

# Текстовая аналитика

Школа прикладного анализа данных Data-Diving



Университетский  
консорциум  
исследователей  
больших данных

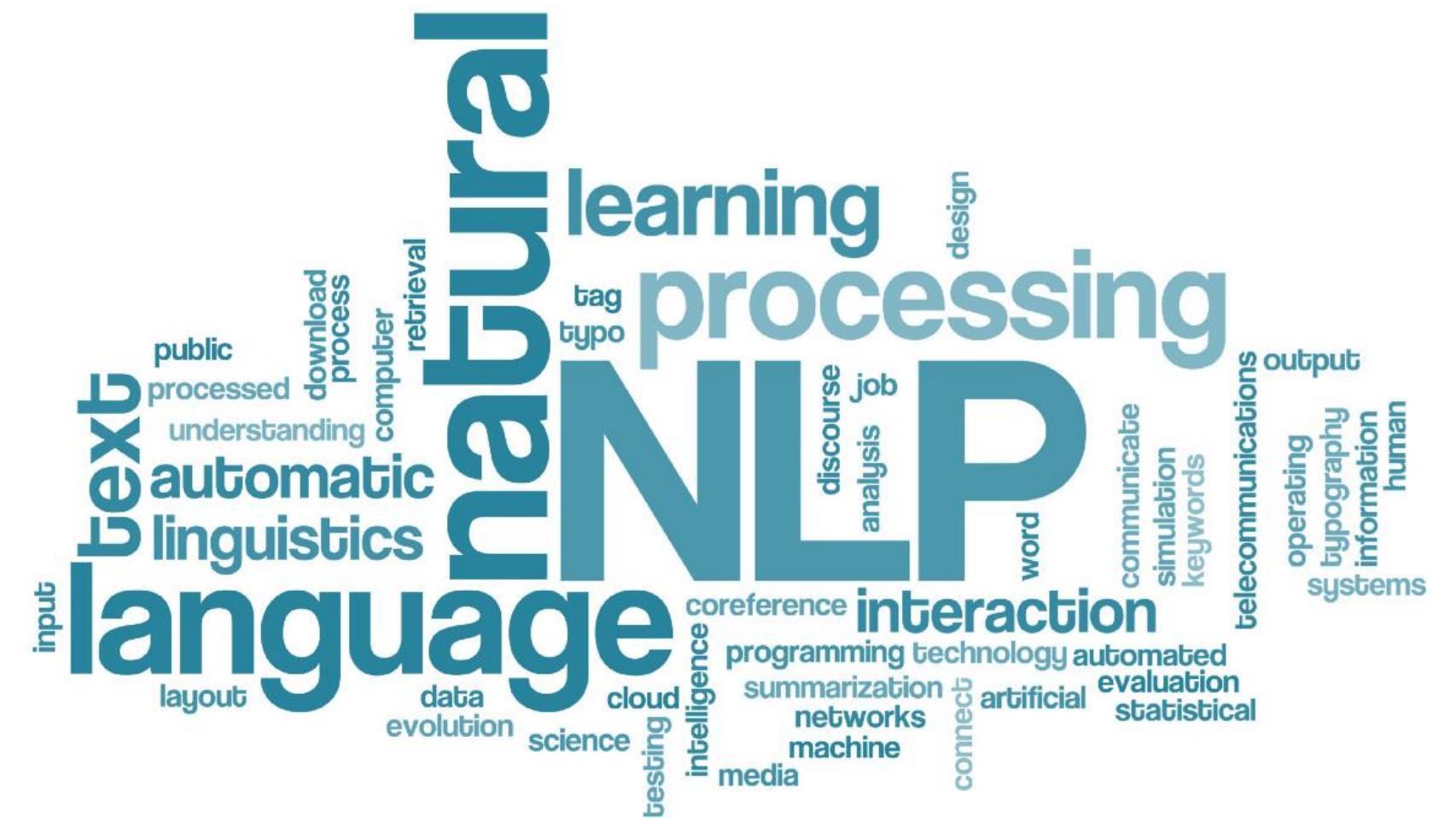
**Басина Полина**

аналитик Центра прикладного  
анализа больших данных

# Что такое обработка естественного языка (NLP)?

**Обработка текстов на естественном языке**  
(*Natural Language Processing, NLP*)

общее направление искусственного интеллекта и математической лингвистики, направленное на изучение методов анализа и синтеза текстов на естественных языках



 Распознавание речи

 Анализ текста

 Синтез речи

 Генерирование текста

Встречались ли Вы в своей  
повседневной жизни с  
примерами использования  
обработки естественного  
языка?

# Популярные задачи

## Формирование ответов на вопросы (Question Answering)



Скажите  
«Привет, Алиса»

## Анализ эмоциональной окраски высказываний



## Машинный перевод

Яндекс Переводчик

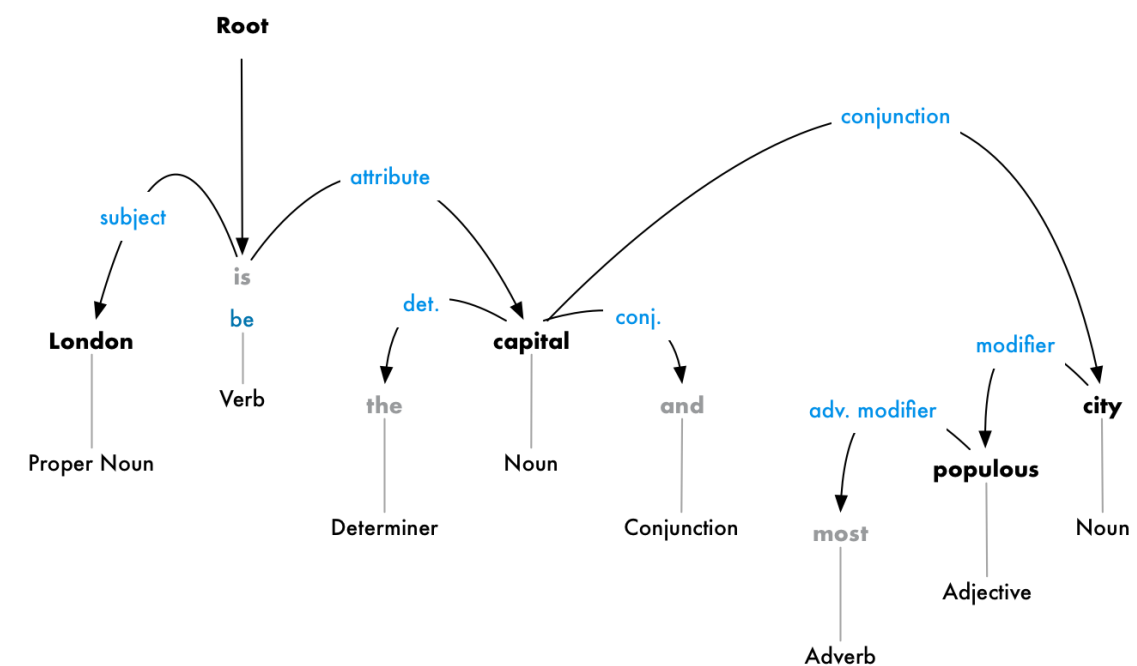
Текст Сайт

АНГЛИЙСКИЙ ↔ РУССКИЙ

great deal of

отличный интернет

## Частеречная разметка

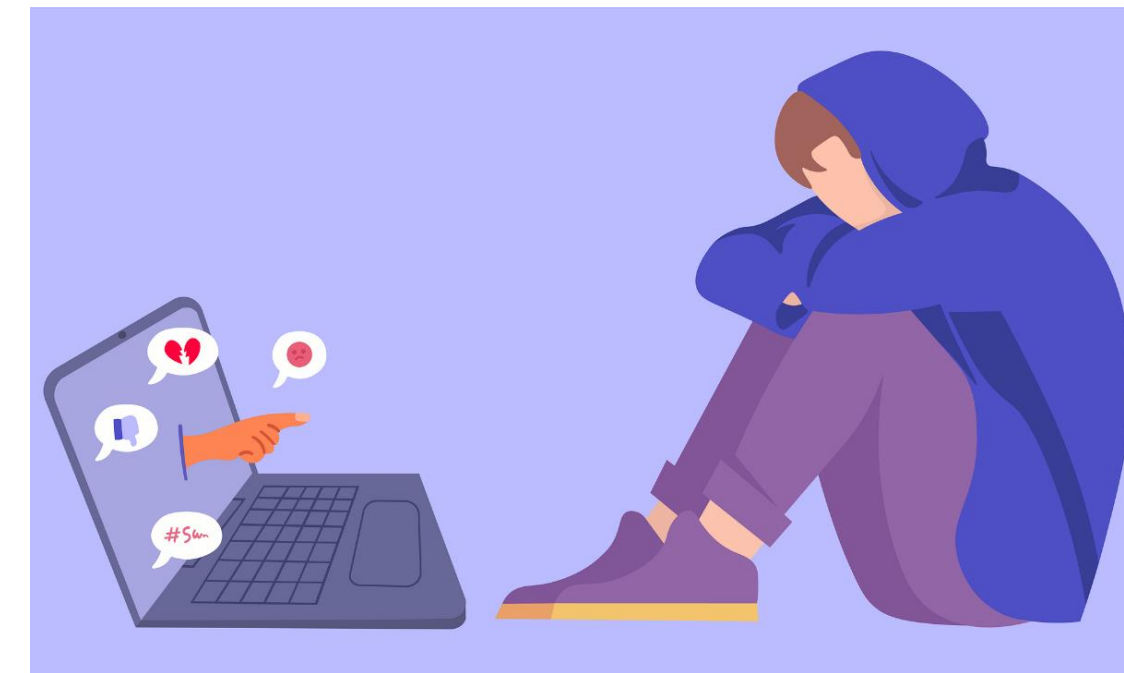


# Популярные задачи

Определение фальшивых новостей



Определение токсичных высказываний

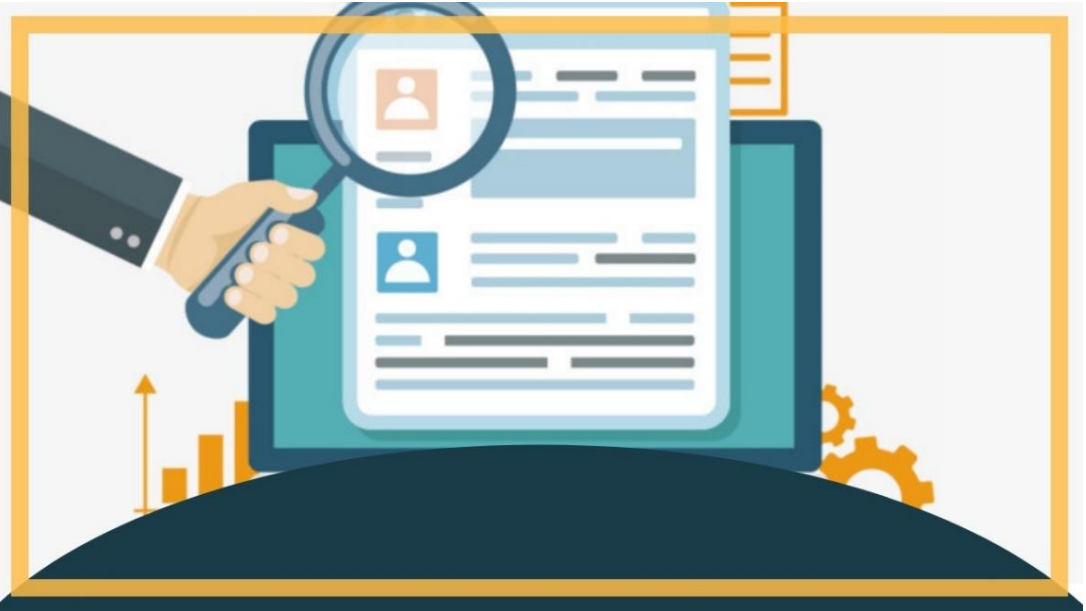


Классификация текстов по различным категориям

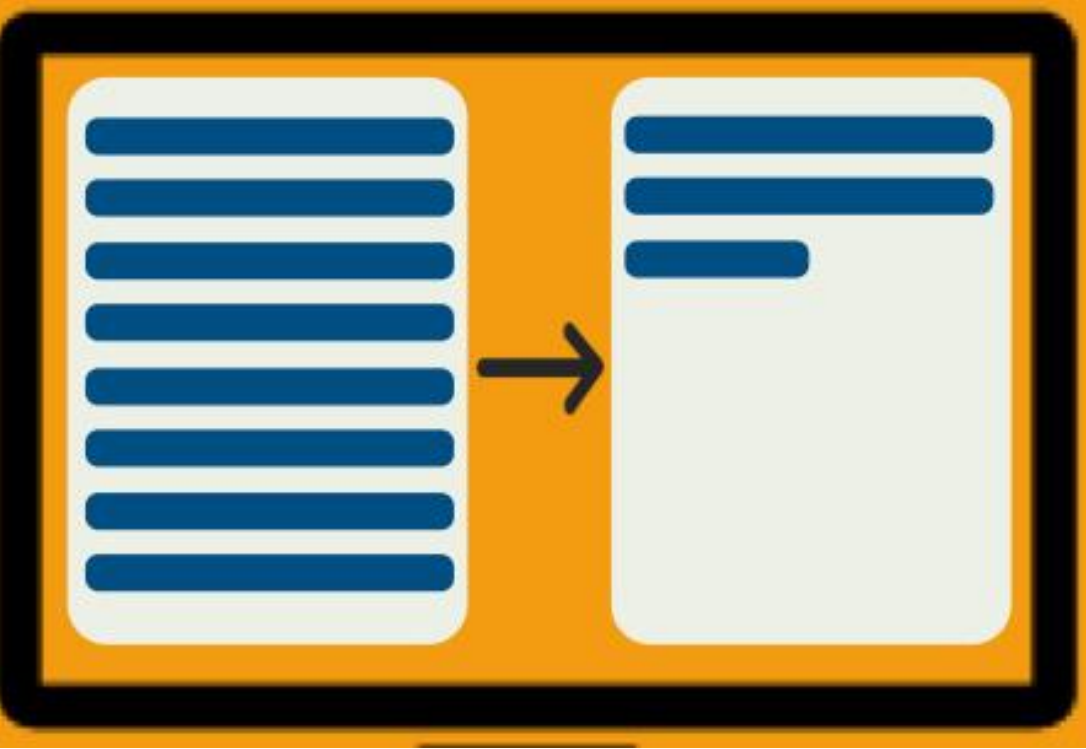


# Популярные задачи

## Поиск дубликатов



## Реферирование текста



## Извлечение фактов




## Извлечение сущностей


Руководитель НИИ транспорта и дорожного хозяйства Михаил Блинкин утверждает, что может помочь Сергею Собянину решить транспортные проблемы столицы, в частности разгрузить МКАД. «Сама постановка задачи — вернуть МКАД исходное значение — мне чрезвычайно импонирует. Я все последние годы писал и говорил о том, что МКАД превратилась у нас в просёлочную дорогу, подъезд к магазинам, даже в городскую улицу межквартальную», — цитирует РБК Блинкина.


Закончить разметку абзаца


СПАНЫ	УПОМИНАНИЯ
✗ [НИИ] [НИИ транспорта и дорожного хозяйства]	Org
✗ [Михаил] [Блинкин]	Person
✗ [Сергею] [Собянину]	Person
✗ [МКАД]	Location
✗ [МКАД]	Location
✗ [МКАД]	Location
✗ [РБК]	Org
✗ [Блинкина]	Person

# Основные термины

 **Токен** – текстовые единицы  
*Символы, слова, словосочетания и предложения.*

 **Документ** – совокупность токенов, которые принадлежат одной смысловой единице  
*Предложение, комментарий или пост пользователя.*

 **Словарь** – совокупность всех токенов, встречающихся в корпусе текстов.

 **Корпус** – генеральная совокупность всех документов.

**Документ 1**  
«Сегодня мы едем на дачу»

**Документ 2**  
«Мы поехали на речку»

**Корпус**

**Токен** – слово

**Словарь:** сегодня, мы, едем, на, дачу, поехали, речку

# Поработаем самостоятельно...

## Документ 1

«Мы занимаемся  
обработкой языка»

## Документ 2

«Язык сложно  
анализировать»

Корпус

Токен – слово

Словарь - ?



# Поработаем самостоятельно...

## Документ 1

«Мы занимаемся  
обработкой языка»

## Документ 2

«Язык сложно  
анализировать»

Корпус

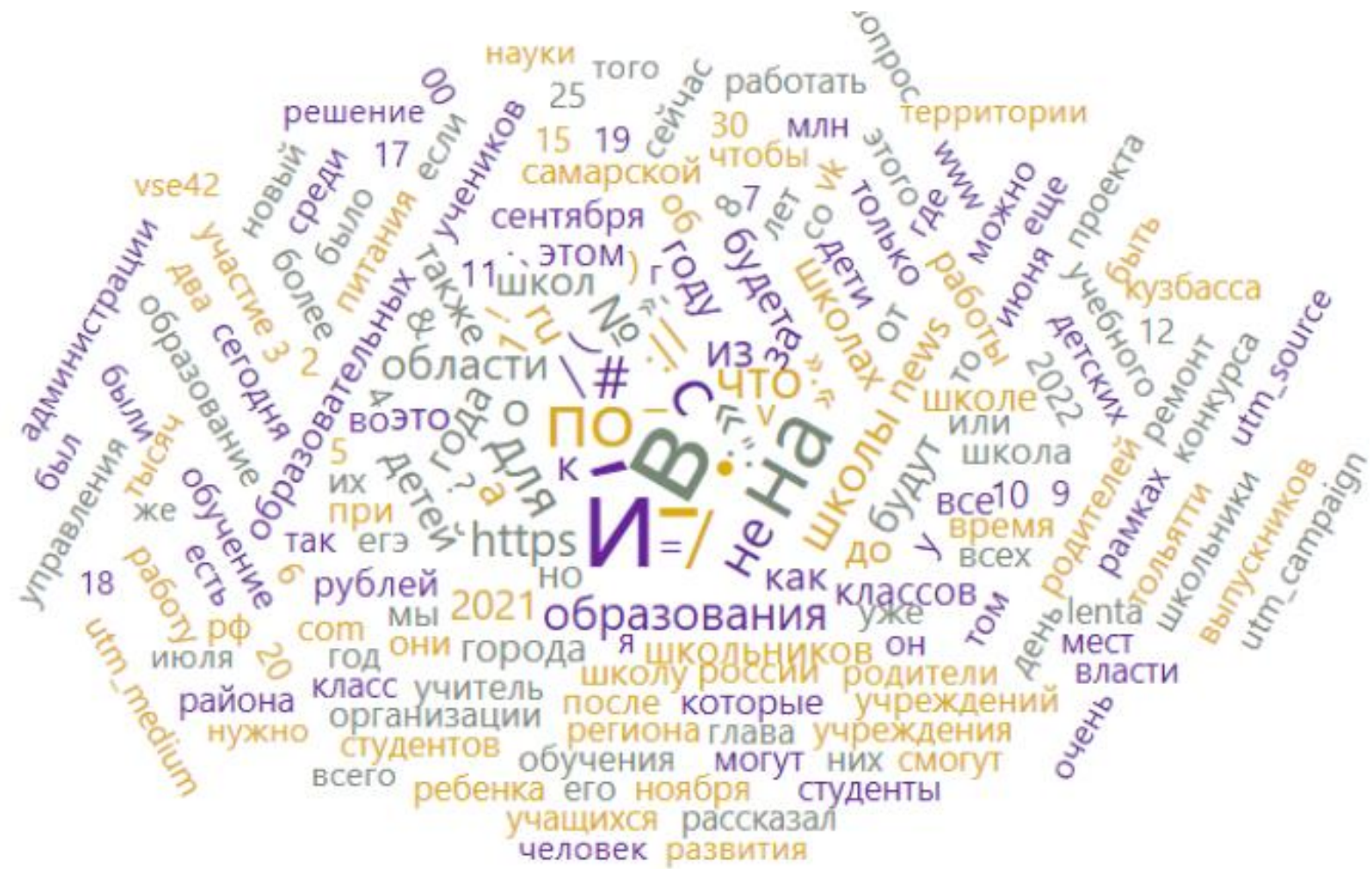
Токен – слово

Словарь:  
Мы, занимаемся,  
обработкой, языка,  
язык, сложно,  
анализировать

# Этапы текстового анализа данных



# Шаг 2. Предварительная обработка



## Шаг 3. Векторизация



### Прямое кодирование

- *вектор - слово*
- *признаки – токены из словаря*

Токен представлен бинарным вектором, который указывает наличие или отсутствие признака; единица ставится тому элементу, который соответствует номеру токена в словаре

**Словарь:** обработка, естественного, языка

обработка 

1	0	0
---	---	---

естественного 

0	1	0
---	---	---

языка 

0	0	1
---	---	---

# Попробуем сами...

## Прямое кодирование

**Словарь:** мама, мыла, раму

мама

--	--	--

мыла

--	--	--

раму

--	--	--

# Попробуем сами...

## Прямое кодирование

**Словарь:** мама, мыла, раму

мама	1	0	0
------	---	---	---

мыла	0	1	0
------	---	---	---

раму	0	0	1
------	---	---	---

## Шаг 3. Векторизация



### Мешок слов BoW

- вектор всего документа
- признаки – токены из словаря

Каждое число соответствует появлению или количеству появлений соответствующего слова в тексте

**Словарь:** обработка (1), естественный (2), язык (3), русский (4), английский (5), я (6), занимаюсь (7), нравится (8)

Документ 1:

Мне нравится  
заниматься обработкой  
естественного языка

	1	2	3	4	5	6	7	8
	1	1	1	0	0	1	1	1

Документ 2:

Мне нравится  
заниматься обработкой  
русского языка

	1	0	1	1	0	1	1	1
--	---	---	---	---	---	---	---	---

Документ 3:

Мне нравится  
заниматься обработкой  
естественного языка

	1	0	0	1	1	1	1	1
--	---	---	---	---	---	---	---	---

# Попробуем сами...



## Мешок слов VoW

**Словарь:** обработка (1), образовательный(2),  
выбинар (3), русский (4), английский (5), я (6),  
смотреть(7), любить(8)

### Документ 1:

Я люблю смотреть  
образовательные  
вебинары

1	2	3	4	5	6	7	8



# Попробуем сами...



## Мешок слов BoW

**Словарь:** обработка (1), образовательный(2), вебинар (3), русский (4), английский (5), я (6), смотреть(7), любить(8)

**Документ 1:**  
Я люблю смотреть образовательные вебинары

1	2	3	4	5	6	7	8
0	1	1	0	0	1	1	1

## Шаг 3. Векторизация

### TF-IDF

- *вектор всего документа*
- *признаки – токены из словаря*

Вес некоторого слова пропорционален частоте употребления этого слова в документе и обратно пропорционален частоте употребления слова во всех документах коллекции

**Словарь:** обработка (1), естественный (2), язык (3), данные (4), фотографий(5), я (6), занимаюсь (7), нравится (8)

#### Документ 1:

Мне нравится заниматься обработкой естественного языка

0.1	0.92	0.92	0	0	0.1	0.1	0.1
-----	------	------	---	---	-----	-----	-----

#### Документ 2:

Мне нравится заниматься обработкой фотографий

0.1	0	0	0	0.92	0.1	0.1	0.1
-----	---	---	---	------	-----	-----	-----


#### Документ 3:


Мне нравится заниматься обработкой данных


0.1	0	0	0.92	0	0.1	0.1	0.1
-----	---	---	------	---	-----	-----	-----

## Шаг 3. Векторизация

### Недостатки классических подходов

 отсутствие учета  
контекста и порядка слов

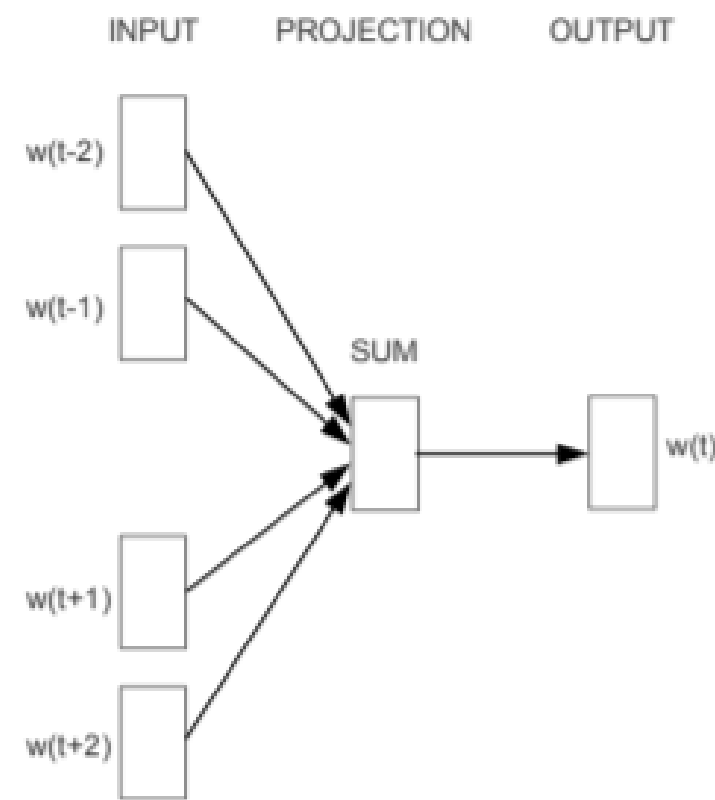
 невозможно эффективно  
представить новые слова,  
отсутствующие в словаре

 высокая размерность  
получаемых векторов

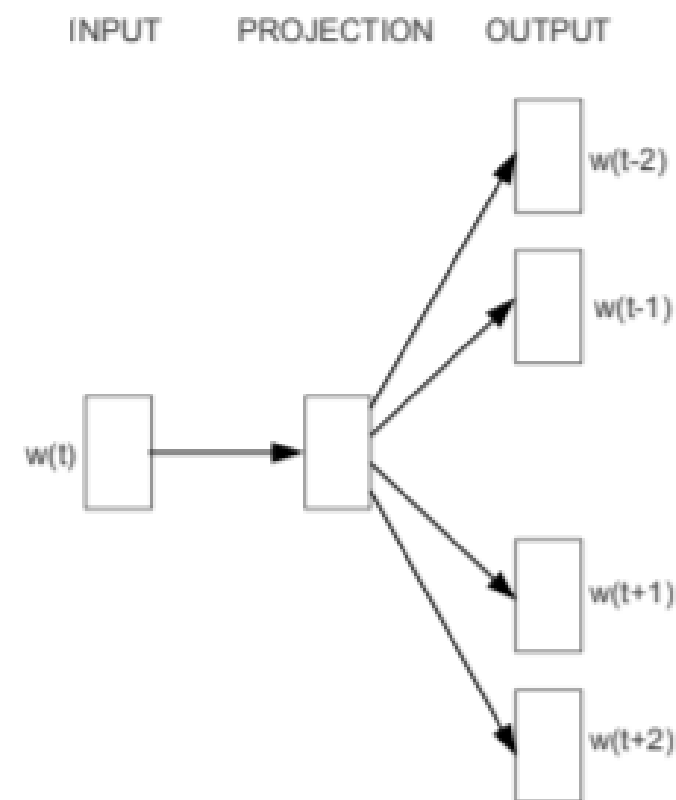
## Шаг 3. Векторизация

### 👉 Статические представления Word2vec, FastText

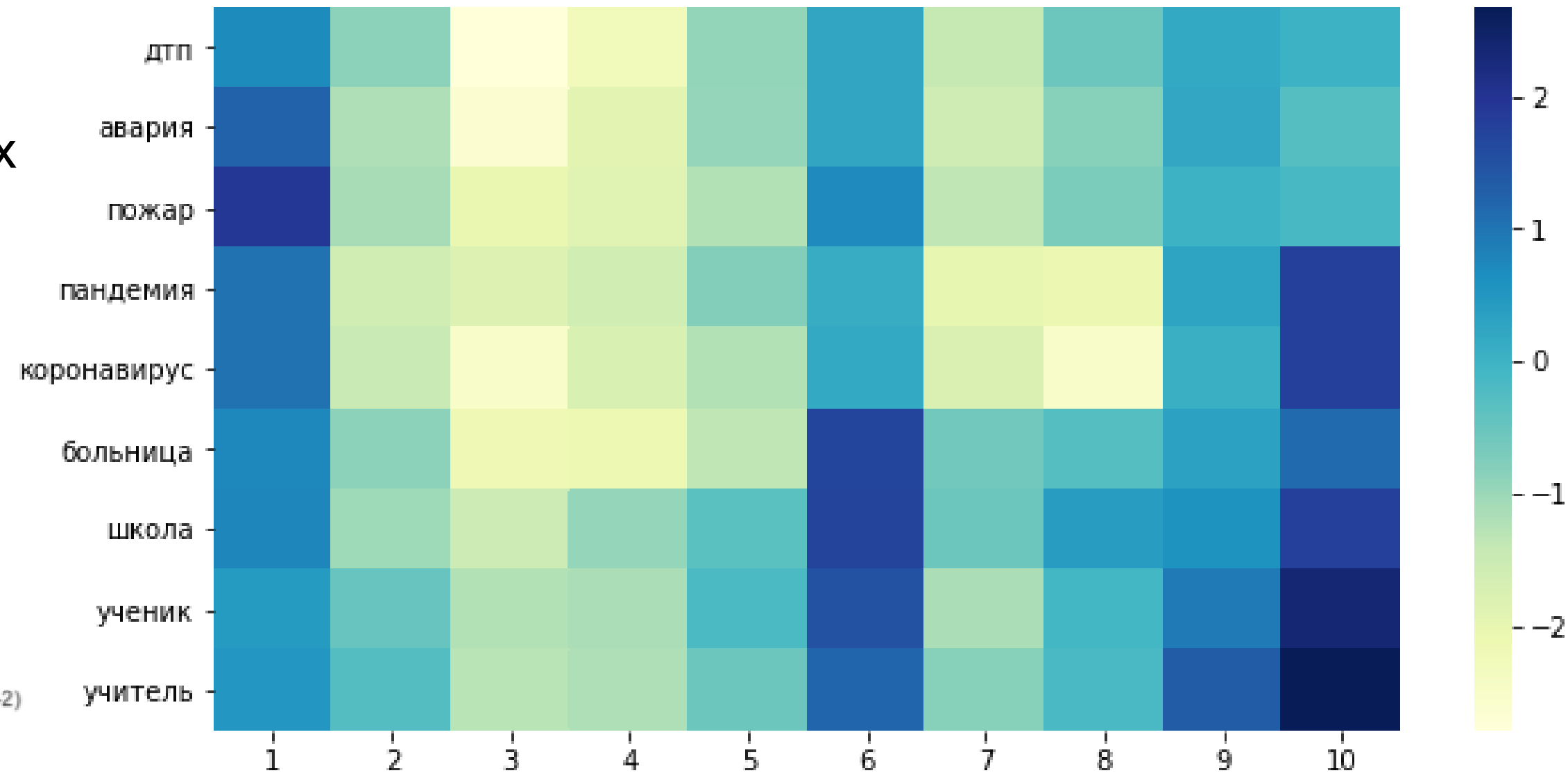
Статические модели построены на принципах дистрибутивной семантики, основная идея которых заключается в том, что значение определяется употреблением, а семантика может быть получена из контекстов, в которых употребляется данное слово



CBOW



Skip-gram



- *вектор – слово*
- *признаки – токены из словаря*

## Шаг 3. Векторизация

### Контекстуализированные вектора ELMo, RuBERT

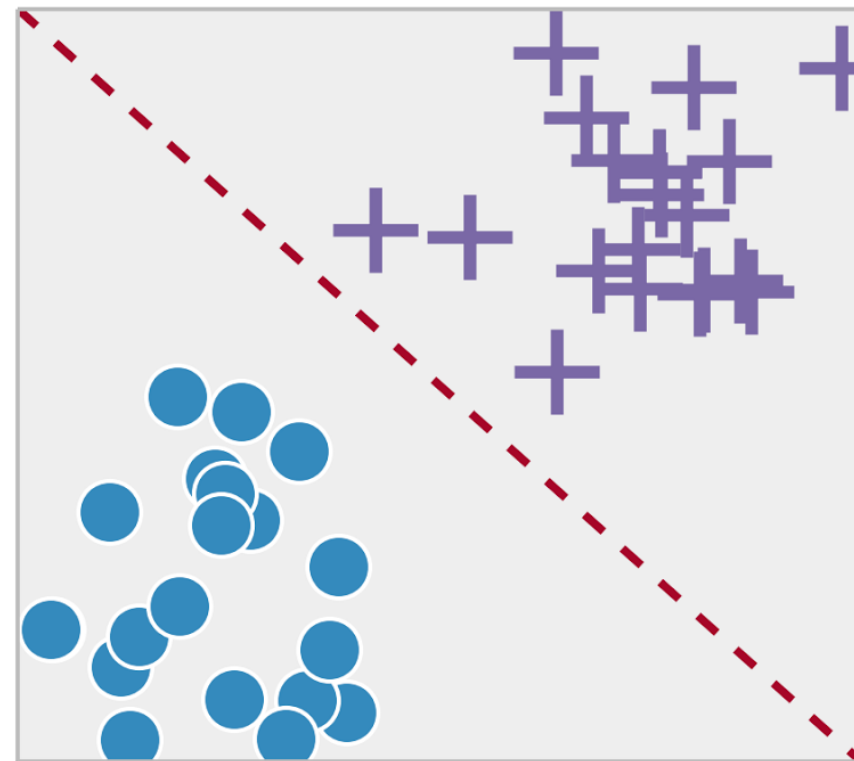
В разных контекстах одна и та же текстовая конструкция может означать различные понятия. В связи с этим, в последние годы получило развитие направление по формированию зависящих от контекста векторных представлений.



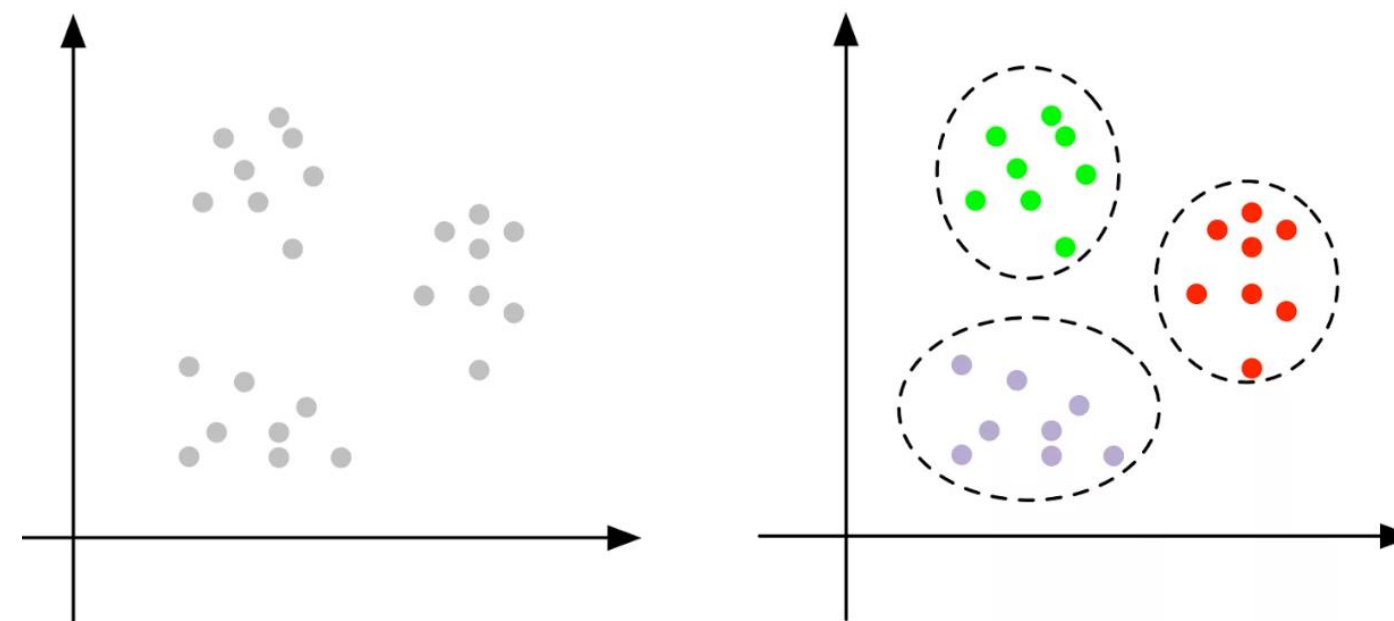
- *вектор* – слово
- *признаки* – токены из словаря

# Шаг 4. Методы интеллектуального анализа данных

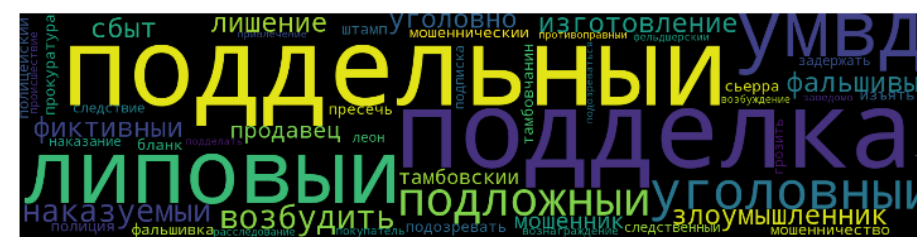
## Классификация



## Тематическое моделирование



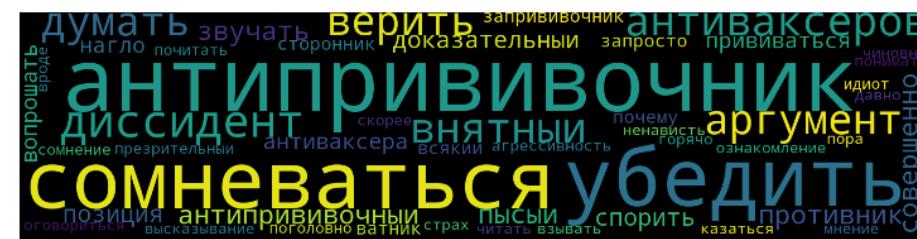
Topic 33



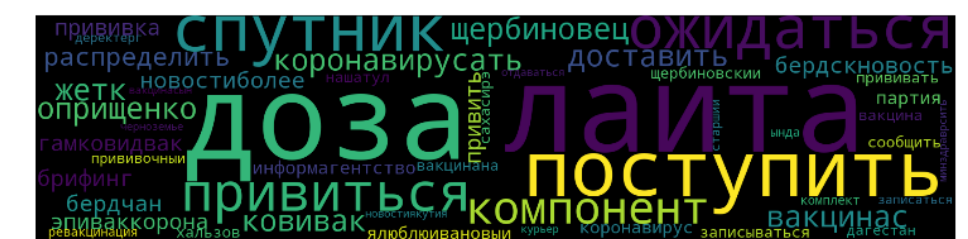
Topic 8



Topic 13



Topic 7



# Инструменты

## Python

4. Предварительная обработка данных

Первое, что необходимо сделать - проверить размерность признаков пространства. Данные не являются признаками - признаками, которые не имеют смысловой информации для дальнейшего анализа.

```
In [101]: df.head()
```

```
Out[101]:
```

name_type	sex	is_adult	abundance	is_adult	is_adult	is_adult	is_adult	is_adult	is_adult
1001	Мужчина	1.0	1000	1000	1000	1000	1000	1000	1000
1002	Женщина	0.0	1000	1000	1000	1000	1000	1000	1000
1003	Мужчина	1.0	1000	1000	1000	1000	1000	1000	1000
1004	Женщина	0.0	1000	1000	1000	1000	1000	1000	1000
1005	Мужчина	1.0	1000	1000	1000	1000	1000	1000	1000

## Orange

## KNIME

## PolyAnalyst

# Текстовая аналитика

Школа прикладного анализа данных Data-Diving



Университетский  
консорциум  
исследователей  
больших данных

**Басина Полина**

аналитик Центра прикладного  
анализа больших данных