# РАСШИФРОВЫВАЕМ ЗАГАДОЧНЫЕ КОДЫ: РУКОПИСЬ ВОЙНИЧА, БИОИНФОРМАТИКА И МОЗГОВАЯ АКТИВНОСТЬ
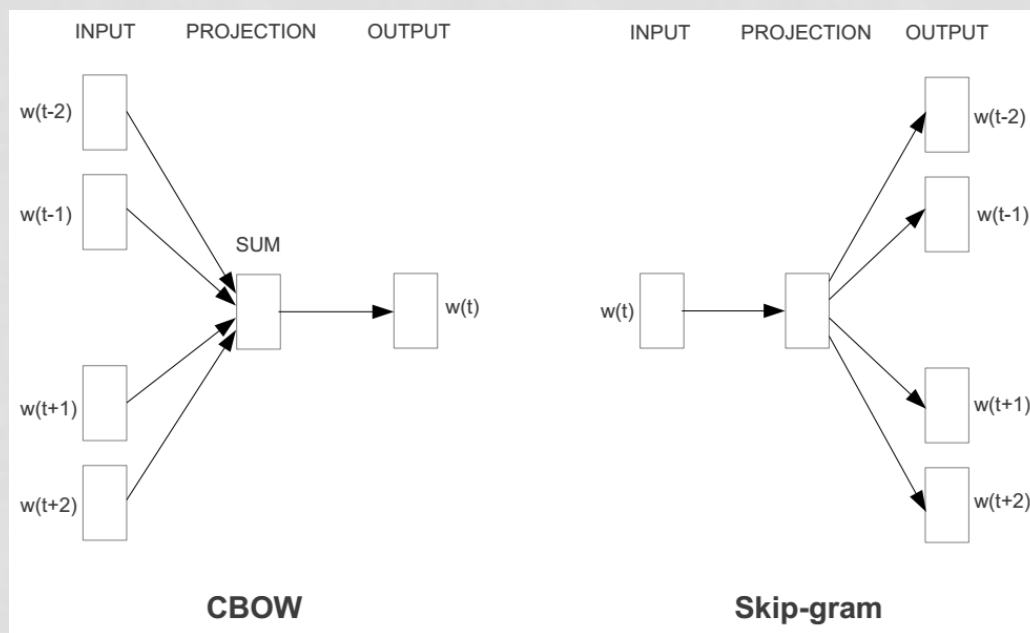
….и как в этом поможет word2wec

ДМИТРИЙ НОВИЦКИЙ

1

# WORD2VEC : A RECALL

- Represent each word with a low-dimensional vector

- Word similarity = vector similarity

- Key idea: Predict surrounding words of every word

- Faster and can easily incorporate a new sentence/document or add a word to the vocabulary
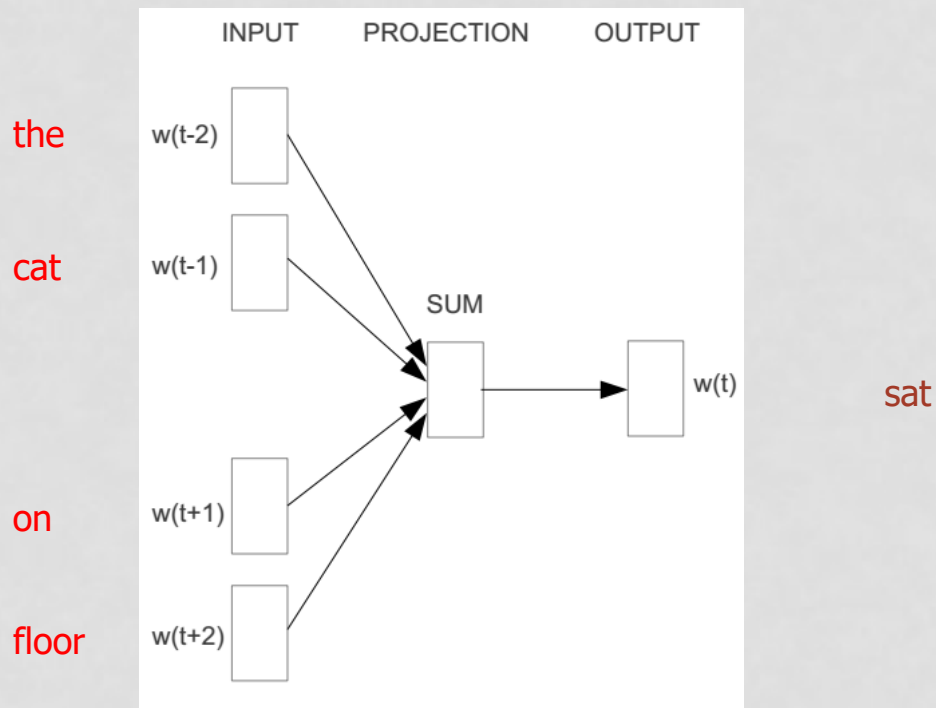
# REPRESENT THE MEANING OF WORD – WORD2VEC

- 2 basic neural network models:
  - Continuous Bag of Word (CBOW): use a window of word to predict the middle word
  - Skip-gram (SG): use a word to predict the surrounding ones in window.

# WORD2VEC – CONTINUOUS BAG OF WORD

- E.g. "The cat sat on floor"
  - Window size = 2



the

cat

on

floor

sat

Input layer

Index of "cat" in vocabulary
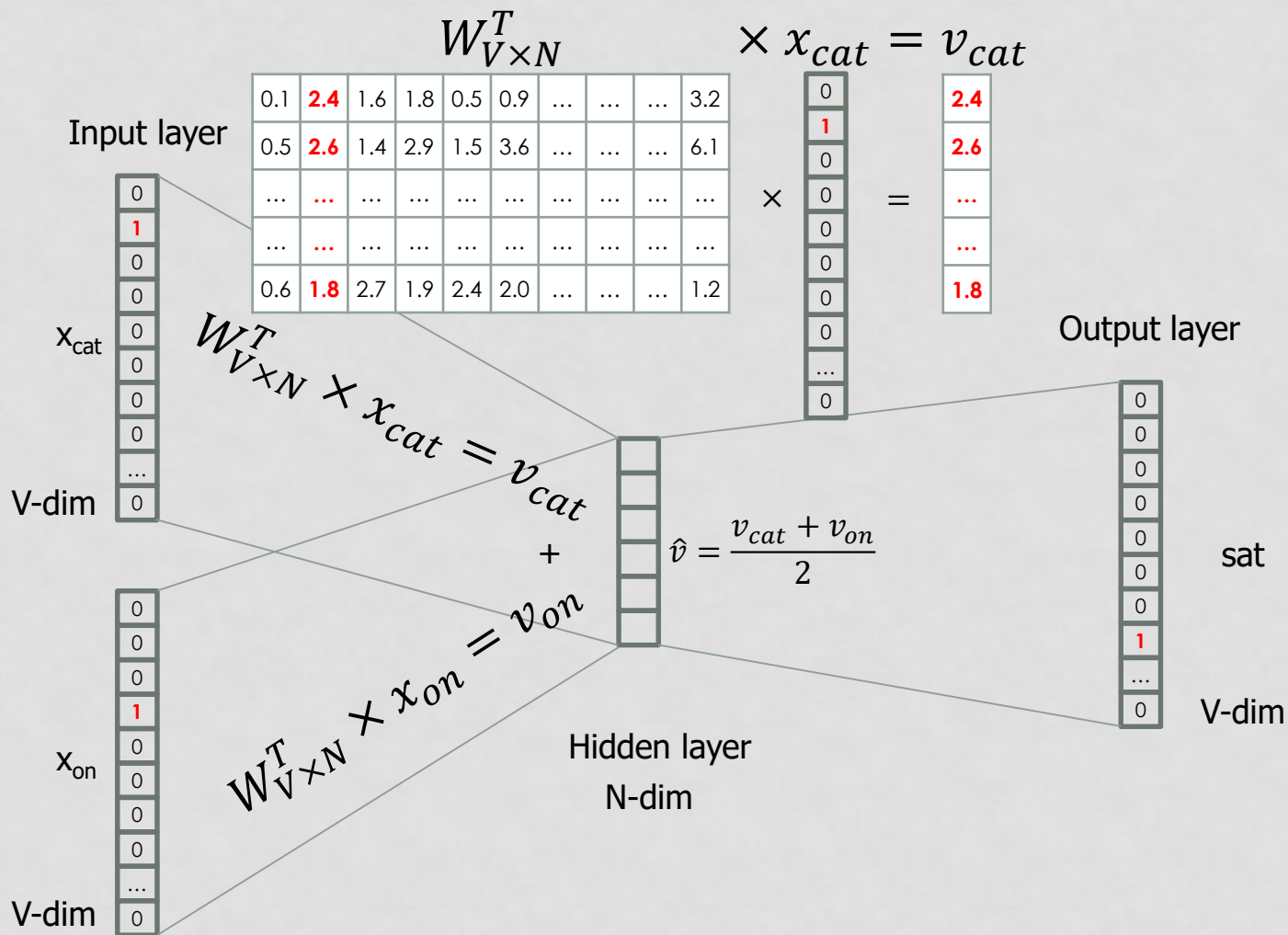
cat

0
1
0
0
0
0
0
0
...
0

one-hot vector

on

0
0
0
1
0
0
0
0
...
0

Hidden layer

Output layer

sat

one-hot vector

0
0
0
0
0
0
0
1
...
0

5

We must learn W and W′

Input layer

Hidden layer

Output layer

cat

| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |

V-dim | 0 |

$W_{V \times N}$

| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |

on

V-dim | 0 |

$W_{V \times N}$

N-dim

$W'_{N \times V}$

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| ... |
| 0 |

sat

V-dim

N will be the size of word vector

6

$$W_{V \times N}^{T} \qquad \times \, x_{cat} = v_{cat}$$

Input layer

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|---------|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$x_{cat}$

$$W_{V \times N}^{T} \times x_{cat} = v_{cat}$$

$+$

V-dim

$$W_{V \times N}^{T} \times x_{on} = v_{on}$$

$x_{on}$

V-dim

Input vector $x_{cat}$:
0
1
0
0
0
0
0
0
...
0

Input vector $x_{on}$:
0
0
0
1
0
0
0
0
...
0

$x_{cat}$ column:
0
1
0
0
0
0
0
0
...
0

$v_{cat}$ column:
**2.4**
**2.6**
**...**
**...**
**1.8**

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

Hidden layer
N-dim

Output layer

sat

V-dim

Output vector:
0
0
0
0
0
0
0
1
...
0

$$W_{V \times N}^T \times x_{on} = v_{on}$$

Input layer

| 0.1 | 2.4 | 1.6 | **1.8** | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|---------|-----|-----|-----|-----|-----|-----|
| 0.5 | 2.6 | 1.4 | **2.9** | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | **...** | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | **...** | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | **1.9** | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$x_{cat}$

V-dim

$W_{V \times N}^T \times x_{cat} = v_{cat}$

$+$

$W_{V \times N}^T \times x_{on} = v_{on}$

$x_{on}$

V-dim

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

Hidden layer
N-dim

Output layer

sat

V-dim

8

Input layer

| 0 |
| 1 |
| 0 |
| 0 |
cat | 0 |
| 0 |
| 0 |
| 0 |
| ... |
V-dim | 0 |

$W_{V \times N}$

Hidden layer

$$W'_{V \times N} \times \hat{v} = z$$

$\hat{v}$

N-dim

Output layer

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| ... |
| 0 |

$\hat{y} = softmax(z)$

$\hat{y}_{\text{sat}}$

V-dim

| 0 |
| 0 |
| 0 |
| 1 |
on | 0 |
| 0 |
| 0 |
| 0 |
| ... |
V-dim | 0 |

$W_{V \times N}$

N will be the size of word vector

9

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |

cat

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

V-dim

| 0 |

We would prefer $\hat{y}$ close to $\hat{y}_{sat}$

$W_{V \times N}$

Hidden layer

Output layer

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

$\hat{v}$

| 0.01 |
|---|
| 0.02 |
| 0.00 |
| 0.02 |

$W'_{V \times N} \times \hat{v} = z$
$\hat{y} = softmax(z)$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| **1** |
| ... |
| 0 |

$\hat{y}_{sat}$

V-dim

| 0 |
|---|
| 0 |
| 0 |
| **1** |

on

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

V-dim

| 0 |

$W_{V \times N}$

N-dim

N will be the size of word vector

| 0.01 |
|---|
| 0.02 |
| 0.01 |
| **0.7** |
| ... |
| 0.00 |

$\hat{y}$

$$W^T_{V \times N}$$

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

Contain word's vectors

Input layer

$x_{cat}$

V-dim

$W_{V \times N}$

$W_{V \times N}$

$x_{on}$

V-dim

Hidden layer
N-dim

$W'_{V \times N}$

Output layer

sat

V-dim

We can consider either W or W' as the word's representation.
Or even take the average.

11

# SOME INTERESTING RESULTS

# Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?  $\longrightarrow$  $d = \arg\max_{x} \dfrac{(w_b - w_a + w_c)^T w_x}{||w_b - w_a + w_c||}$

man:woman :: king:?

| | | |
|---|---|---|
| + | king | [ 0.30 0.70 ] |
| - | man | [ 0.20 0.20 ] |
| + | woman | [ 0.60 0.30 ] |

queen [ 0.70 0.80 ]

# WORD ANALOGIES

# МАНУСКРИПТ ВОЙНИЧА

# EVA TRANSCRIPTION

- To train this model, I had to parse and extract the transcription from the EVA (European Voynich Alphabet) to be able to feed the Voynich sentences into the word2vec model. This EVA transcription has the following format:

- <f1r.P1.1;U> fya!ys.ykal.ar.ytaiin.shol.shory.***!r*s.y.kor.s holdo*- #

- <f1r.P1.2;H> sory.ckhar.o!r.y.kair.chtaiin.shar.are.cthar.c thar.dan!-

# T-SNE VISUALIZATION

# RESULTS

- >>> w2v_model.most_similar("octhey")
- [('qoekaiin', 0.6402825713157654),
- ('otcheody', 0.6389687061309814),
- ('ytchos', 0.566596269607544),
- ('ocphy', 0.5415685176849365),
- ('dolchedy', 0.5343093872070312),
- ('aiicthy', 0.5323750376701355),
- ('odchecthy', 0.5235849022865295),
- ('okeeos', 0.5187858939170837),
- ('cphocthy', 0.515974388694763),
- ('oteor', 0.5050544738769531)]

# „ASTRONOMICAL WORDS"

# ВВЕДЕНИЕ: ЧТО ТАКОЕ БИОИНФОРМАТИКА

- математические методы компьютерного анализа генома, транскриптома, протеома (омикс- биоинформатика).
- разработка алгоритмов и программ для предсказания пространственной структуры биополимеров– РНК и белок - структурная биоинформатика ~ ФОЛДНИНГ
- [1]моделирование белковых каскадов,предсказание функции белка, регуляторных контуров и т.

# SHOTGUN & NEXT GEN. SEQUENCING

| Strand | Sequence |
| --- | --- |
| Original | AGCATGCTGCAGTCATGCTTAGGCTA |
| First shotgun sequence | AGCATGCTGCAGTCATGCT--------------------------TAGGCTA |
| Second shotgun sequence | AGCATG------------------------CTGCAGTCATGCTTAGGCTA |
| Reconstruction | AGCATGCTGCAGTCATGCTTAGGCTA |

# ПРИМЕР БЕЛКОВОЙ ПОСЛЕДОВАТЕЛЬНОСТИ

- Обнаружение внутривидового и межвидового полиморфизма.
- Таксономия
- Молекулярные часы

# ОСНОВНАЯ СТАТЬЯ

- Continuous Distributed Representation of Biological Sequences for Deep Proteomics and Genomics
- Ehsaneddin Asgari,Mohammad R. K. Mofrad
- PLOS ONE November 10, 2015
- https://doi.org/10.1371/journal.pone.0141287

α2 loop
α2-β8 loop

PPIA

PPIE

PPIC

PPIG

PPWD1

PPIL2

NKTR

SDCCAG-10

RANBP2*

PPIL6*

PPIL4*

СЕМЕЙСТВА БЕЛКОВ

24

# РАЗБИВКА БЕЛКОВОЙ ПОСЛЕДОВАТЕЛЬНОСТИ

Original Sequence

$(1)\vec{M}\,(2)\vec{A}\,(3)\vec{F}\,SAEDVLKEYDRRRRMEAL..$

Splittings

1) MAF, SAE, DVL, KEY, DRR, RRM, ..
2) AFS, AED, VLK, EYD, RRR, RME, ..
3) FSA ,EDV, LKE, YDR, RRR, MEA, ..

# РАСПРЕДЕЛЕНИЕ БЕЛКОВ В ПРОСТРАНСТВЕ 2Х КОМПОНЕНТ
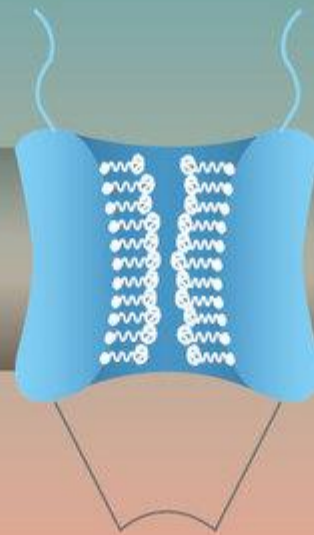# ЦВЕТ ОБОЗНАЧАЕТ ЗНАЧЕНИЕ СООТВ. ПРИЗНАКА
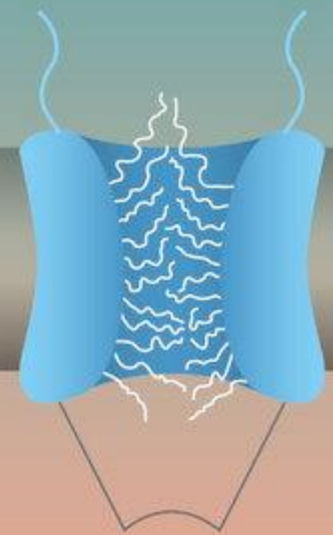
# PHENYLALANINE-GLYCINE NUCLEOPORINS (FG-NUPS)

# ПОСЛЕДОВАТЕЛЬНОСТИ FG-NUP VS СТРУКТУРИРОВАННЫЕ БЕЛКОВЫЕ ПОСЛЕДОВАТЕЛЬНОСТИ

# РЕАЛИЗАЦИЯ

- https://github.com/peter-volkov/biovec
- https://github.com/ehsanasgari/Deep-Proteomics

# Electrocorticography (ECoG)

**Electrodes fo epidural recordings.**

A)

3D trajectory of the hand in time $(x, y, z)(t)$

Brain signals in time $ECoG(t)$

B)

$$Y = \begin{bmatrix} y(t_1) \\ y(t_2) \\ \vdots \\ y(t) \\ \vdots \\ y(t_n) \end{bmatrix} \qquad \underline{X} = \begin{bmatrix} \underline{x}(t_1) \\ \underline{x}(t_2) \\ \vdots \\ \underline{x}(t) \\ \vdots \\ \underline{x}(t_n) \end{bmatrix}$$

Multiway ECoG Data

Artifacts filtration

Decimation along the temporal modality

# СПАСИБО ЗА ВНИМАНИЕ!