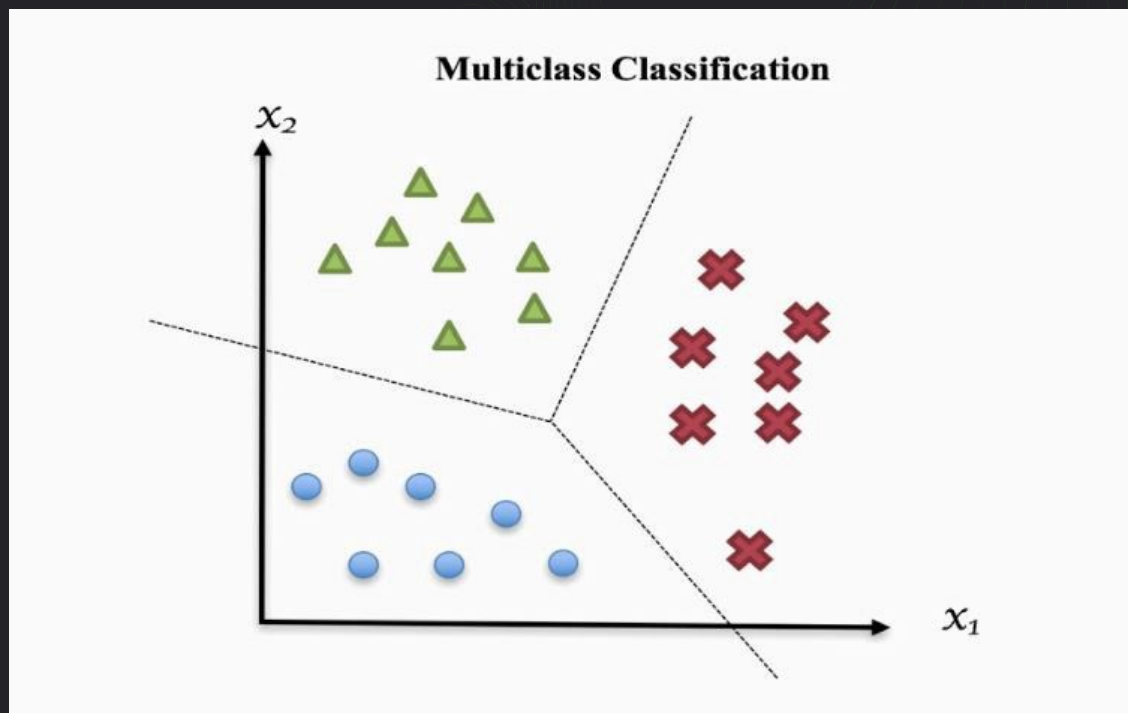
# What task is solving Word2vec

$P(w_{t-2}|w_t)$ $P(w_{t-1}|w_t)$ $P(w_{t+1}|w_t)$ $P(w_{t+2}|w_t)$

... I saw a cute grey cat playing in the garden ...

$w_{t-2}$  $w_{t-1}$  $w_t$  $w_{t+1}$  $w_{t+2}$

context words    central word    context words

example from Lena Voita NLP course
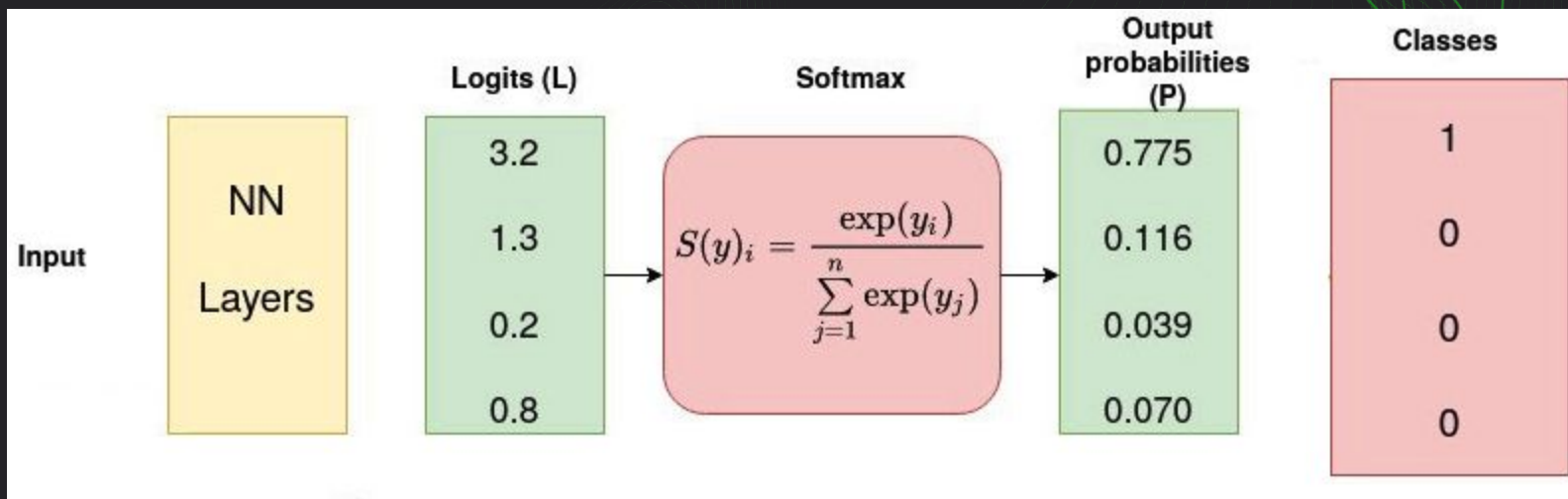
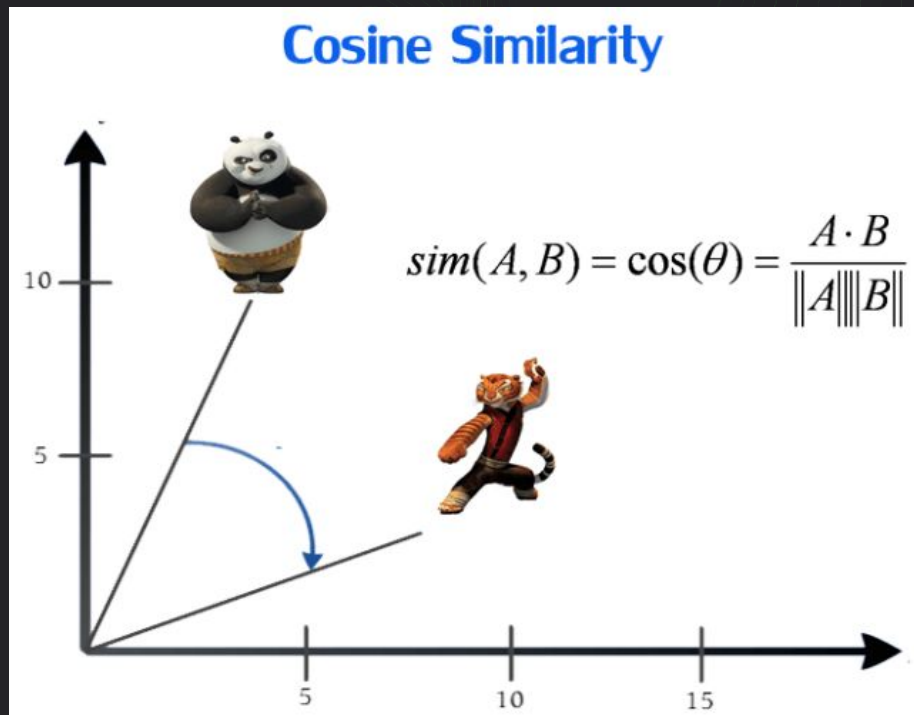# What task is solving Word2vec

# Cross-entropy

# Cross-entropy

$$CCE(p, t) = -\sum_{c=1}^{C} t_{o,c} \log (p_{o,c})$$

# Distributional semantics

1. A bottle of _____ is on the table.

2. Everybody likes _____.

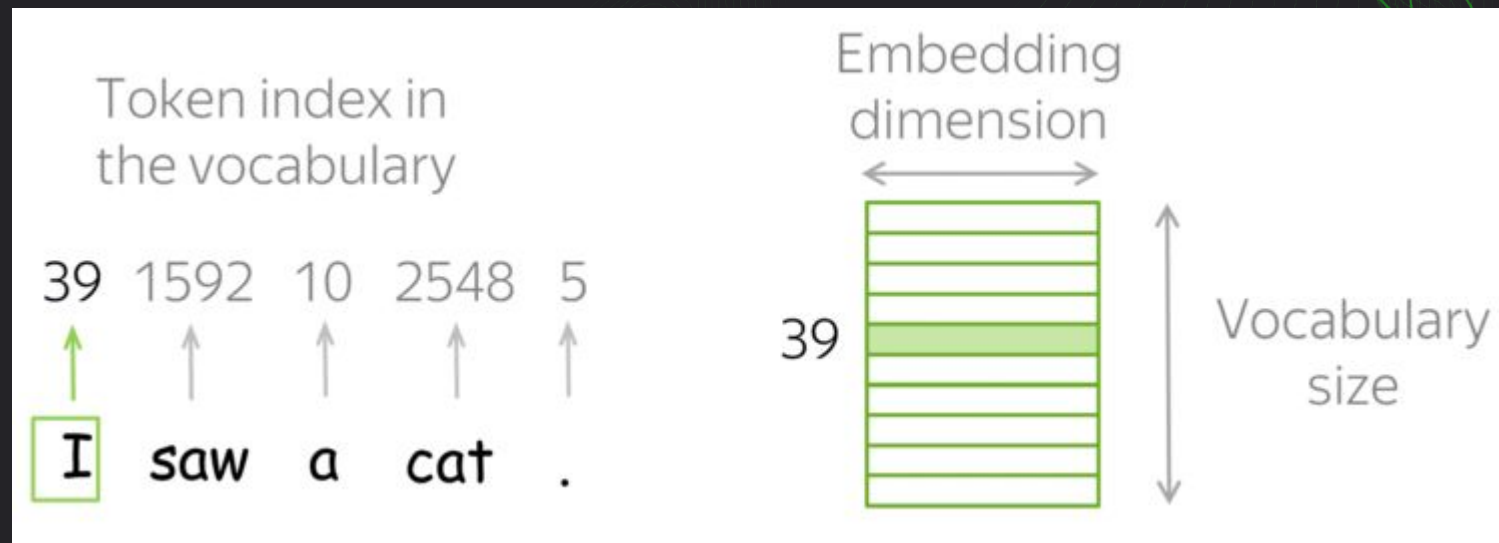3. Don`t have _____ before you drive.

4. We make _____ out of corn.

example from Jacob Eisenstein's NLP notes

# How to measure distances between vectors?

# What we want from word embeddings?

# What we want from word embeddings?



$$[0 \quad 0 \quad 0 \quad 1 \quad 0] \times \begin{bmatrix} 8 & 2 & 1 & 9 \\ 6 & 5 & 4 & 0 \\ 7 & 1 & 6 & 2 \\ 1 & 3 & 5 & 8 \\ 0 & 4 & 9 & 1 \end{bmatrix} = [1 \quad 3 \quad 5 \quad 8]$$

One-hot vector    Embedding Weight Matrix    Hidden layer output

# Word2vec



a) Learns Analogy

b) Similar Words have same angles

# Word2vec



example from Chris McCormick blog

# Word2vec

# Word2vec

How to get probabilities for occurrence of words?
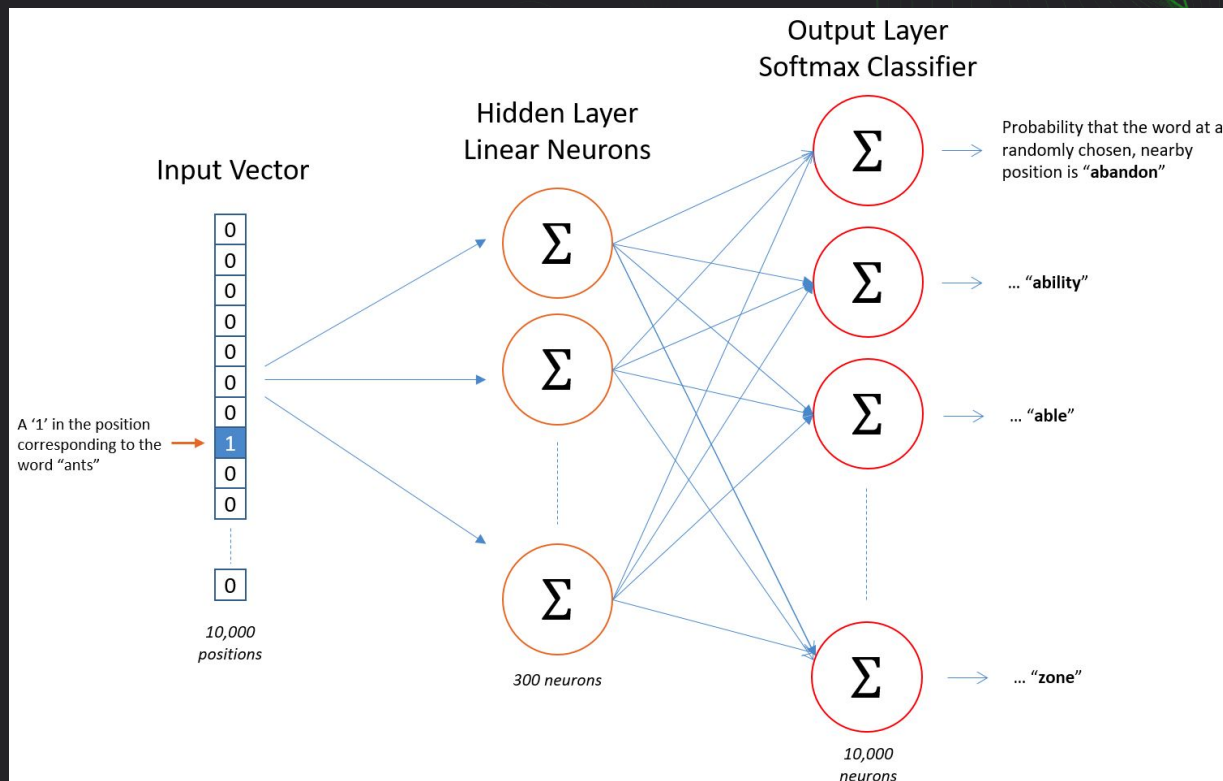
$$f(s)_i = \frac{e^{s_i}}{\sum_j^C e^{s_j}} \qquad CE = -\sum_i^C t_i log(f(s)_i)$$

Softmax → Cross-Entropy Loss

# Word2vec approaches for training

# CBOW



## CBOW - Working

Hope can set you free.

$V_{5 \times 1}$, one hot vector of "Hope"

$W_{3 \times 5}$

3 nodes in hidden layer

Compare and Update weights

$W'_{5 \times 3}$

Actual Target

$V_{5 \times 1}$, one hot vector of "Set"

$W_{3 \times 5}$

$V_{5 \times 1}$, predicted one hot vector of "Can"

# CBOW

# Skip-gram



Skip-Gram: from central predict context (one at a time)

# How to learn

# How to learn

1) For each word we would have vector of context

2) Represent each context as a $d$ dimensional vector

3) Initialize all vectors to random weights

4) Arrange vectors in two matrices W and C

# How to learn

$$\log p(c \mid w; \theta) = \frac{\exp v_c \cdot v_w}{\sum_{c' \in C} \exp v_{c'} \cdot v_w}$$

- predict context word(s)

- from word w

# How to learn

... I saw a cute grey cat playing in the garden ...

**1.** Take dot product of $v_{cat}$ with **all** $u$     **2.** exp     **3.** sum all

$$u_{w1}^T v_{cat} \longrightarrow \exp(u_{w1}^T v_{cat})$$
$$u_{w3}^T v_{cat} \longrightarrow \exp(u_{w3}^T v_{cat})$$
$$\vdots$$
$$\boxed{u_{cute}^T v_{cat}} \longrightarrow \exp(u_{cute}^T v_{cat})$$
$$\vdots \quad ①$$
$$u_{wn}^T v_{cat} \longrightarrow \exp(u_{wn}^T v_{cat})$$

$$\boxed{\sum_{w \in V} \exp(u_w^T v_{cat})}$$

②

cat    cute

$v$        $u$

example from Lena Voita NLP course

# How to learn

... I saw a cute grey cat playing in the garden ...

**4.** get loss (for this one step)

$$J_{t,j}(\theta) = -u_{cute}^T v_{cat} + \log \sum_{w \in V} \exp(u_w^T v_{cat})$$

1      2

**5.** evaluate the gradient, make an update

$$v_{cat} := v_{cat} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{cat}}$$

$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$

# What`s the problem?

# Negative sampling



Dot product of $v_{cat}$:
- with $u_{cute}$ - increase,
- with <u>all other</u> $u$ - decrease

$\longrightarrow$

Dot product of $v_{cat}$:
- with $u_{cute}$ - increase,
- with <u>a subset of other</u> $u$ - decrease

Negative samples: randomly selected K words

$u_{w1}^T v_{cat}$
$u_{w3}^T v_{cat}$  decrease
$\vdots$
$u_{cute}^T v_{cat}$  increase
$\vdots$
$u_{wn}^T v_{cat}$  decrease

$v$  $u$

$u_{w_{i_1}}^T v_{cat}$
$u_{w_{i_2}}^T v_{cat}$  decrease
$u_{cute}^T v_{cat}$  increase
$u_{w_{i_{K-1}}}^T v_{cat}$
$u_{w_{iK}}^T v_{cat}$  decrease

$v$  $u$

Parameters to be updated:

- $v_{cat}$
- $u_w$ for all $w$ in the vocabulary   $|V|$ + 1 vectors

Parameters to be updated:

- $v_{cat}$
- $u_{cute}$ and $u_w$ for $w$ in K negative examples   K + 2 vectors

example from Lena Voita NLP course

# OOV words

# Fasttext

- Take not only words, but n-grams in this words

- harder to compute

- longer to train

- bigger models

- well works for morphologically rich languages

# Fasttext

- Skip-gram model as base model

- take each word and n-grams for it (from 3 to 6)

- for reducing space we use hashing trick

- negative sampling is our everything

# Convolutions

In deep learning, a convolutional neural network (CNN) is a class of artificial neural network, most commonly applied to analyze visual imagery. They are also known as artificial neural networks that slide along input features and provide translation equivariant responses known as feature maps. Counter-intuitively, most convolutional neural networks are only equivariant, as opposed to invariant, to translation.

# Convolutions

# Convolutions

# Convolution for texts

<pad>     I     like     to     swim.     <eos>     <pad>

# Convolution for texts



&lt;pad&gt;      I      like      to      swim.      &lt;eos&gt;      &lt;pad&gt;

# Convolution for texts



<pad>  I  like  to  swim.  <eos>  <pad>

# Convolution for texts



<pad>    I    like    to    swim.    <eos>    <pad>

# Convolution for texts



<pad>   I   like   to   swim.   <eos>   <pad>

# Convolution for texts

Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|------|------|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

# Convolution for texts

Text

| | | | | |
|---|---|---|---|---|
| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

Apply filters of size 3 and that have 4 channels

| | | | |
|---|---|---|---|
| 3 | 2 | 1 | -1 |
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

# Convolution for texts

Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|-----|-----|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

Apply filters of size 3 and that have 4 channels

| 3 | 2 | 1 | -1 |
|---|---|---|----|
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

Result

| <IL | |
|-----|---|
| ILT | |
| LTS | |
| TS< | |
| S<< | |

# Convolution for texts

## Text

| | | | | |
|---|---|---|---|---|
| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

## Apply filters of size 3 and that have 4 channels

| | | | |
|---|---|---|---|
| 3 | 2 | 1 | -1 |
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

## Result

| | |
|---|---|
| <IL | 4.2 |
| ILT | |
| LTS | |
| TS< | |
| S<< | |

# Convolution for texts

## Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|------|------|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

## Apply filters of size 3 and that have 4 channels

| 3 | 2 | 1 | -1 |
|---|---|---|----|
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

## Result

| <IL | 4.2 |
|-----|-----|
| ILT | 2.4 |
| LTS | |
| TS< | |
| S<< | |

# Convolution for texts

Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|-----|-----|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

Apply filters of size 3 and that have 4 channels

| 3 | 2 | 1 | -1 |
|---|---|---|----|
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

Result

| <IL | 4.2 |
|-----|-----|
| ILT | 2.4 |
| LTS | 2.5 |
| TS< | |
| S<< | |

# Convolution for texts

Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|-----|-----|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

Apply filters of size 3 and that have 4 channels

| 3 | 2 | 1 | -1 |
|---|---|---|----|
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

Result

| <IL | 4.2 |
|-----|-----|
| ILT | 2.4 |
| LTS | 2.5 |
| TS< | 3.2 |
| S<< | |

# Convolution for texts

Text

| <pad> | 0.3 | 0.4 | -0.2 | -0.6 |
|-------|-----|-----|------|------|
| I | 0.5 | 0.1 | -0.3 | 0.4 |
| like | -0.1 | 0.5 | 0.8 | -0.2 |
| to | 0.2 | -0.3 | -0.4 | -0.5 |
| swim. | -0.7 | 0.5 | 0.9 | 0.1 |
| <eos> | -0.4 | 0.4 | 0.1 | -0.5 |
| <pad> | 0.1 | -0.6 | -0.3 | 0.2 |

Apply filter of size 3, that have 4 channels

| 3 | 2 | 1 | -1 |
|---|---|---|----|
| 1 | 0 | 2 | 1 |
| -1 | 1 | 1 | -2 |

Result

| <IL | 4.2 |
|-----|-----|
| ILT | 2.4 |
| LTS | 2.5 |
| TS< | 3.2 |
| S<< | -2.4 |

# Pooling for convolutions



example from Lena Voita NLP course