# Проект: Статистический анализ данных

Вы аналитик популярного сервиса аренды самокатов GoFast. Вам передали данные о некоторых пользователях из нескольких городов, а также об их поездках. Проанализируйте данные и проверьте некоторые гипотезы, которые могут помочь бизнесу вырасти.

Чтобы совершать поездки по городу, пользователи сервиса GoFast пользуются мобильным приложением. Сервисом можно пользоваться:

- . без подписки
  - .1. абонентская плата отсутствует;
  - .2. стоимость одной минуты поездки 8 рублей;
  - .3. стоимость старта (начала поездки) -50 рублей;
- . с подпиской Ultra
  - .1. абонентская плата 199 рублей в месяц;
  - .2. стоимость одной минуты поездки 6 рублей;
  - .3. стоимость старта бесплатно.

## Описание данных

В основных данных есть информация о пользователях, их поездках и подписках.

Пользователи — users\_go.csv

user_id	уникальный идентификатор пользователя
name	имя пользователя
age	возраст

city	город
subscription_type	тип подписки (free, ultra)

# Поездки — <u>rides\_go.csv</u>

user_id	уникальный идентификатор пользователя
distance	расстояние, которое пользователь проехал в текущей сессии (в метрах)
duration	продолжительность сессии (в минутах) — время с того момента, как пользователь нажал кнопку «Начать поездку» до момента, как он нажал кнопку «Завершить поездку»
date	дата совершения поездки

## Подписки — <u>subscriptions\_go.csv</u>

subscription_type	тип подписки
minute_price	стоимость одной минуты поездки по данной подписке
start_ride_price	стоимость начала поездки
subscription_fee	стоимость ежемесячного платежа

## Шаг 1. Загрузка данных

- 1.1 Считайте CSV-файлы с данными с помощью библиотеки pandas и сохраните их в датафреймы. Пути к файлам:
- . /datasets/users\_go.csv
- . /datasets/rides\_go.csv
- . /datasets/subscriptions\_go.csv

1.2 Выведите первые строки каждого набора данных. Изучите общую информацию о каждом датафрейме.

#### Шаг 2. Предобработка данных

- 2.1 Приведите столбец date к типу даты pandas.
- 2.2 Создайте новый столбец с номером месяца на основе столбца date.
- 2.3 Проверьте наличие пропущенных значений и дубликатов в датафреймах. Обработайте их, если такие значения присутствуют.

#### Шаг 3. Исследовательский анализ данных

Опишите и визуализируйте общую информацию о пользователях и поездках:

- 3.1 частота встречаемости городов;
- 3.2 соотношение пользователей с подпиской и без подписки;
- 3.3 возраст пользователей;
- 3.4 расстояние, которое пользователь преодолел за одну поездку;
- 3.5 продолжительность поездок.

#### Шаг 4. Объединение данных

- 4.1 Объедините данные о пользователях, поездках и подписках в один датафрейм. Для этого воспользуйтесь методом merge().
- 4.2 Создайте ещё два датафрейма из датафрейма, созданного на этапе 4.1:
- . с данными о пользователях без подписки;
- . с данными о пользователях с подпиской.
- 4.3 Визуализируйте информацию о расстоянии и времени поездок для пользователей обеих категорий.

#### Шаг 5. Подсчёт выручки

- 5.1 Создайте датафрейм с агрегированными данными о поездках на основе датафрейма с объединёнными данными из шага 4: найдите суммарное расстояние, количество поездок и суммарное время для каждого пользователя за каждый месяц.
- 5.2 В этот же датафрейм добавьте столбец с помесячной выручкой, которую принёс каждый пользователь. Для этого обратитесь к информации об условиях оплаты для подписчиков и тех, у кого нет подписки. Продолжительность каждой поездки в каждой строке исходного датафрейма для подсчёта стоимости округляется до следующего целого числа: например, значения 25.3, 25.5 и 26.0 должны быть преобразованы к 26.

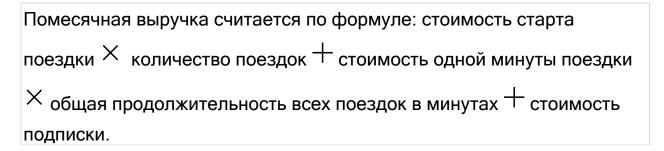
#### Подсказка

Продолжительность каждой поездки должна быть округлена с помощью метода «потолок» библиотеки numpy — пр. сеіі():

```
import numpy as np

# каждое значение из столбца duration округляется с помощью «потолка»:

rides_df['duration'] = np.ceil(rides_df['duration'])
```



#### Шаг 6. Проверка гипотез

Продакт-менеджеры сервиса хотят увеличить количество пользователей с подпиской. Для этого они будут проводить различные акции, но сначала нужно выяснить несколько важных моментов.

6.1 Важно понять, тратят ли пользователи с подпиской больше времени на поездки? Если да, то пользователи с подпиской могут быть «выгоднее» для компании. Проверьте гипотезу. Используйте

исходные данные о продолжительности каждой сессии — отдельно для подписчиков и тех, у кого нет подписки.

- 6.2 Расстояние одной поездки в 3130 метров оптимальное с точки зрения износа самоката. Можно ли сказать, что среднее расстояние, которое проезжают пользователи с подпиской за одну поездку, не превышает 3130 метров? Проверьте гипотезу и сделайте выводы.
- 6.3 Проверьте гипотезу о том, будет ли помесячная выручка от пользователей с подпиской по месяцам выше, чем выручка от пользователей без подписки. Сделайте вывод.
- 6.4 Представьте такую ситуацию: техническая команда сервиса обновила сервера, с которыми взаимодействует мобильное приложение. Она надеется, что из-за этого количество обращений в техподдержку значимо снизилось. Некоторый файл содержит для каждого пользователя данные о количестве обращений до обновления и после него. Какой тест вам понадобился бы для проверки этой гипотезы?

Шаг 7 (необязательное задание). Распределения

7.1 Отделу маркетинга GoFast поставили задачу: нужно провести акцию с раздачей промокодов на один бесплатный месяц подписки, в рамках которой как минимум 100 существующих клиентов должны продлить эту подписку. То есть по завершении периода действия подписки пользователь может либо отказаться от неё, либо продлить, совершив соответствующий платёж.

Эта акция уже проводилась ранее и по итогу выяснилось, что после бесплатного пробного периода подписку продлевают  $10\,\%$  пользователей. Выясните, какое минимальное количество промокодов нужно разослать, чтобы вероятность не выполнить план была примерно  $5\,\%$ . Подберите параметры распределения, описывающего эту ситуацию, постройте график распределения и сформулируйте ответ на вопрос о количестве промокодов.

## Подсказка

Нужно использовать биномиальное распределение, которое описывает указанную ситуацию. Затем подобрать подходящее значение параметра *N* для заданного параметра p=0.1 с помощью графиков и метода cdf(), сформулировать вывод и ответ.

7.2 Отдел маркетинга рассылает клиентам push-уведомления в мобильном приложении. Клиенты могут открыть его или не открывать. Известно, что уведомления открывают около 40 % получивших клиентов. Отдел планирует разослать 1 млн уведомлений. С помощью аппроксимации постройте примерный график распределения и оцените вероятность того, что уведомление откроют не более 399,5 тыс. пользователей.

## Подсказка

Эта ситуация тоже описывается биномиальным распределением. Но считать каждое отдельное значение достаточно долго. Вы можете воспользоваться нормальной аппроксимацией биномиального распределения и cdf() для быстрой оценки.