«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ

«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»

**ИТОГОВАЯ АТТЕСТАЦИОННАЯ РАБОТА**

# Прогнозирование и анализ оттока пользователей по сообщениям в соцсети

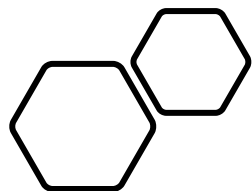Выполнил:

Крылов  Сергей
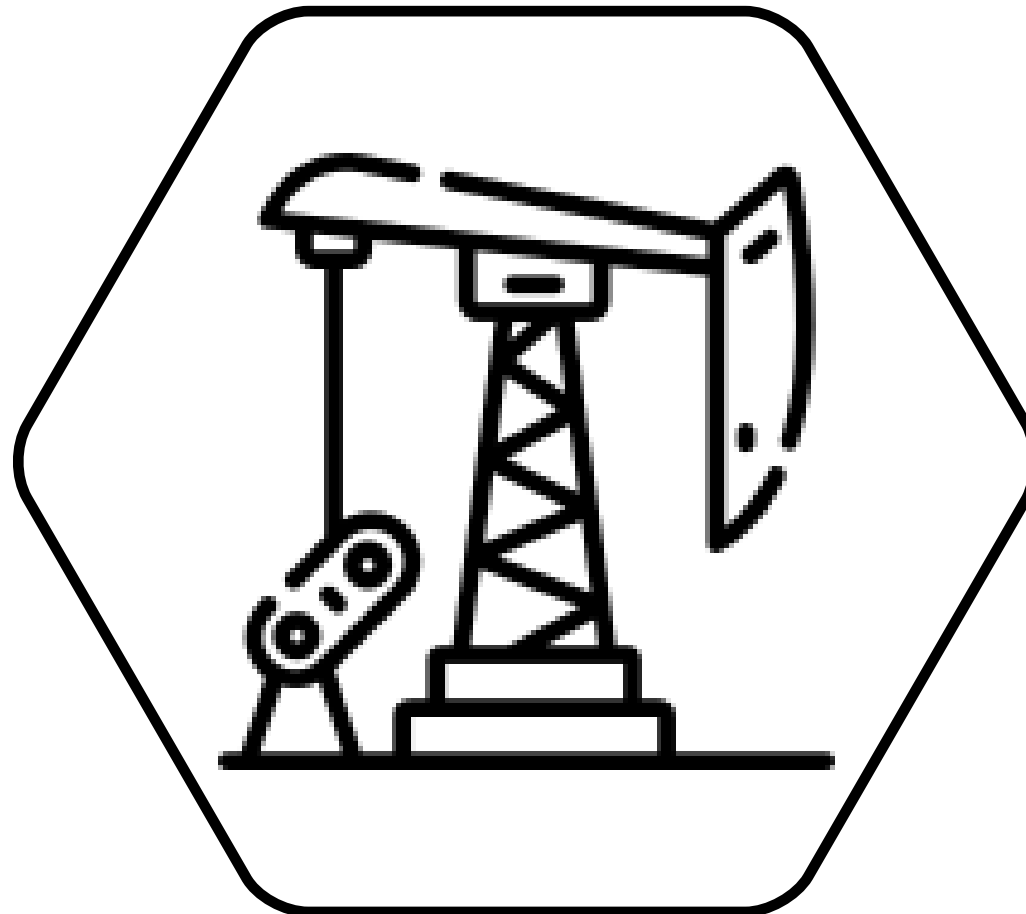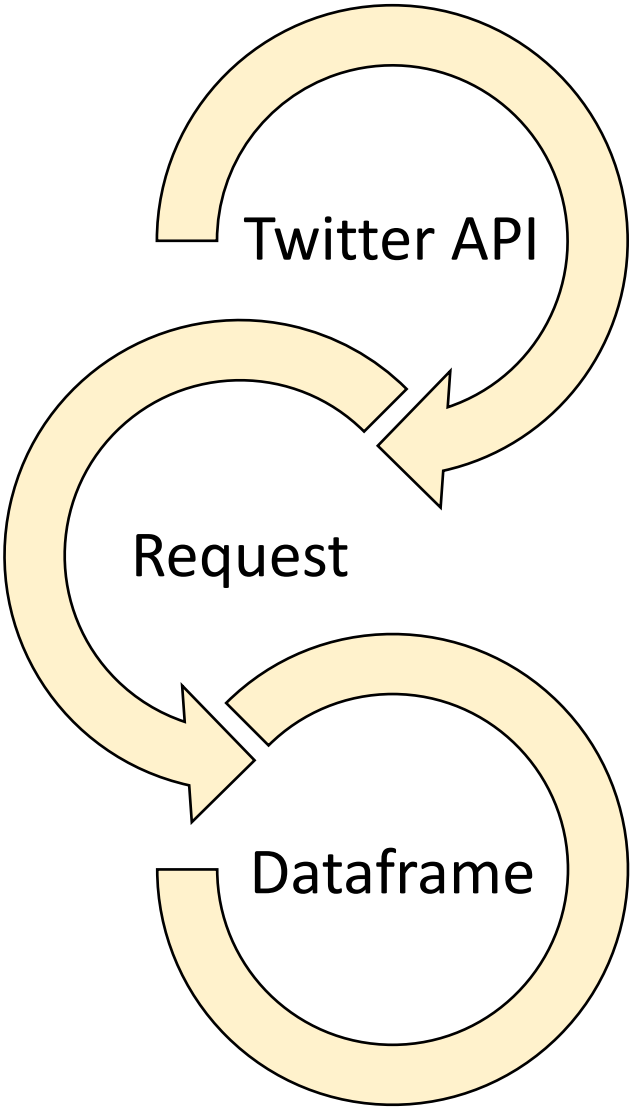
Руководитель:

Селезнев Артем

# Актуальность проблемы



Клиент

Реакция в Twitter

Сбор данных

API Twitter

Анализ данных

Предсказание оттока

Оператор связи

# Алгоритм работы



**Получение данных**

**Exploratory Data Analysis**

**Создание фичей**

**Кластеризация**

**Rule-based approach**

**Классификация**

# Получение данных

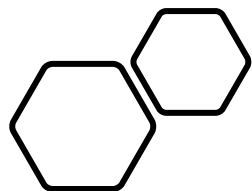| | Получение данных | EDA | Создание фичей | Кластеризация | Rule-based approach | Классификация |
|---|---|---|---|---|---|---|

Twitter API

Request

Dataframe

tweepy.OAuthHandler(tokens)
set_access_token
tweepy.API(auth)

tweepy.Cursor(api.search, q=searchString, lang='en')

| | screen_name | date_time | location | text |
|---|---|---|---|---|
| 0 | TheSkubis | 2021-03-30 23:59:19 | Pennsylvania, USA | @VerizonSupport I have issues with closed capt... |
| 1 | VerizonSupport | 2021-03-30 23:56:55 | | @_carolinek This could be due to regional rest... |
| 2 | VerizonSupport | 2021-03-30 23:55:30 | | @tvmurray We'll be happy to help with anything... |

# Exploratory Data Analysis

# *E*xploratory*D*ata*A*nalysis
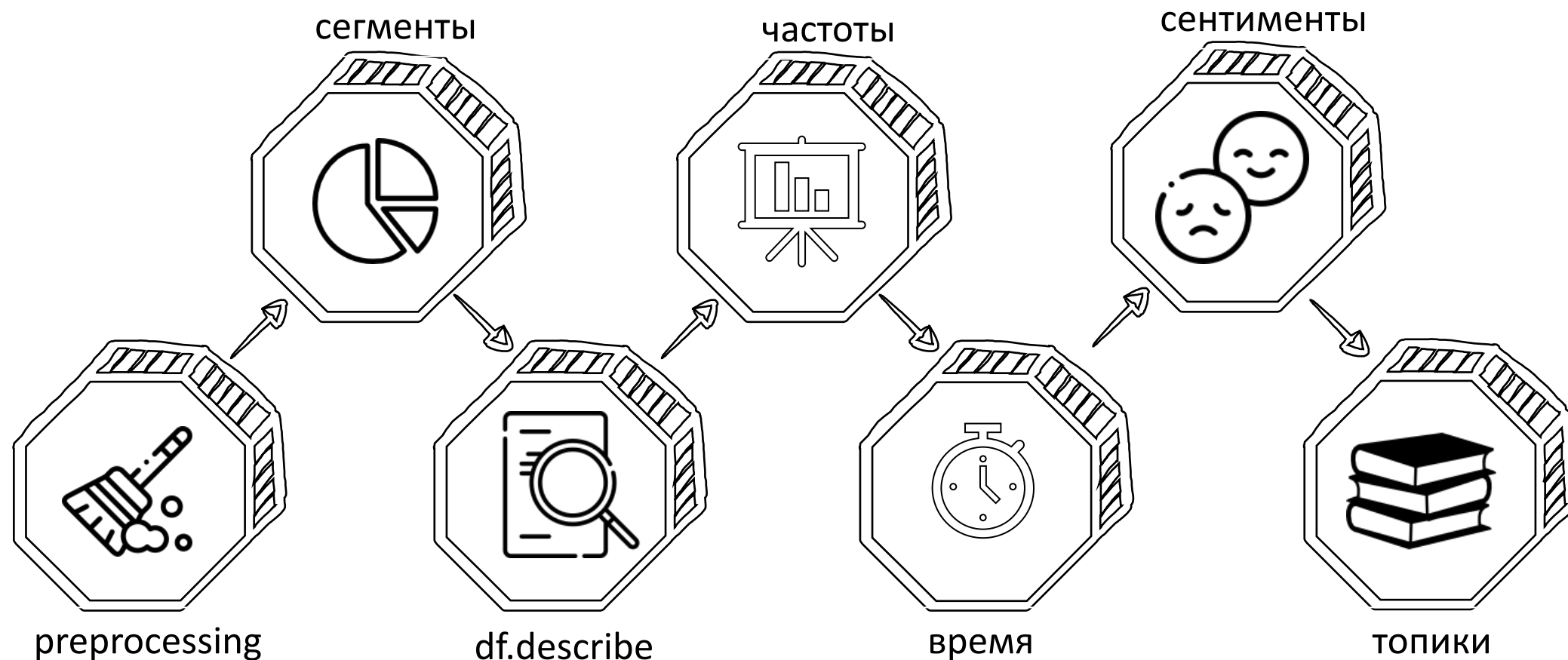


сегменты · частоты · сентименты

preprocessing · df.describe · время · топики

## **Preprocessing**

# Сегменты

# df.describe

## Анализ частот слов по сентиментам

## Время

- НУЖНО БОЛЬШЕ ДАННЫХ

<div style="background-color:red; text-align:center; font-size:3em; padding:1em;">ACTUAL</div>

# Сентименты

## **Топики**

1  help get got keep team verizon full per pas attempt

2  help cut time great get funny number given experiencing thanks

3  wireless verizon need best direct way get feel please work

4  using customer account hear app true tvision try need att

5  sent thanks told elc hey would keep needed tried order

6  help send month want take look code zip detail please

7  need give service classroom get hello please type work today

8  account month vacation billed suspension verizon att reaching call since

9  customer internet please thanks service working local victim tornado provider

10  help please happy send jersey follow detail today assist look

11  check device get new hour call detail data pay wireless

12  want make sure experience service need help fuck getting sorry

13  phone know verizon around hanging minute anyone spare internet let

14  verizon thank great live fios lost service hear att home

15  would love please look team meet back help get follow

16  tweet prank every tap year donating time verizon customer million

17  please verizon hello could one get send customer issue use

18  guy day supposed verizon back say long sprint past ordered

19  help always family happy service love card ever fios still

20  auto get cut said like thank maybe hold paying hell

# Создание фичей

| Получение данных | EDA | Создание фичей | Кластеризация | Rule-based approach | Классификация |
|---|---|---|---|---|---|



**TF-IDF**

**Отбор топ-80%**

**location_proba**

**Вероятность города**

**topic**

**Номер топика из topic_modeling**

**время**

Применение sin/cos трансформации

**RFM-analysis**

**Проекция RFM на Recency, frequency, polarity**

**fact**

Textblob. subjectivity

# Кластеризация

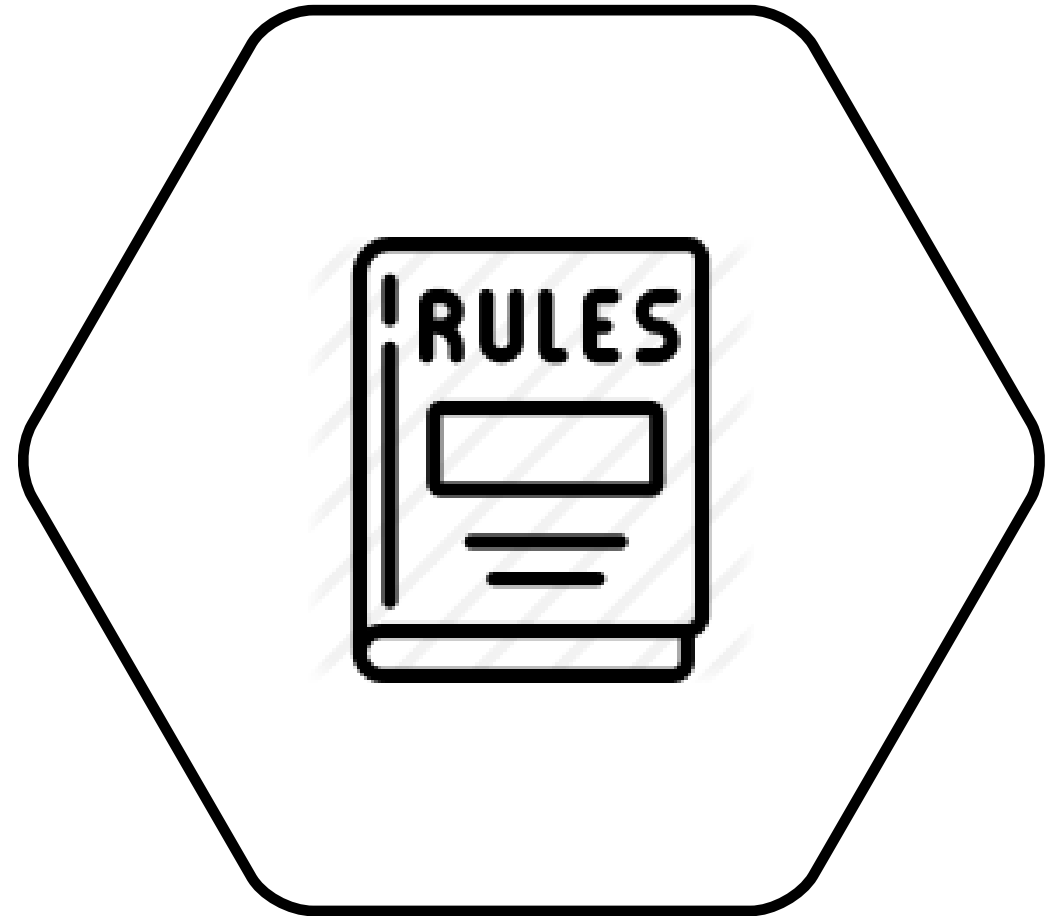|  | silhouette | roc_auc | accuracy | recall |
|---|---|---|---|---|
| **euclidean_churn** | 0.707732 | 0.346743 | 0.681733 | 0.000000 |
| **chebyshev_churn** | 0.644369 | 0.654215 | 0.427495 | 0.888889 |
| **sqeuclidean_churn** | 0.670944 | 0.311303 | 0.612053 | 0.000000 |

ACTUAL

Class separation using first two principal components

Ruled-based approach

| Rule | | churny words | churny reasons |
|---|---|---|---|
| **'from'** | -0.8 | switch , transfer | service, better, experts, disney |
| **'to'** | +0.8 | free, come | network, wifi, price, tower, coverage |
| **'with'** | +0.5 | come, change | family, horrible, awfull, slow, free of charge |
| **'like'** | +0.5 | leave, stay | worst, worse, bad, price, money |
| **'Disney' & 'Verizon'** | +0.5 | welcome, goodbye | can\'t stand, promise |
| **sentiment** | Count(carrier) x polarity | | |
| **score < 0** | Отток_от.append(оператор) | | |
| **Score > 0** | Отток_к.append(оператор) | | |

Churn Out VS. Churn In for per carrier

Churn Out VS. Churn In for per carrier

## Хороший пример

```
@JewdyGold @VerizonSupport Unacceptable, which is why I canceled my service with Verizon.
scores:  {'verizon': 0.0}
Subjects: ['Verizon']
Reason: ['@JewdyGold @VerizonSupport Unacceptable, which is why I canceled my service with Verizon.']
Conclusion: Churn from  ['verizon']  to  []
```

## Пример не очень

```
@NotLacking_ @OMGItsBirdman @verizon please help!!!! I can't see the pictures!!!
scores:  {'verizon': 0.0}
Subjects: ['Verizon']
Reason: []
Conclusion: Churn from  ['verizon']  to  []
```
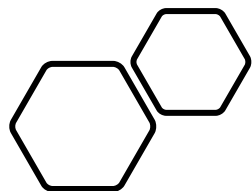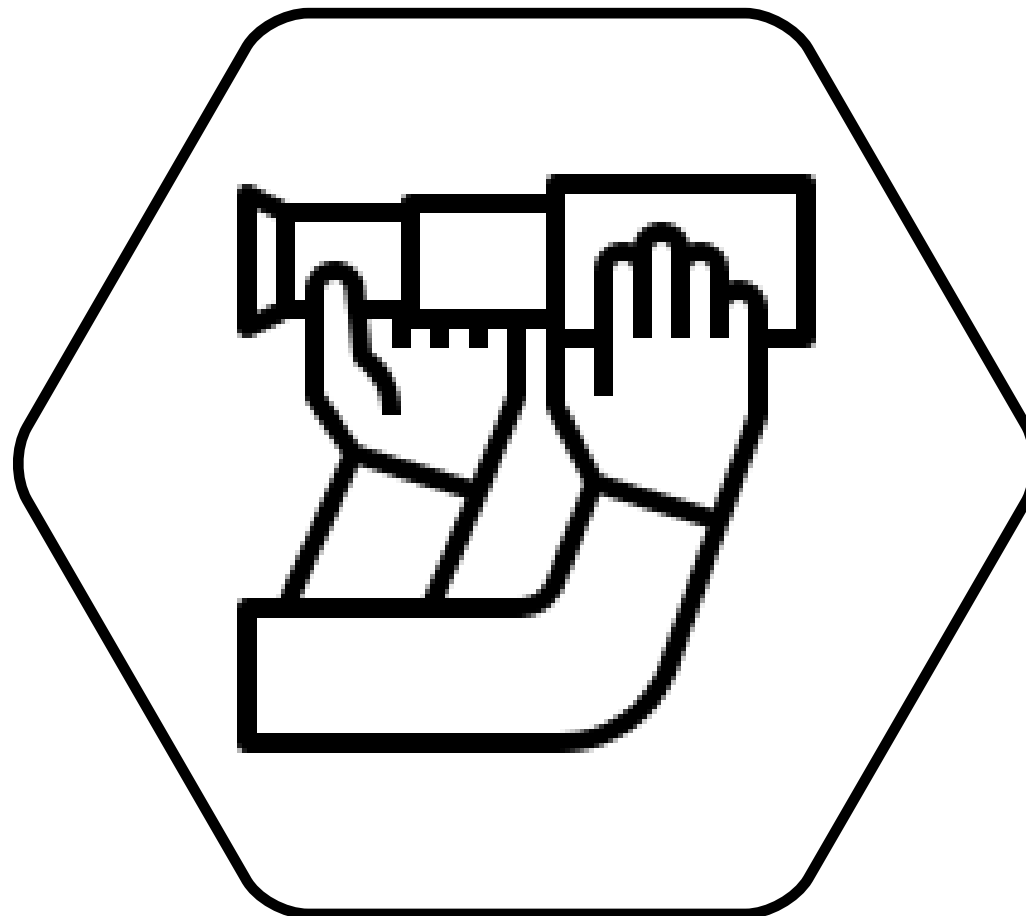
# Облако слов причин оттока

# Классификация

| model | roc_auc |
|---|---|
| **Naïve Bayes** | 0.546875 |
| **Logistic regression** | 0.750625 |
| **LogRegression + grid_search** | 0.797500 |
| **LogRegression + MinMaxScaler** | 0.792917 |
| **LogRegression + grid_search + MinMaxScaler** | 0.817500 |
| **Decision Tree** | 0.815208 |
| **LightGBM** | 0.857708 |

# LightGBM



|  | precision | recall | score | support |
|---|---|---|---|---|
| **0** | 0.99 | 0.95 | 0.97 | 75 |
| **1** | 0.89 | 0.97 | 0.93 | 32 |
| **accuracy** |  |  | 0.95 | 107 |
| **macro avg** | 0.94 | 0.96 | 0.95 | 107 |
| **weighted avg** | 0.96 | 0.95 | 0.95 | 107 |

|  | train | test |
|---|---|---|
| **Cross_validate** | 1 | 0.99 |

спасибо за внимание!