

---

*NLP && production*

# Эволюция задач и алгоритмов на текстах

Штех Геннадий @ NAUMEN • 07.04.2018

---

# Задачи бизнеса

## Введение

1. **Введение:** что такое НЛП, какие задачи решает?
2. **Куда идёт индустрия** и какие неочевидные задачи из этого возникают?

## Секция реального опыта

1. **Сфера информационной доступности:** польза для бизнеса, науки, законодательства

# Алгоритмы НЛП

## Что происходит в мире

1. **Актуальные классические подходы:** обработка текста, topic modeling, полнотекстовый поиск
2. **Развитие методов:** дистрибутивная семантика, \*2вес, нейронные сети

## Что исследуем мы

1. **Информационный поиск**
2. **Семантика документов**
3. **Группировка документов**
4. **Рекомендательные системы**

# Инструментарий

## Требования к инструментам

1. Классические проблемы NLP in production
2. Какие решения применяли мы

## Очень полезные ссылки

1. Тезисный обзор открытых проектов, о которых надо знать NLP-инженеру

---

# Задачи

*Историческая ретроспектива*

---

## Список

1. Поиск
2. Машинный перевод
3. Извлечение информации: заполнение форм
4. Распознавание речи
5. Кластеризация, каталогизация, автотеги́рование
6. Классификация, сентимент

# Поиск

1. Запросная строка с полнотекстовым (Wiki)
2. Четкие фильтры (по дате, количеству цитат, исключающие какие-то слова)

## Машинный перевод

1. Помощь в составлении словарей
2. Статистический перевод фраз
3. Перевод документов/интерфейсов



## Извлечение информации

1. Извлечение событий из потока новостей по шаблону
2. Автоматизированное заполнение форм, описаний, номенклатуры
3. Составление “карты знаний” с отношениями на объектах или понятиях “объект1-отношение-объект2”

## Распознавание речи

1. Речевой ввод aka “Диктограф”
2. Автоматизация колл-центров
3. Субтитры и синхронный перевод

## Задачи на текстах

### Кластеризация

1. Маркетинговые исследования отзывов
2. Информационная разведка

### Каталогизация

1. Составление банков знаний
2. Структуризация отчетности

## Задачи на текстах

### Классификация

1. Спам
2. Маршрутизация обращений
3. Оценка тональности отзыва

### Автотегирование

1. Сортировка потока новостей
2. Краткая аннотация темы содержания

## **Куда двигается индустрия**

1. Улучшение существующих техник методом “всё становится лучше с нейросетями”
2. End-to-End подходы (благодаря нейросетям, опять же)
3. Улучшение интерфейсов: взрыв “чат-ботов”
4. Q&A

## **Куда двигается индустрия**

1. Стремление “спрямить” потоки информации
2. Комбинации методов
3. Больше контекста в поисковых сессиях: “диалоговые системы”

## Почему я так думаю

Или список решенных  
Naumen задач

1. Семантический поиск документов
2. Группировка и каталогизация научно-технических библиотек
3. Поиск заимствований и повторов
4. Диалоговые интерфейсы

---

# Алгоритмы

*Историческая ретроспектива и справка*

---



## О типичном NLP-конвейере

### Предобработка текста (нормализация)

- Токенизация, стемминг/лемматизация
- “Красненькая шапочка” -> ['красный', 'шапка']

### Представление

- one-hot, Bag of Words, tf-idf, nPMI, hashing trick

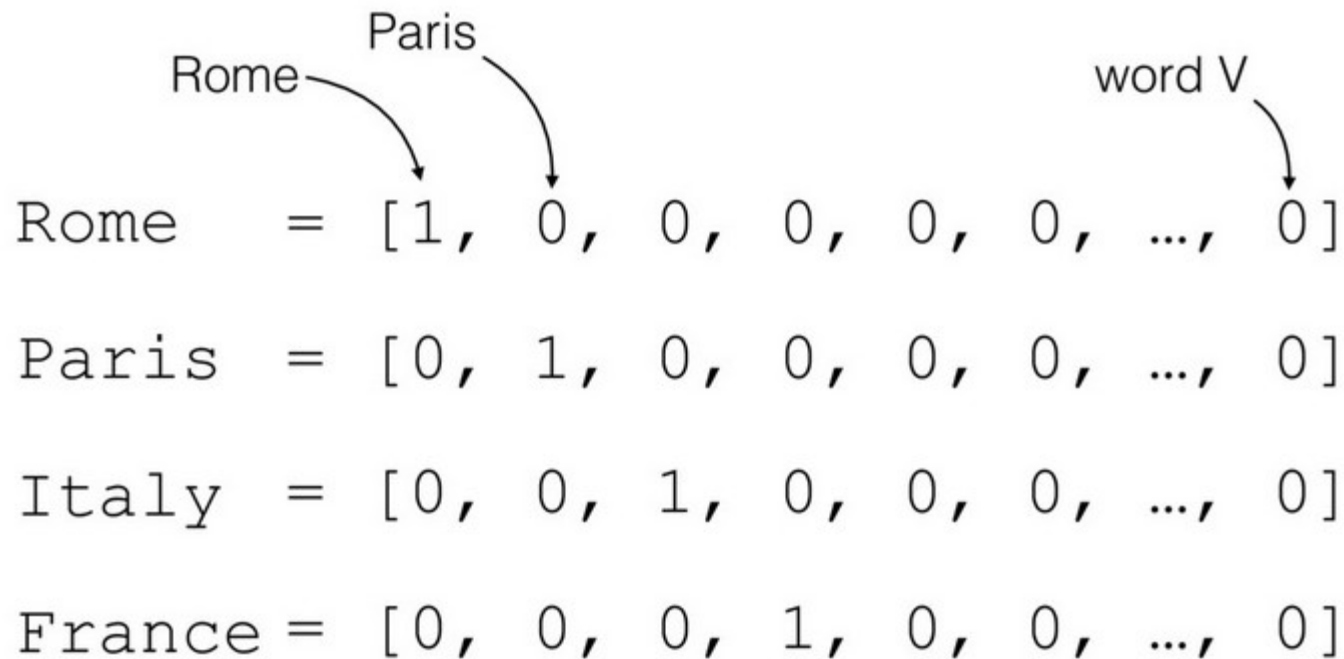
### Трансформация

- Matrix decomposition, topic modeling, embeddings

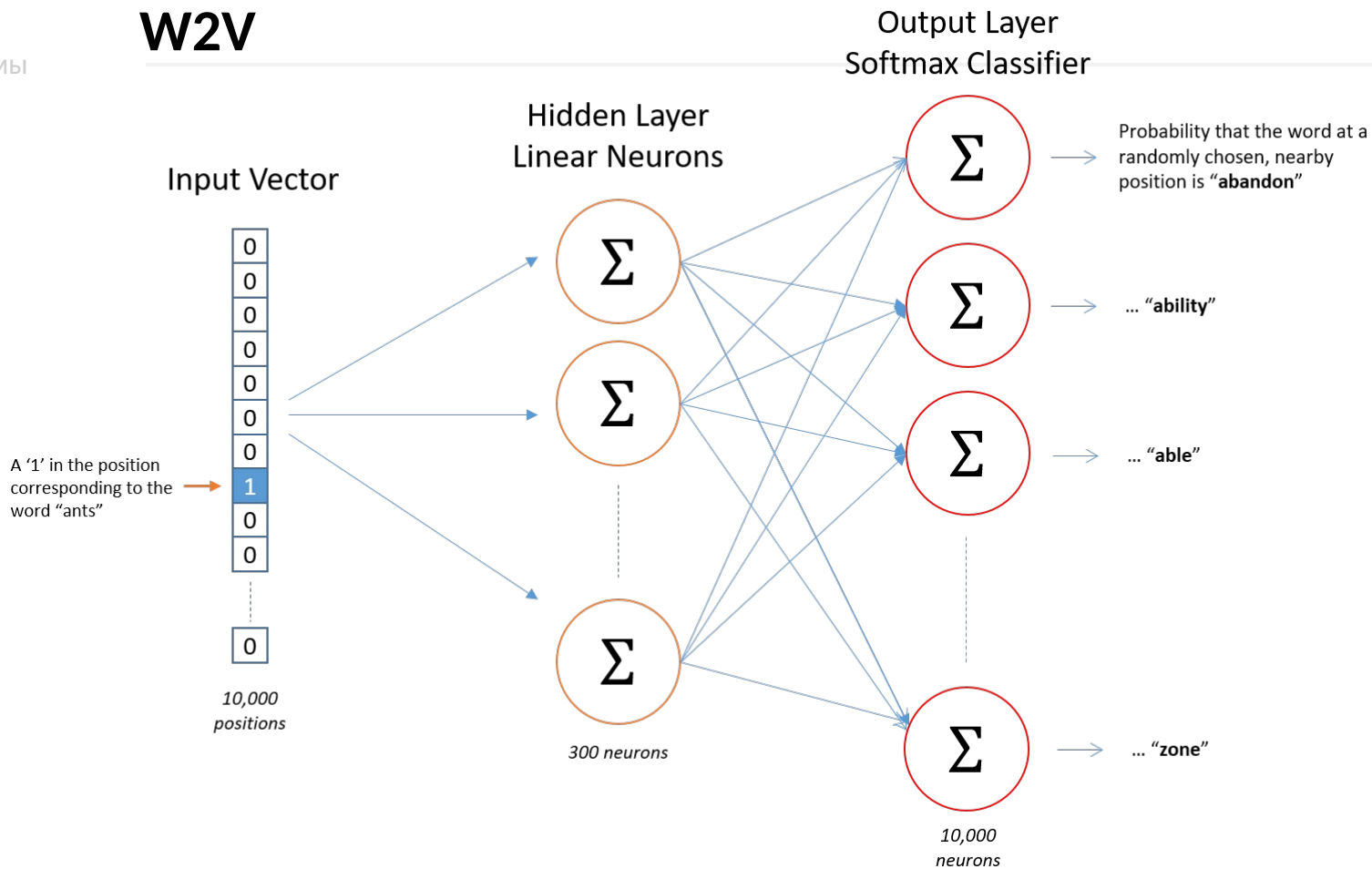
### Конечный алгоритм для решения задачи

- Классификация, кластеризация, регрессия, etc

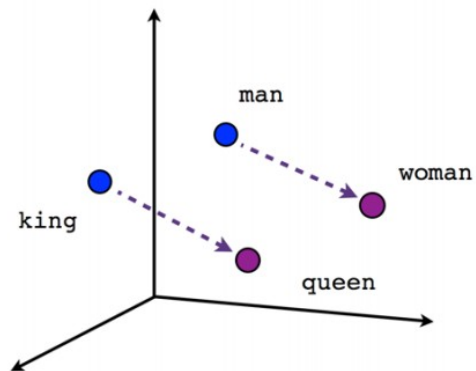
# One-Hot Encoding



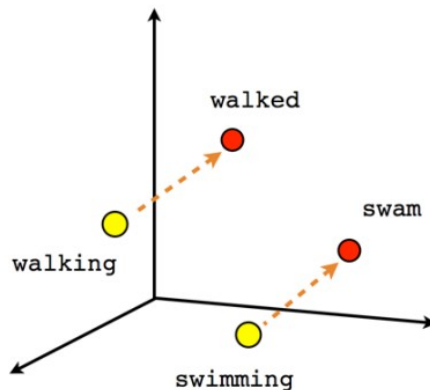
# W2V



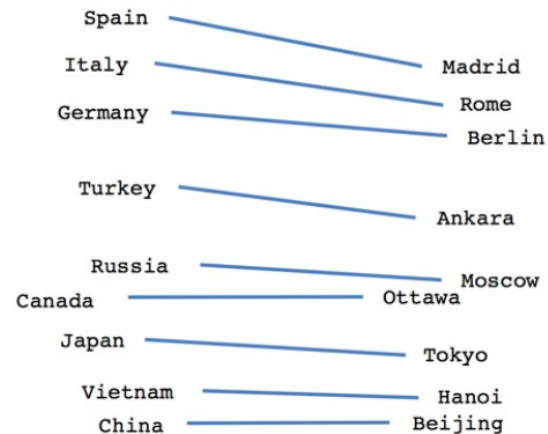
# Embeddings



Male-Female



Verb tense



Country-Capital

# Тематическое моделирование

Алгоритмы

Topics

Documents

Topic proportions & assignments

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

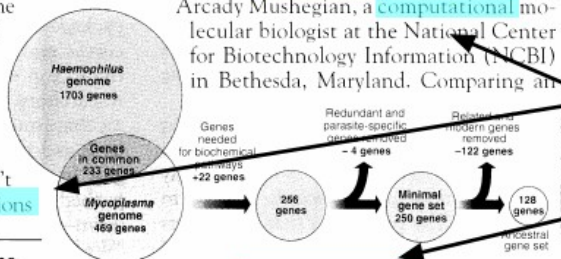
data 0.02  
number 0.02  
computer 0.01  
...

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

## Пример тематической модели

#106: "gpu" + "performance" + "core" + "cpu" + "hardware" + "cuda" + "than" + ...

#107: "state" + "electron" + "quantum" + "algorithm" + "method" + "functional" + ...

#108: "square" + "circle" + "least" + "rectangle" + "now" + "find" + "like" + ...

#109: "correction" + "kappa" + "r2" + "energy" + "conditional" + "crf" + "exchange" + ...

#110: "f" + "h" + "g" + "p" + " $\cdot$ " + "s" + "function" + ...

#111: "method" + "system" + "solve" + "problem" + "solver" + "equation" + "linear" + ...

#112: "state" + "action" + "reward" + "q" + "agent" + "policy" + "value" + ...

...

---

# Ветки развития

*Алгоритмы*

---

## Направления исследований

### Специализированные и мультимодальные эмбединги

01.

**Часто нужно уметь строить** не только эмбединги слов и документов: изображения, метаданные, пользователи

02.

**Иногда нужно захватывать структурные свойства:** иерархии смыслов, иерархии документов



## Направления исследований

### Вероятностные методы, решающие вспомогательные задачи

#### 01.

**Вероятностные модели требуют меньше данных**, легко интерпретируемы, позволяют пользователям воздействовать на работу алгоритма

#### 02.

Иногда в терминах вероятностных моделей **легче формулировать задачи**, например DPP позволяет добиться оптимального баланса разнообразия и релевантности выдачи

---

# Направления исследований

## Разрешение синонимии и омонимии, больше лингвистики

Качественные эмбеддинги слов,  
“понимающие” лингвистические  
особенности в конечном итоге  
улучшают работу на всех задачах

---

# State-Of-The-Art

*Алгоритмы*

---

## SOTA RNN

### LSTM + Attention

**Рекуррентные сети могут “забывать” важные детали** или с трудом решать задачи, которые плохо похожи друг на друга

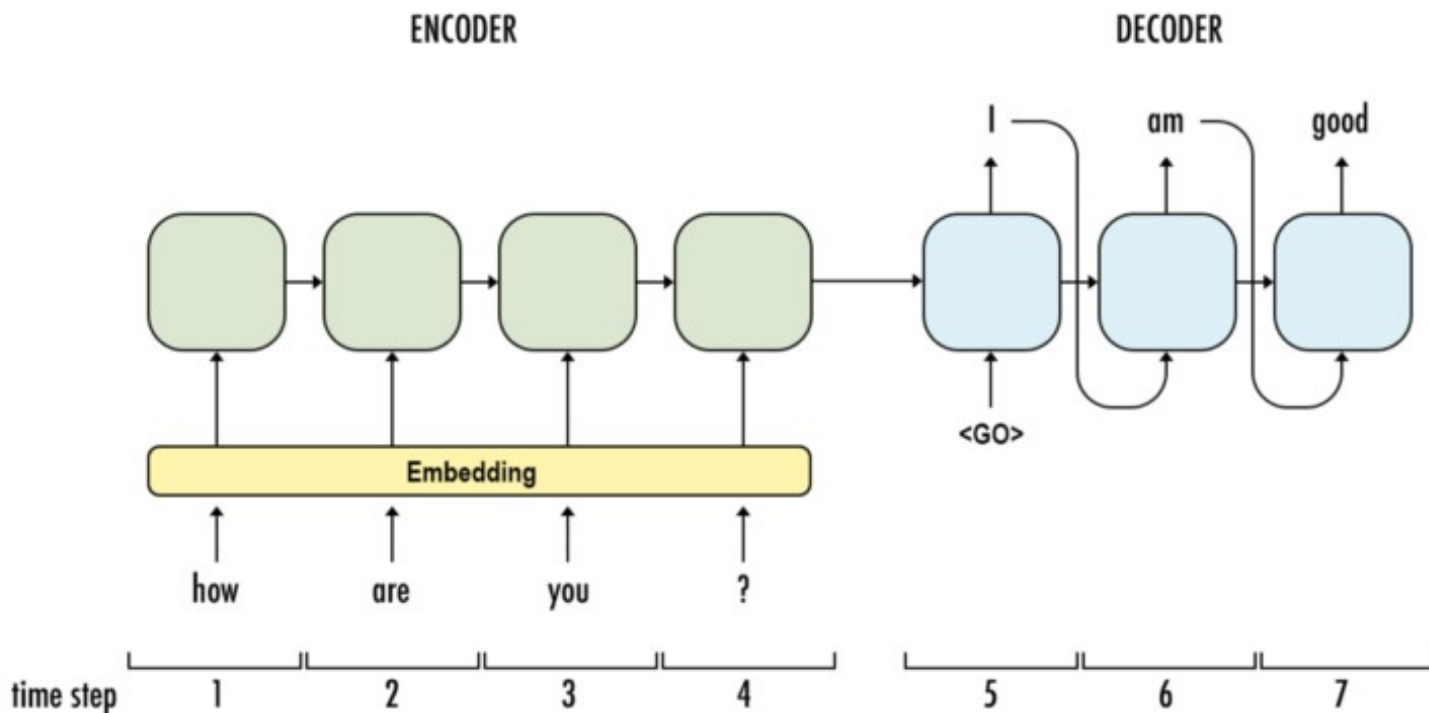
**Сети с вниманием способны “вспомнить”** подходящие части последовательности, даже если они были далеко, причем именно те, которые соответствуют текущей задаче

Различные модификации внимания **можно подобрать под задачу**

Некоторые модификации очень **похожи на свёрточные сети**

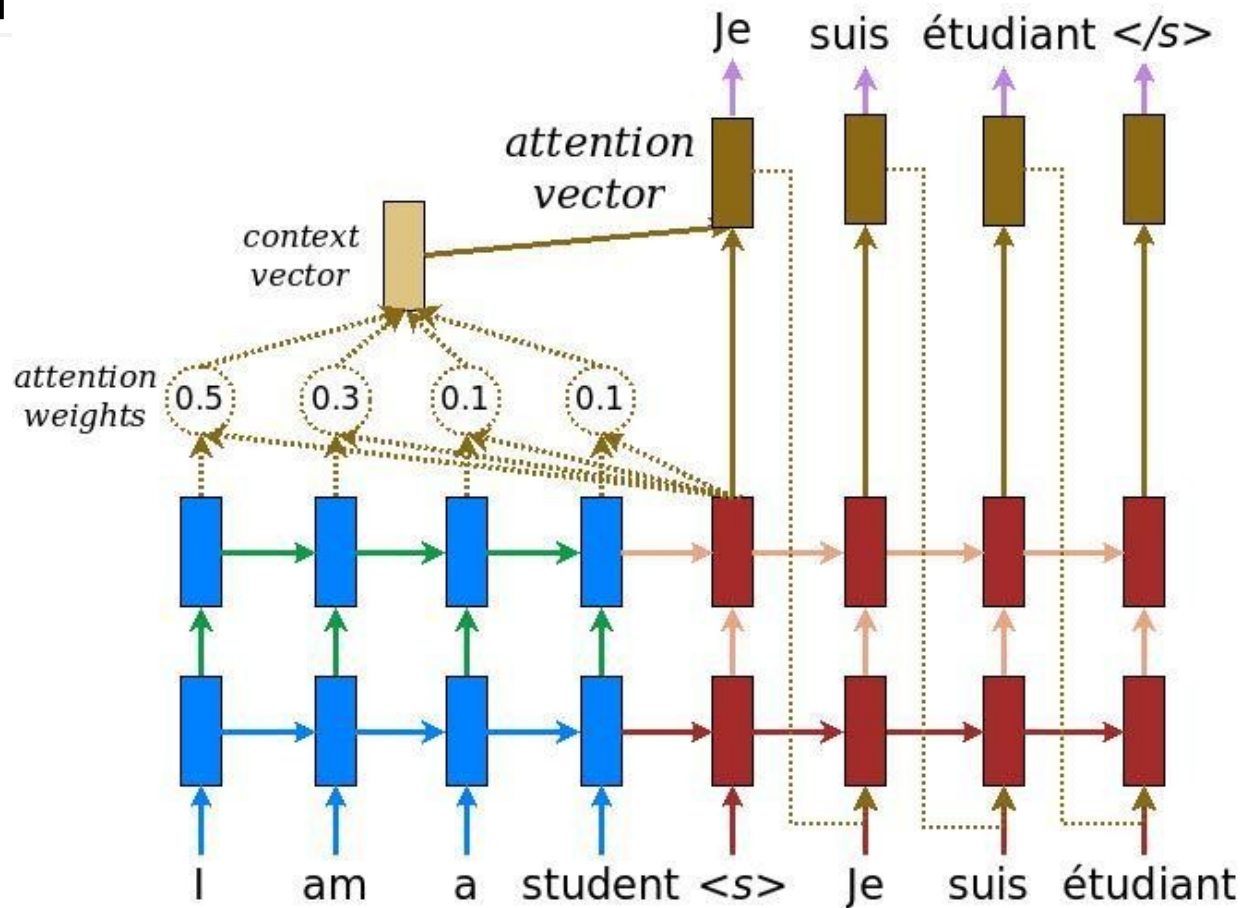
# LSTM

Алгоритмы



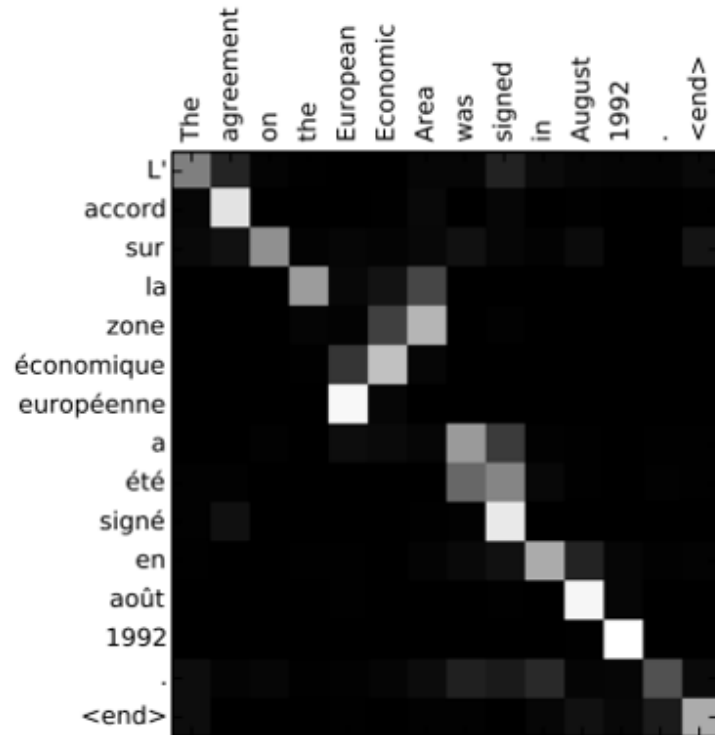
# ATTN

— Алгоритмы

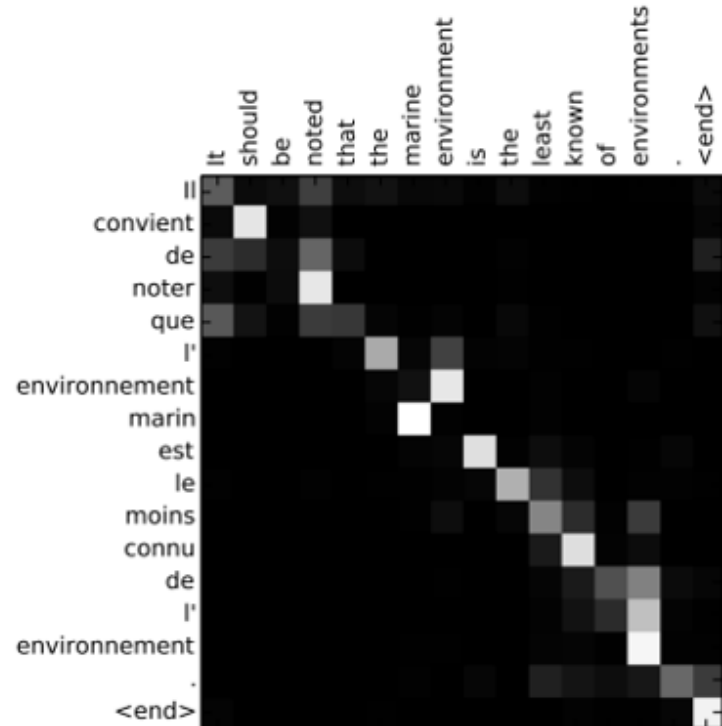


# Attention on translation task

Алгоритмы



(a)



(b)

## SOTA RNN

### Pointer Network

01.

Позволяют выбирать  
наилучший набор  
объектов из  
перечисленных в  
нужном порядке

02.

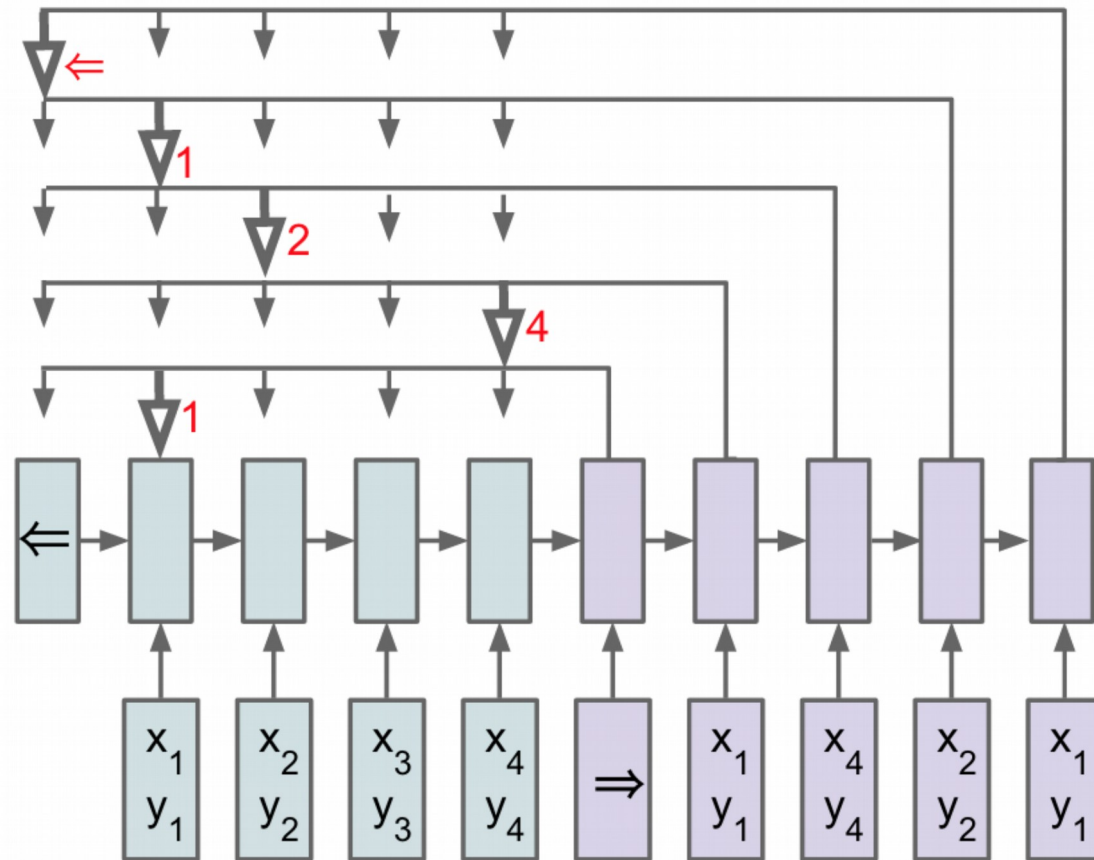
Впервые были  
применены к  
задачам  
**Коммивояжера**  
(графы) и поиска  
выпуклой  
оболочки(геометрия)

03.

Впервые были  
применены к задачам  
**Коммивояжера** (графы)  
и поиска выпуклой  
оболочки(геометрия)



# SOTA RNN



# SOTA Embeddings

## Fasttext

01.

Может работать с  
**минимальной**  
**предобработкой** текста

02.

**Способен строить эмбединги даже**  
**неизвестных слов** по аналогии  
орфографической структуры

# SOTA Embeddings

## Starspace

01.

**Является скорее нетривиальным применением** почти классических эмбедингов в комплекте с эффективной реализацией

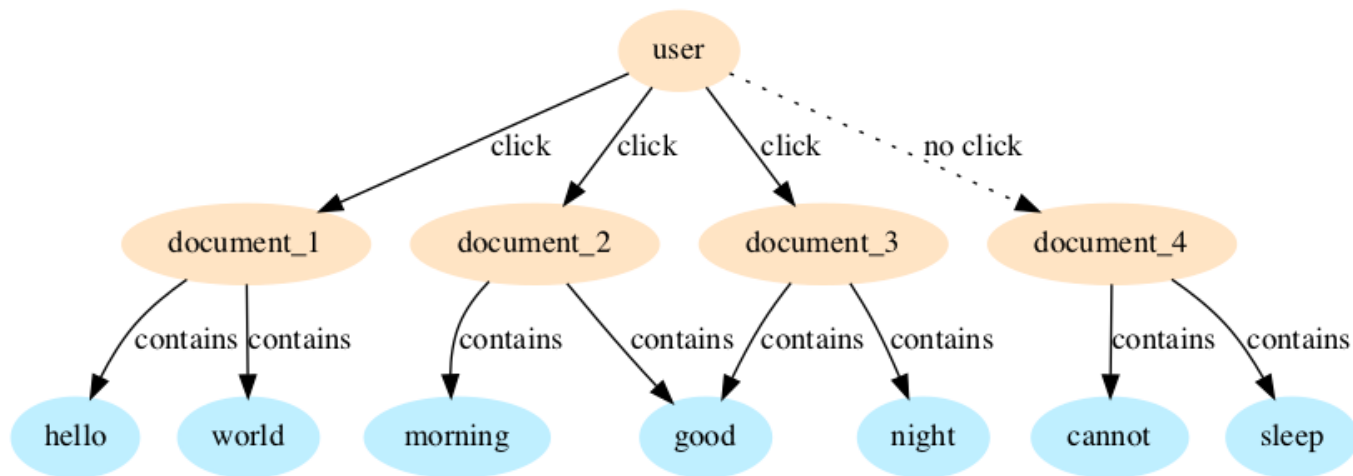
02.

**Мультимодальные эмбединги**, строит представление для всех элементов документа

03.

**Можно строить карты знаний** “понятие1-отношение-понятие2”

## Starspace



# SOTA Embeddings

## Poincaré

01.

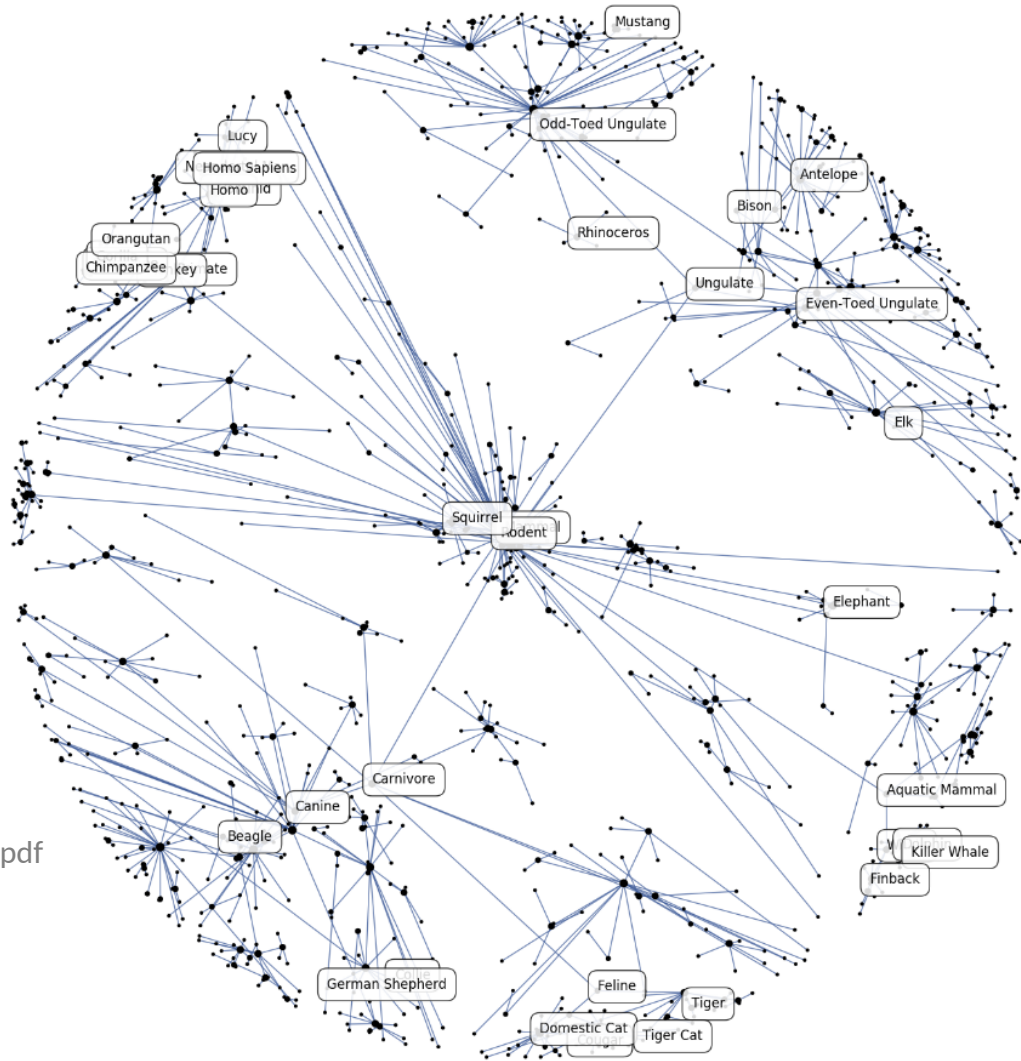
Позволяет захватить иерархическую структуру понятий “понятие1-обобщает-понятие2”

02.

Работает в нелинейных пространствах Poincaré Balls

# Poincaré Embeddings for Learning Hierarchical Representation

<https://arxiv.org/pdf/1705.08039.pdf>



---

# Наши направления

*Алгоритмы*

---

## Сегментация документов

### ARTM Continuous Topic Regularizer

Для поиска по большим документам **полезно нарезать их на непрерывные сегменты**, раскрывающие определенную тему или говорящие об одном(и ровно одном) факте

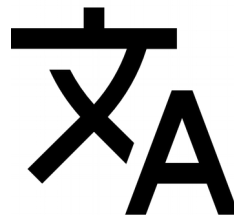




## Мультиязычные эмбединги

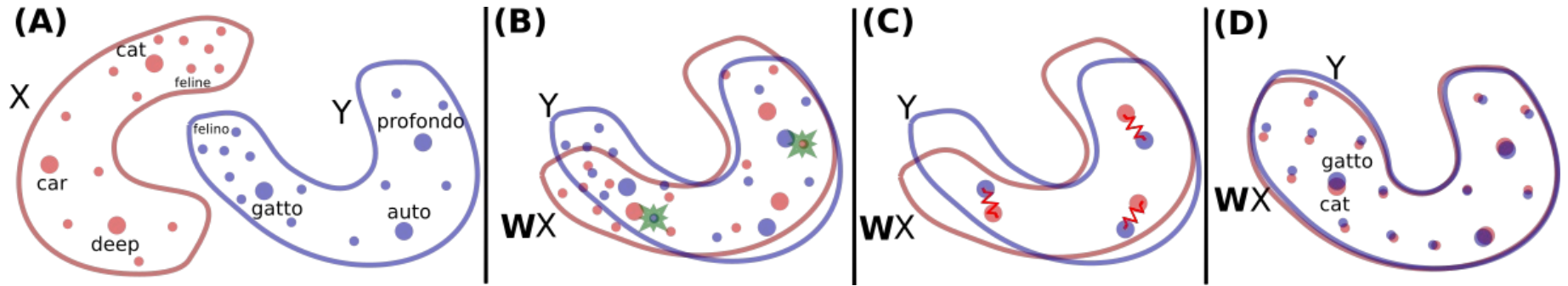
\*2vec

Геометрические свойства  
линейных эмбедингов  
позволяют **без большого  
количества данных получить  
общее пространство слов** для  
разных языков



# Мультиязычные эмбединги

— Алгоритмы



## Актуальные алгоритмы

### Представление

- Tf-idf, nPMI, hashing trick

### Факторизация (декомпозиция)

- PCA, LSI, LSA, pLSA

### Тематическое моделирование

- pLSA, LDA, HDP, ARTM

### Поиски

- BM25, HNSW, LSH

### Эмбединги

- word2vec, glove, paragraph2vec, fasttext, starspace, poincaré

### Нейросетевые подходы

- LSTM, GRU, Attention, siamese network, similarity learning

## Полезный NLP-софт

### Предобработка текста (нормализация, токенизация)

- pymorphy2(ru), snowball stemmer(en), Stanford NLP(en)

### Фреймворки

- sklearn, NLTK, gensim, spaCy

### Узкоспециализированные фреймворки

- BigARTM, Vowpal Wabbit, Fasttext, Starspace, Muse, faiss, annoy, NMSLib, lucene, sphinx, elastic

### Нейросетевые фреймворки

- Pytorch, Keras

---

# Контакты

**Штех Геннадий \***  
**@ NAUMEN**

[gshtekh@naumen.ru](mailto:gshtekh@naumen.ru)

**Gennady Shtekh**

[shtechgen@gmail.com](mailto:shtechgen@gmail.com)

[t.me/sht3ch](https://t.me/sht3ch)

[github.com/ShT3cH](https://github.com/ShT3cH)

\*R&D Data Usage Department Executive

---