

# ОПРЕДЕЛЕНИЕ ЖАНРА МУЗЫКИ НА ОСНОВЕ ЕЁ ТЕКСТА

Кряжев Макар

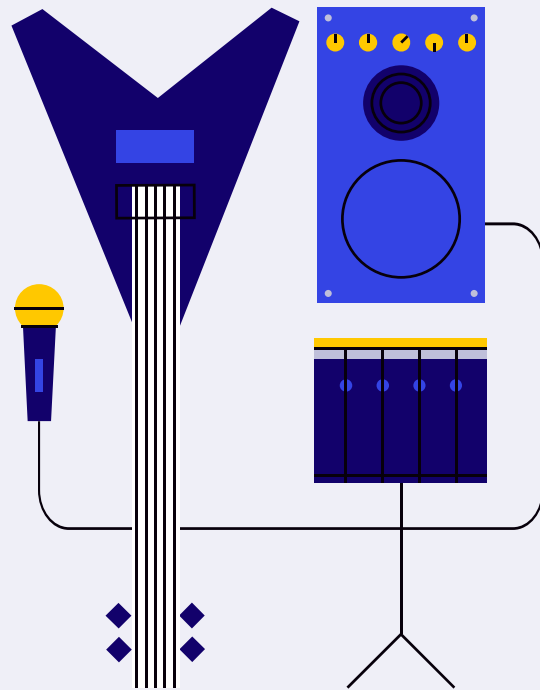
Дю Вадим

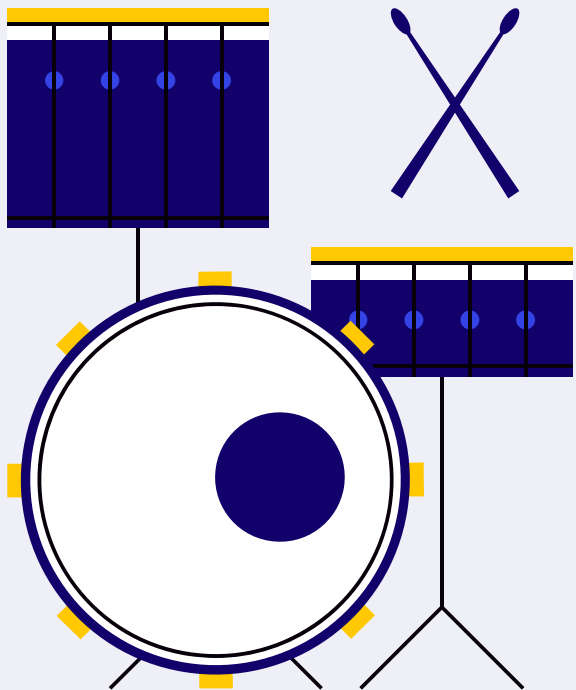
Талпа Григорий

Сухова Екатерина

Ерохин Дмитрий

Куратор: Тимофей Соборнов





# СБОР И РАЗМЕТКА ДАННЫХ

# Первоначальный план

1. Скрапинг текстов и жанров песен с [genius.com](https://genius.com)
2. Profit

Но...

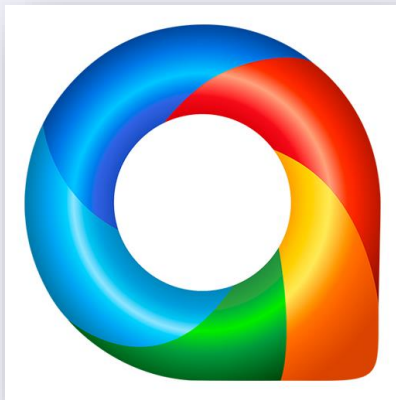
“Жадный” Genius API, отсутствие альтернативного источника информации

# Принцип работы

Основная  
информация



Текст песни



Жанр



# Недостатки

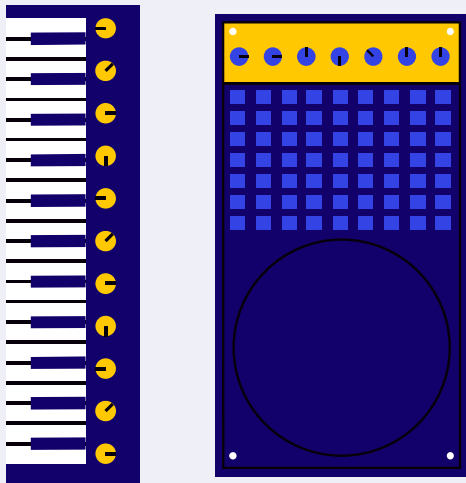
- Неэффективность (перебор id, 3 источника)
- Медлительность (150 песен/25 минут)

# Готовый датасет

- > ≈ 6 миллионов строк
- > 8 признаков
- > 6 жанров (rap, rock, country, rb, pop, misc)

# Готовый датасет

<b>title</b>	Название. Большинство записей – песни, но также присутствуют книги, поэмы и т.д.
<b>tag</b>	Жанр. Немузыкальные записи имеют жанр misc.
<b>artist</b>	Автор(ы) произведения
<b>year</b>	Год выпуска
<b>views</b>	Количество просмотров страницы на genius.com
<b>features</b>	Артисты, участвовавшие в создании
<b>lyrics</b>	Текст произведения
<b>id</b>	Идентификатор произведения на genius.com



# ОБРАБОТКА



# Обработка

- **Предварительная обработка**

  - Балансировка классов

  - Удаление пунктуации

- **Необходимая подготовка**

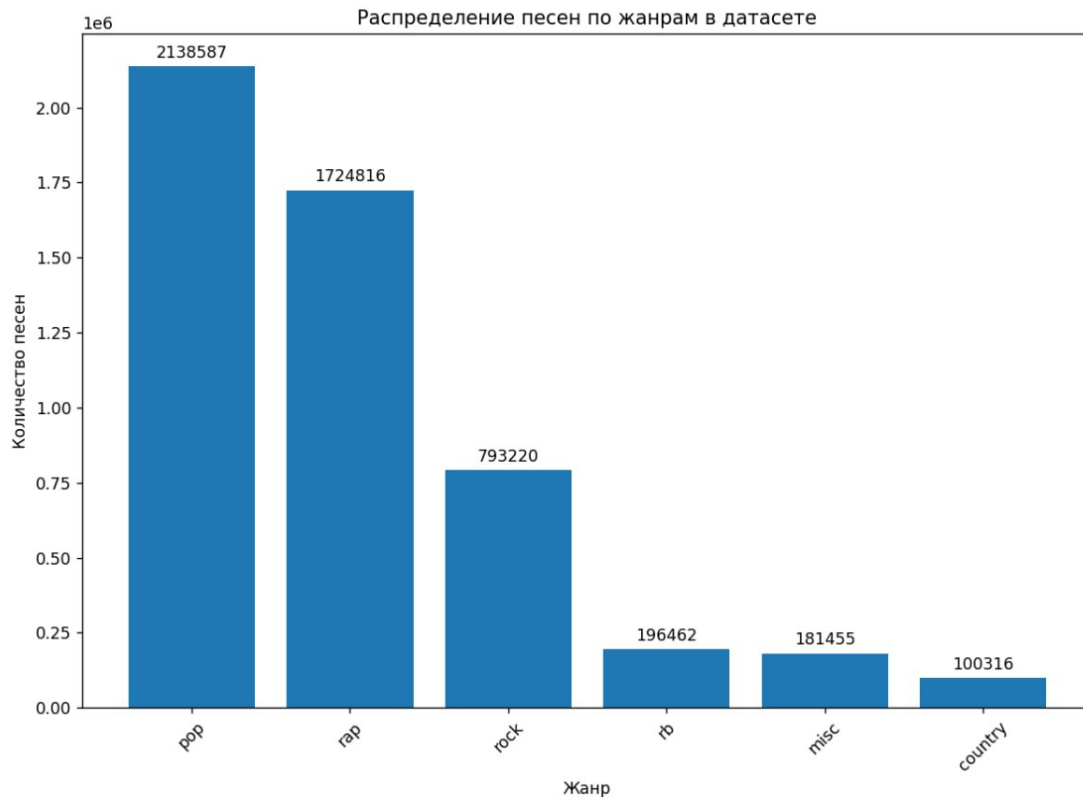
  - Токенизация

  - Лемматизация

  - Векторизация (B-O-W, Tf-Idf, word2vec,...)

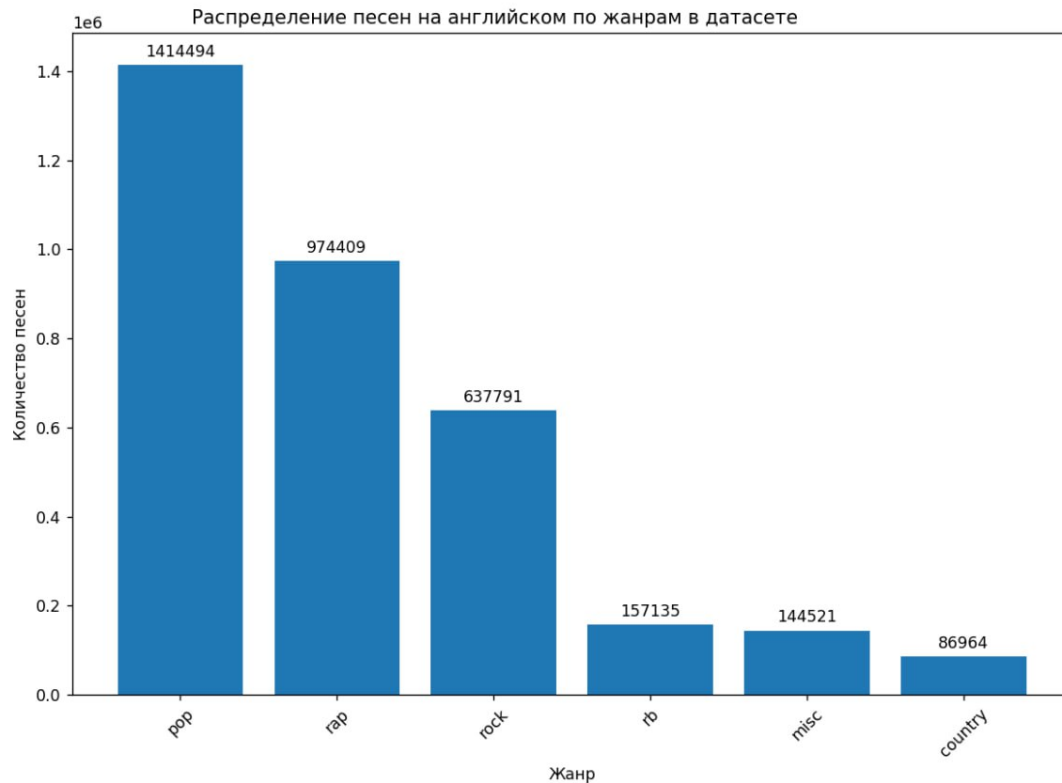
# Обработка

Убираем все  
песни кроме  
английских



# Балансировка классов

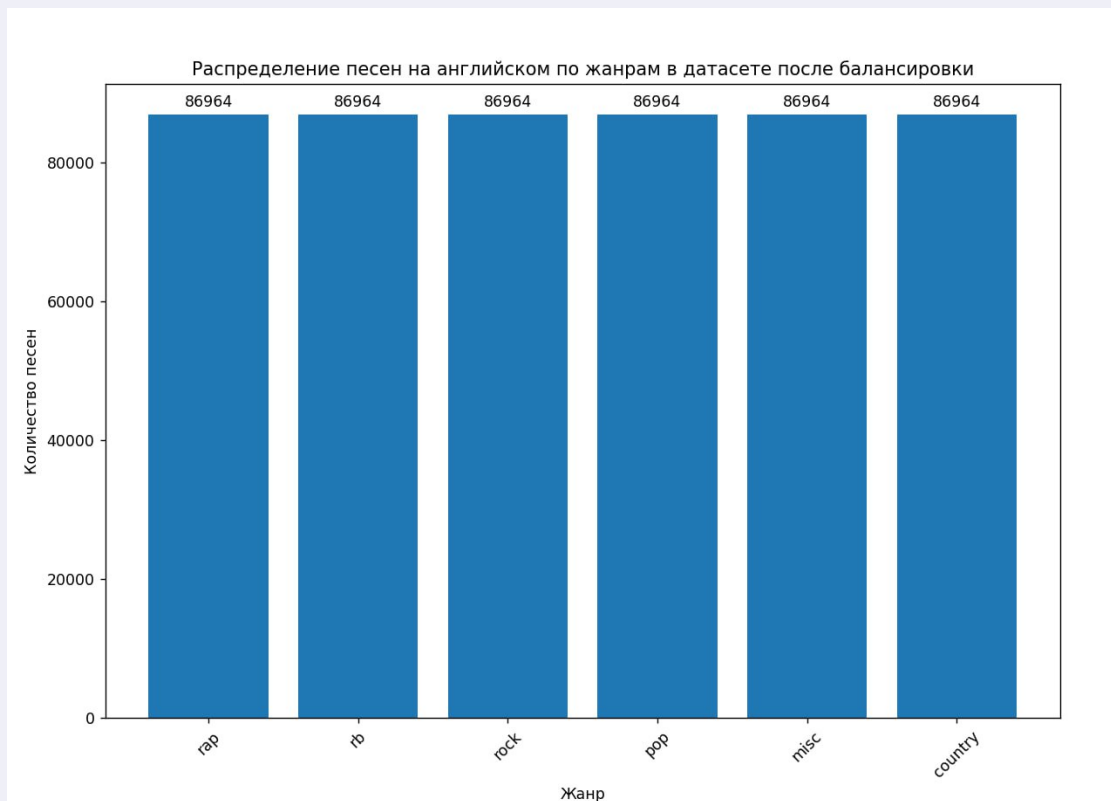
Уменьшаем  
количество  
данных

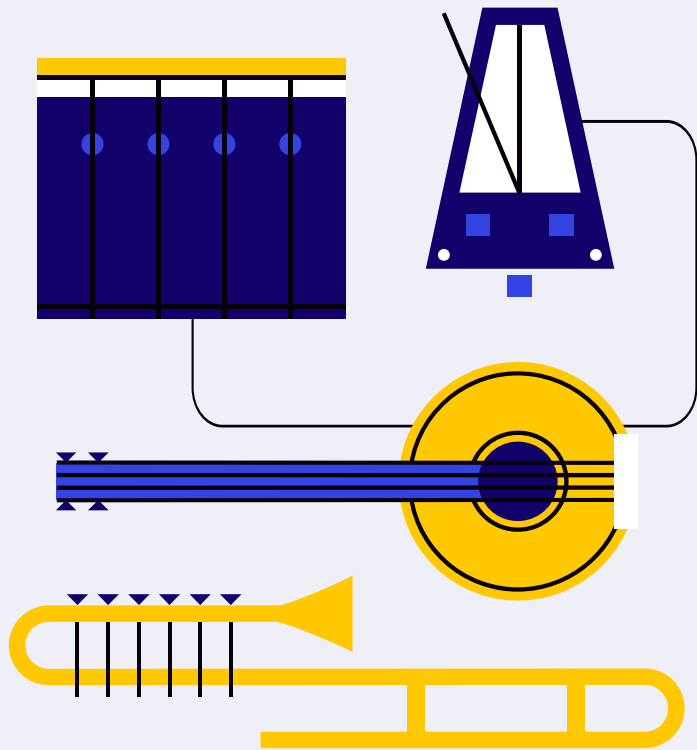


# Балансировка классов

Балансируем  
классы  
относительно  
минимального

≈ 87 000 строк в  
каждом жанре





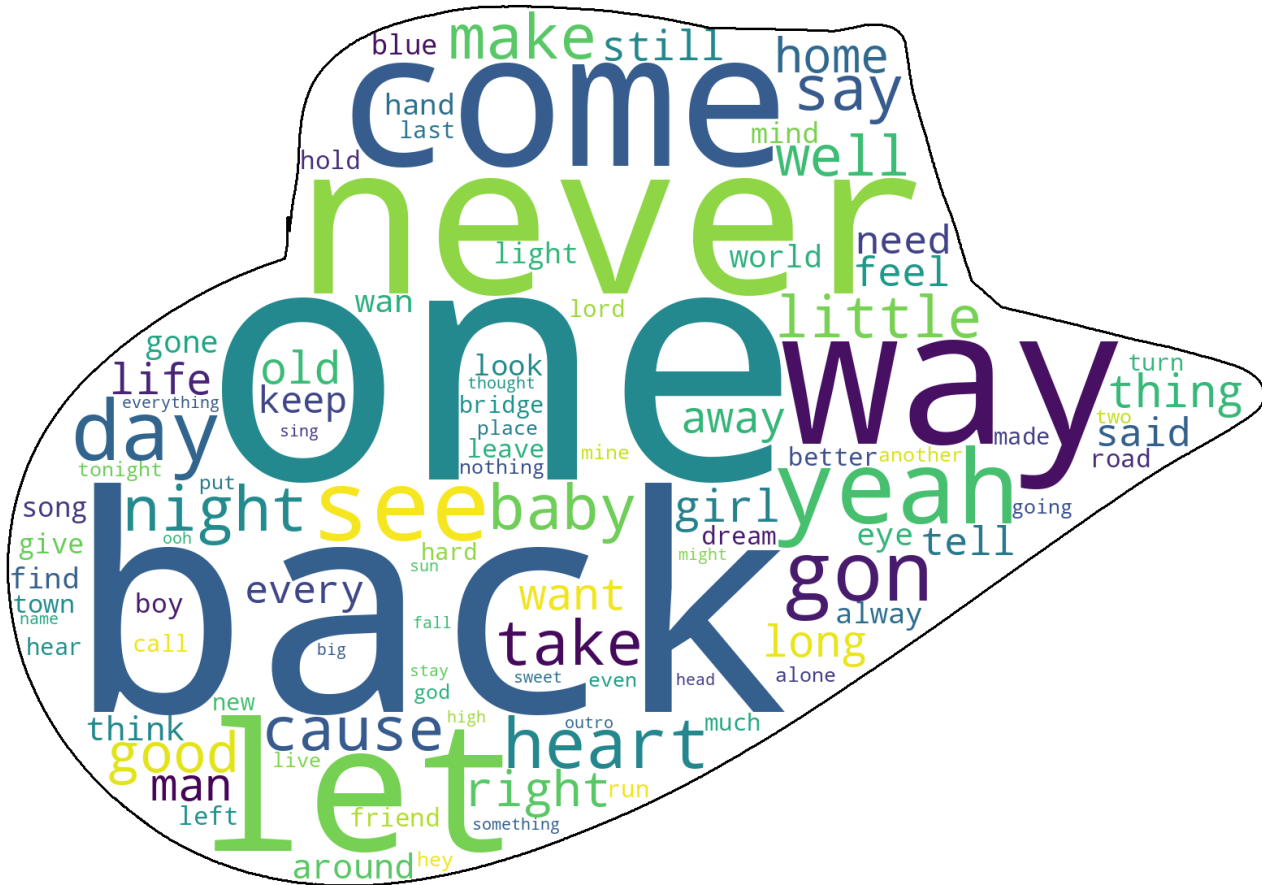
**EDA**

# Облака слов для каждого жанра

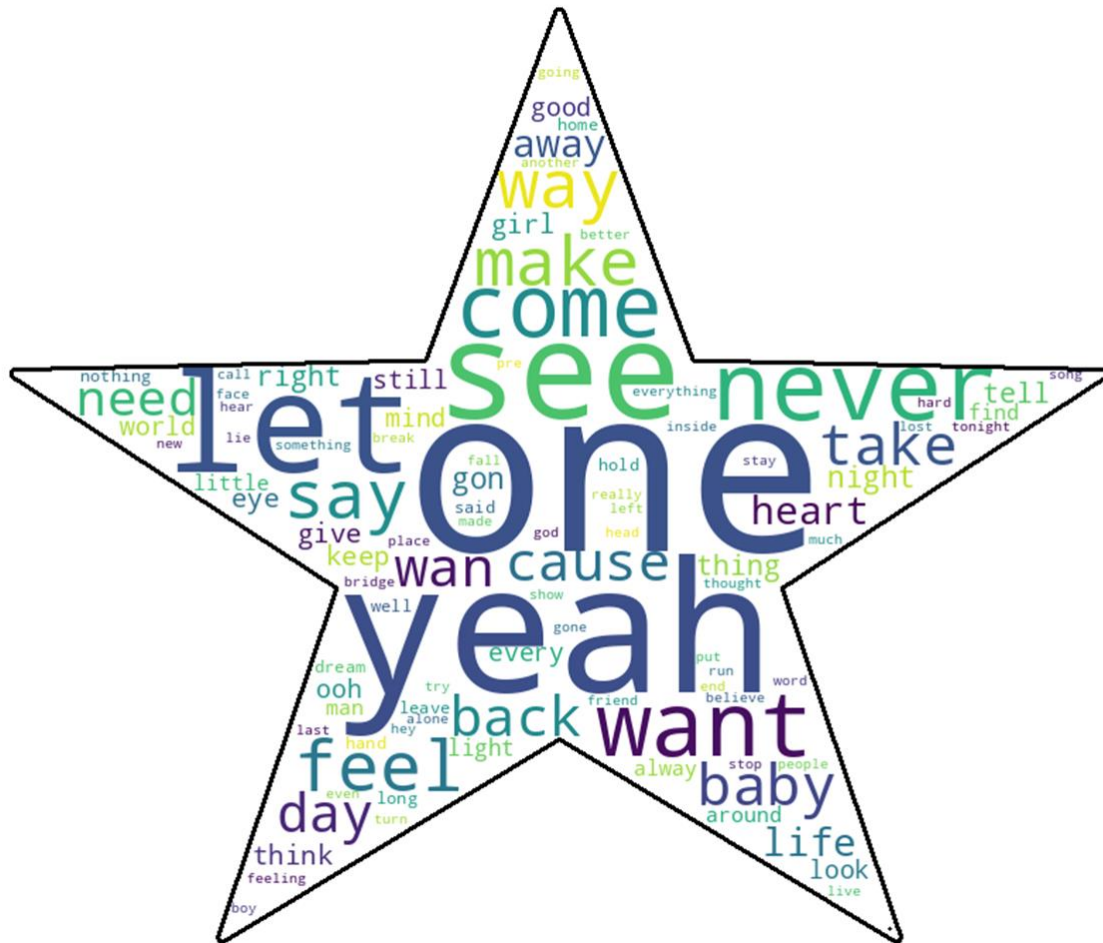
При построении не были использованы следующие слова:

- > verse
- > love
- > know
- > chorus
- > time
- > got

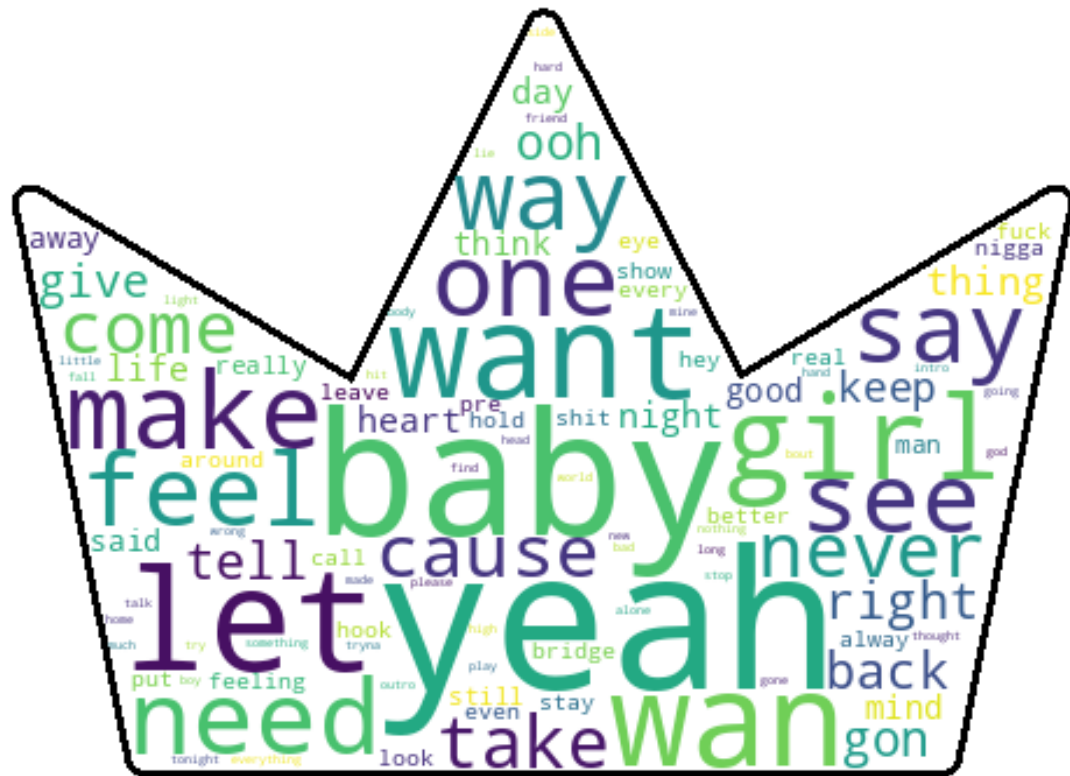
country



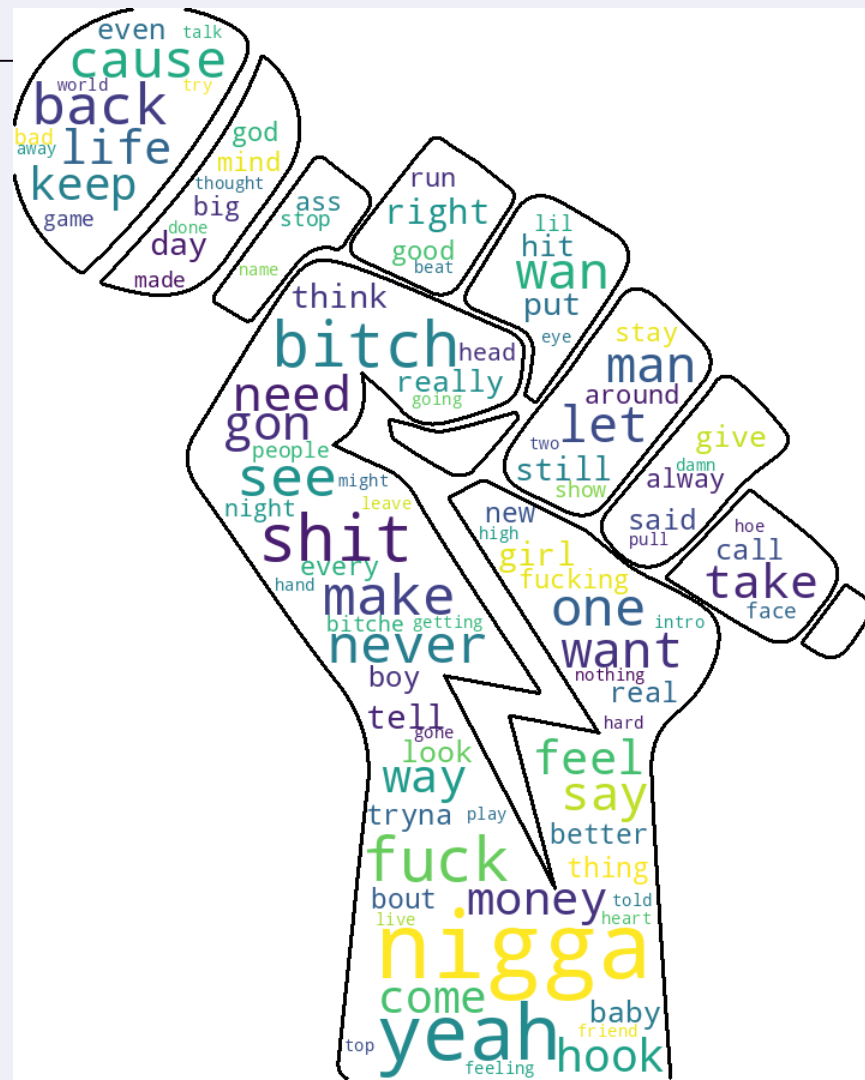
pop





**rb**

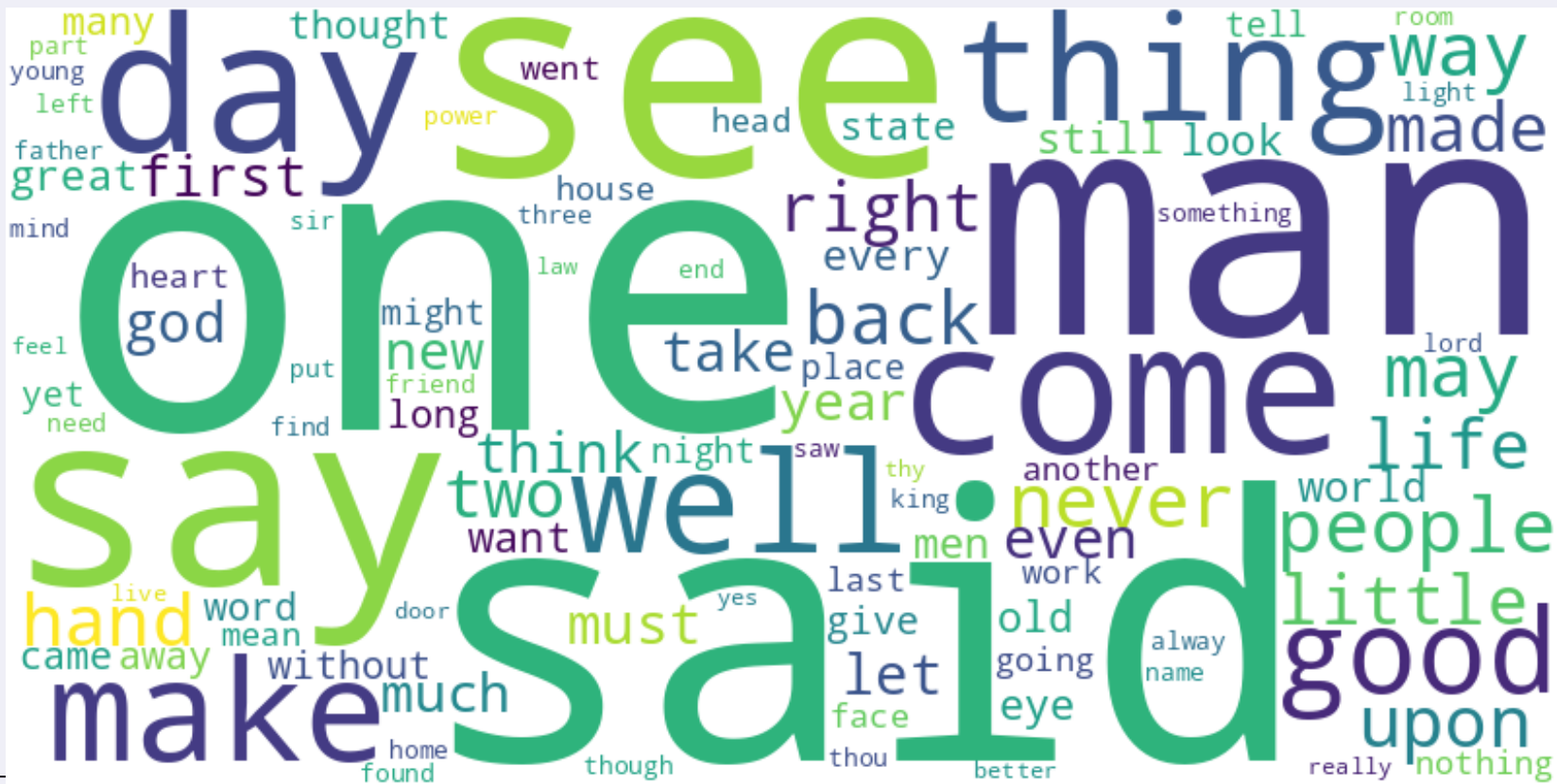
rap



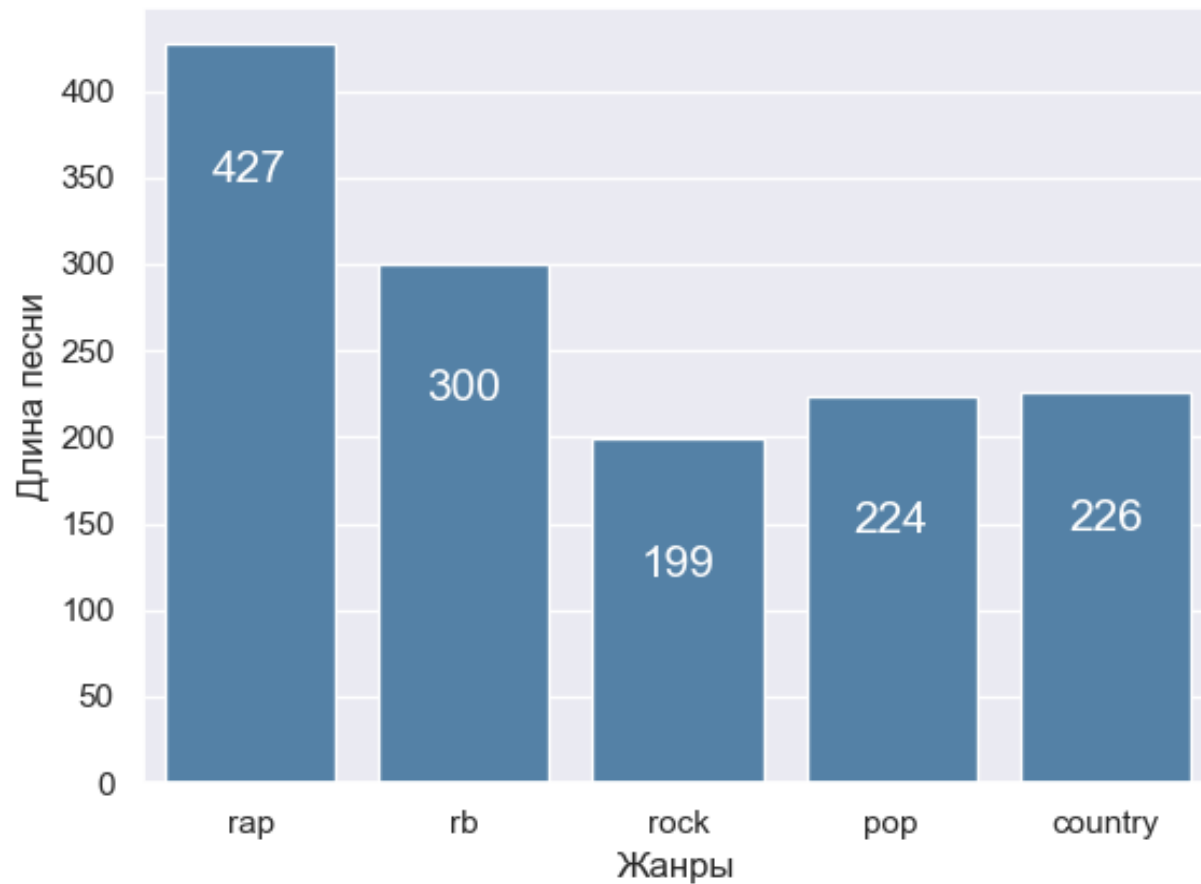
# rock



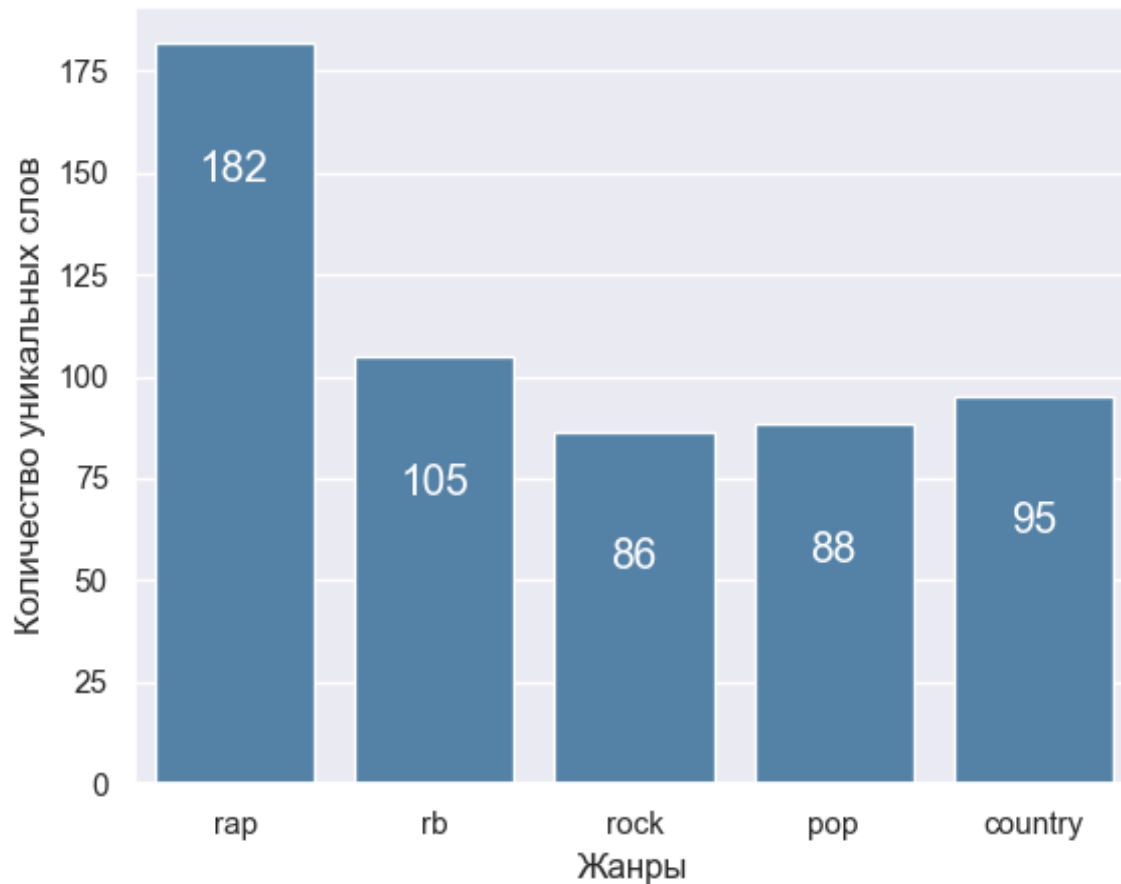
# misc



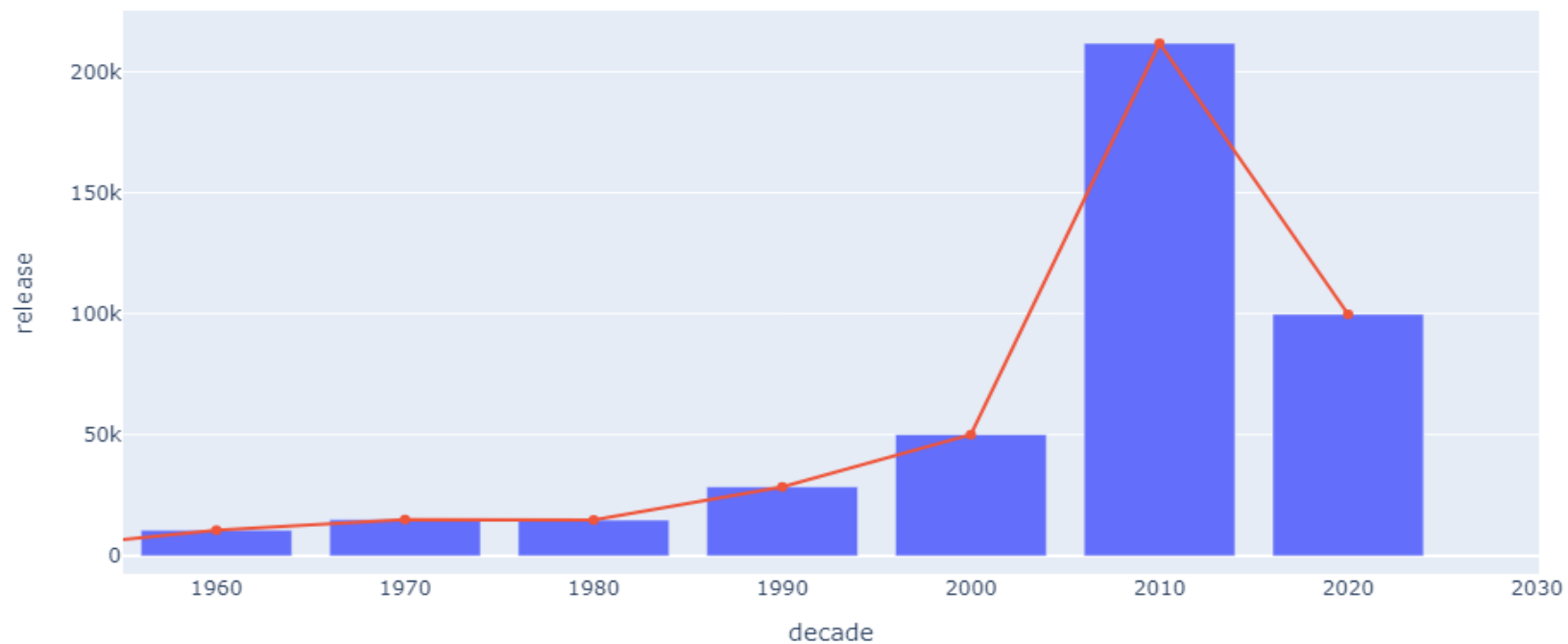
Среднее количество слов в песне по жанрам



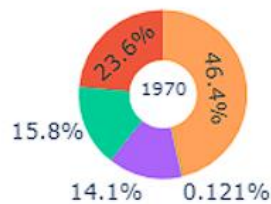
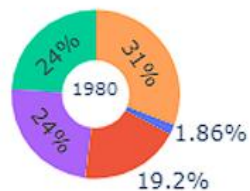
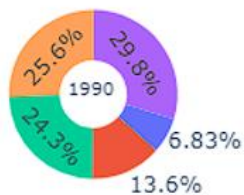
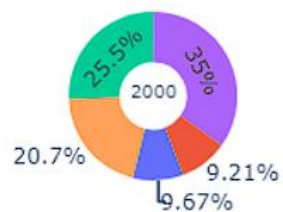
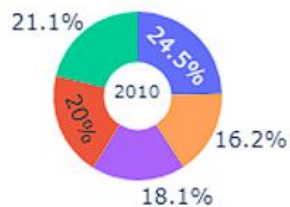
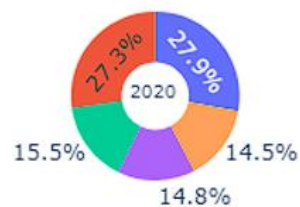
Среднее Количество уникальных слов в песне по жанрам



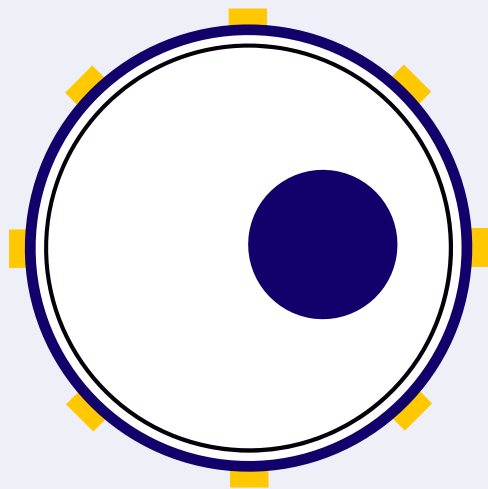
## Количество песен по декадам



## Распределение жанров по декадам







# ОБУЧЕНИЕ МОДЕЛЕЙ

# Модели

- > Logistic Regression
- > Multinomial Naive Bayes
- > Random Forest
- > Gradient Boosting
- > BERT

# Logistic Regression (TF-IDF)

	precision	recall	f1-score	support
country	0.64	0.70	0.67	21525
misc	0.76	0.73	0.74	21783
pop	0.37	0.32	0.34	21788
rap	0.78	0.77	0.78	21879
rb	0.58	0.59	0.59	21761
rock	0.51	0.54	0.52	21710
accuracy			0.61	130446
macro avg	0.60	0.61	0.61	130446
weighted avg	0.60	0.61	0.61	130446

# Multinomial Naive Bayes (ngram)

	precision	recall	f1-score	support
country	0.60	0.66	0.63	17250
misc	0.76	0.67	0.72	17387
pop	0.34	0.21	0.26	17448
rap	0.70	0.74	0.72	17521
rb	0.48	0.64	0.55	17424
rock	0.47	0.47	0.47	17327
accuracy			0.56	104357
macro avg	0.56	0.56	0.56	104357
weighted avg	0.56	0.56	0.56	104357

# **BERT** (distilbert-base-uncased)

DistilBertForSequenceClassification

- > f1-score: 0.189
- > accuracy: 0.189



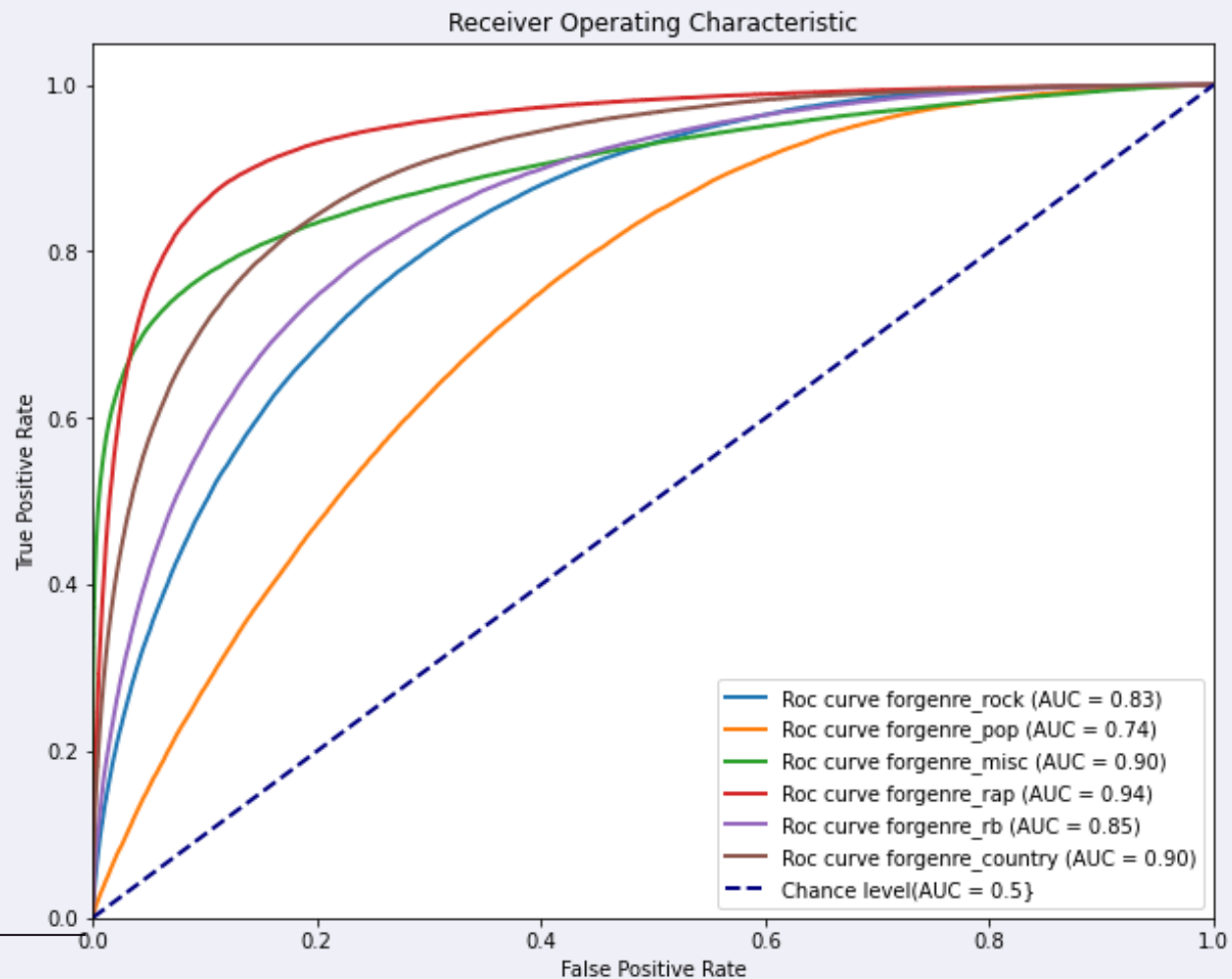
## Особенность задачи определения жанра песни

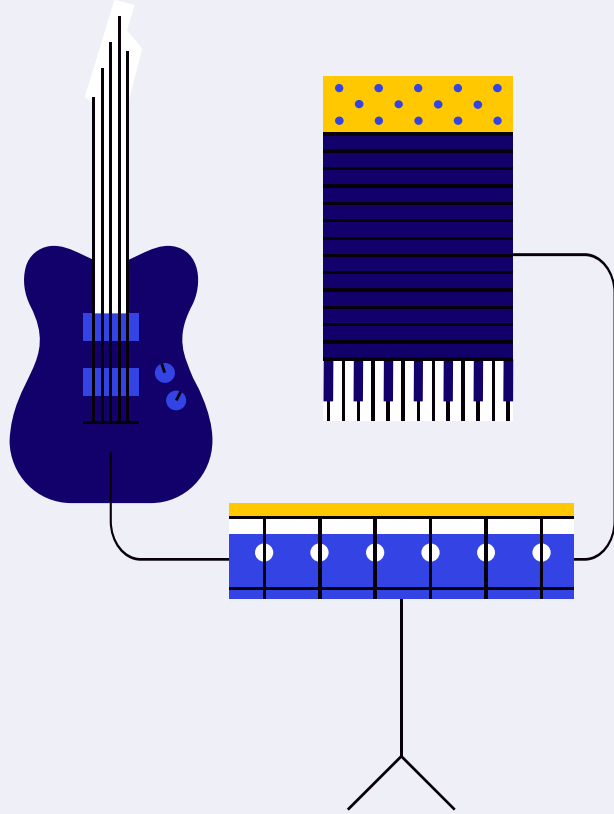
Если модель пытается определить класс песни, которую можно отнести к нескольким жанрам, она работает неэффективно

**Итог:** нельзя относить песню только к 1 классу



# Multiclass classification





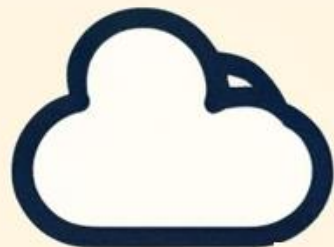
**MVP**



# Streamlit



Yandex cloud



Sending model



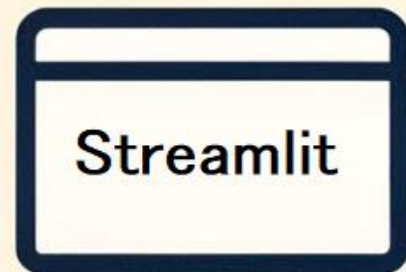
BACK-End

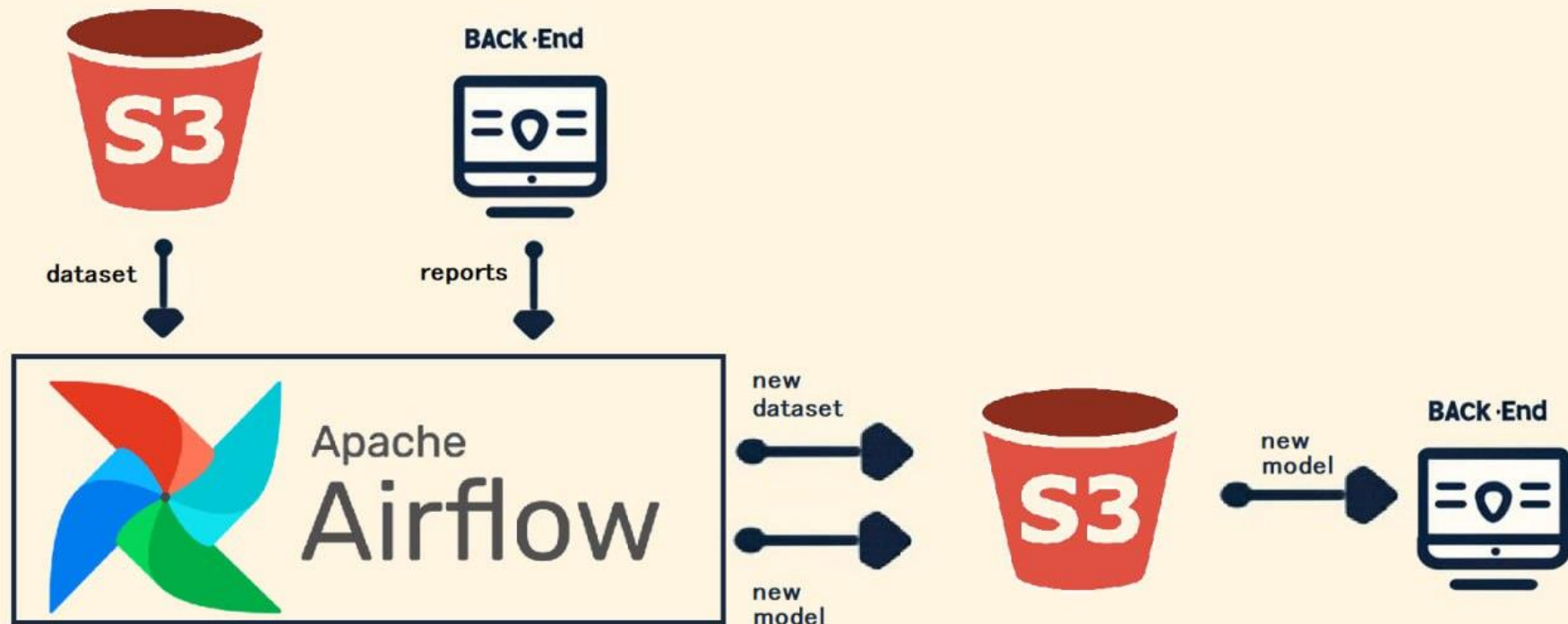


Sending list  
of answers



Request





**get\_csv\_op**

■ success

PythonOperator

**merge\_csv\_op**

■ success

PythonOperator

**train\_model\_op**

■ success

PythonOperator

**upload\_all**

■ success

PythonOperator

# Перспективы

- › Подбор похожих песен по тексту
- › Анализ всей песни, а не только текста
- › Подбор похожих песен по звучанию
- › И многое другое...

