

Исследование данных музыкального сервиса - сравнение музыкальных предпочтений пользователей двух городов

Содержание

- 1 Обзор данных
- 2 Предобработка данных
 - 2.1 Стиль заголовков
 - 2.2 Пропуски значений
 - 2.3 Дубликаты
- 3 Проверка гипотез
 - 3.1 Сравнение поведения пользователей двух столиц
 - 3.2 Музыка в начале и в конце недели
 - 3.3 Жанровые предпочтения в Москве и Петербурге
- 4 Итоги исследования

Сравнение Москвы и Петербурга окружено мифами. Например:

- Москва — мегаполис, подчинённый жёсткому ритму рабочей недели;
- Петербург — культурная столица, со своими вкусами.

На данных Яндекс.Музыки сравним поведение пользователей двух столиц.

Цель исследования — проверка трёх гипотез:

- Активность пользователей зависит от дня недели. Причём в Москве и Петербурге это проявляется по-разному.
- В понедельник утром в Москве преобладают одни жанры, а в Петербурге — другие. Так же и вечером пятницы преобладают разные жанры — в зависимости от города.
- Москва и Петербург предпочитают разные жанры музыки. В Москве чаще слушают поп-музыку, в Петербурге — русский рэп.

Ход исследования

Данные о поведении пользователей хранятся в файле `yandex_music_project.csv`. О качестве данных ничего не известно. Поэтому перед проверкой гипотез понадобится обзор данных.

Необходимо проверить данные на ошибки и оценить их влияние на исследование. Затем, на этапе предобработки следует поискать возможность исправить самые критичные ошибки данных.

Таким образом, исследование пройдёт в три этапа:

- Обзор данных.
- Предобработка данных.
- Проверка гипотез.

Обзор данных

Составим первое представление о данных Яндекс.Музыки.

Основной инструмент аналитика — `pandas`. Импортируем эту библиотеку.

```
In [1]: # импорт библиотеки pandas
import pandas as pd

Прочитаем файл yandex_music_project.csv из папки /datasets и сохраним его в переменной df :

In [2]: # чтение файла с данными и сохранение в df
df = pd.read_csv('/datasets/yandex_music_project.csv')

Выведем на экран первые десять строк таблицы:

In [3]: # получение первых 10 строк таблицы df
df.head(10)
```

	userID	Track	artist	genre	City	time	Day
0	FFB692EC	Kamigata To Boots	The Mass Missile	rock	Saint-Petersburg	20:28:33	Wednesday
1	55204538	Delayed Because of Accident	Andreas Rönnerberg	rock	Moscow	14:07:09	Friday
2	20EC38	Funiculi funiculà	Mario Lanza	pop	Saint-Petersburg	20:58:07	Wednesday
3	A3DD03C9	Dragons in the Sunset	Fire + Ice	folk	Saint-Petersburg	08:37:09	Monday
4	E2DC1FAE	Soul People	Space Echo	dance	Moscow	08:34:34	Monday
5	842029A1	Преданная	IMPERVITOR	rusrap	Saint-Petersburg	13:09:41	Friday
6	4CB90AA5	True	Roman Messer	dance	Moscow	13:00:07	Wednesday
7	F03E1C1F	Feeling This Way	Polina Griffith	dance	Moscow	20:47:49	Wednesday
8	8FA1D3BE	И вновь продолжается бой	NaN	ruspop	Moscow	09:17:40	Friday
9	E772D5C0	Pessimist	NaN	dance	Saint-Petersburg	21:20:49	Wednesday

Получим общую информацию о таблице:

```
In [4]: # получение общей информации о данных в таблице df
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 65079 entries, 0 to 65078
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  -
0   userID      65079 non-null  object
1   Track       63848 non-null  object
2   artist      57876 non-null  object
3   genre       63881 non-null  object
4   City        65079 non-null  object
5   time        65079 non-null  object
6   Day         65079 non-null  object
dtypes: object(7)
memory usage: 3.5+ MB

Итак, в таблице семь столбцов. Тип данных во всех столбцах — object .
```

Согласно документации к данным:

- `userID` — идентификатор пользователя;
- `Track` — название трека;
- `artist` — имя исполнителя;
- `genre` — название жанра;
- `City` — город пользователя;
- `time` — время начала прослушивания;
- `Day` — день недели.

В названиях колонок видны три нарушения стиля:

1. Строчные буквы сочетаются с прописными.
2. Встречаются пробелы.
3. Название колонки `userID` написано не в `snake_case`.

Количество значений в столбцах различается. Значит, в данных есть пропущенные значения.

Выводы

В каждой строке таблицы — данные о прослушанном треке. Часть колонок описывает саму композицию: название, исполнителя и жанр. Остальные данные рассказывают о пользователе: из какого он города, когда он слушал музыку.

Предварительно можно утверждать, что, данных достаточно для проверки гипотез. Но встречаются пропуски в данных, а в названиях колонок — расхождения с хорошим стилем.

Пропуски в данных наблюдаются только в колонках, описывающих саму композицию, что может свидетельствовать об ошибках при извлечении данных о композициях, либо о том, что в системе поля `Track`, `artist` и `genre` являются необязательными и теоретически возможна ситуация, когда все три поля будут пропущены одновременно. *Целесообразно проконсультироваться с разработчиками.*

С точки зрения постановки задачи наиболее важным является столбец `genre`. Предварительно, он содержит меньше всего пропусков данных.

Чтобы двигаться дальше, устраним проблемы в данных.

Предобработка данных

Исправим стиль в заголовках столбцов, исключим пропуски. Затем проверим данные на дубликаты.

Стиль заголовков

Выведем на экран названия столбцов:

```
In [5]: # перечень названий столбцов таблицы df
df.columns

Out[5]: Index(['userID', 'Track', 'artist', 'genre', 'City', 'time', 'Day'], dtype='object')
```

Приведём названия в соответствие с хорошим стилем:

- несколько слов в названии запишите в «змеином_регистре»,
- все символы сделайте строчными,
- уберите пробелы.

Для этого переименуем колонки:

- `'userID' → 'user_id'`;
- `'Track' → 'track'`;
- `'City' → 'city'`;
- `'Day' → 'day'`.

```
In [6]: # переименование столбцов
df = df.rename(
    columns={
        'userID' : 'user_id',
        'Track' : 'track',
        'City' : 'city',
        'Day' : 'day'
    }
)
```

Проверим результат:

```
In [7]: # проверка результатов - перечень названий столбцов
df.columns

Out[7]: Index(['user_id', 'track', 'artist', 'genre', 'city', 'time', 'day'], dtype='object')
```

Пропуски значений

Сначала посчитаем, сколько в таблице пропущенных значений:

```
In [8]: # подсчёт пропусков
df.isna().sum()

Out[8]: user_id      0
track      1231
artist     7203
genre      1198
city        0
time        0
day         0
dtype: int64
```

Не все пропущенные значения влияют на исследование. Так в `track` и `artist` пропуски не важны, достаточно заменить их явными обозначениями.

Но пропуски в `genre` могут мешать сравнению музыкальных вкусов в Москве и Санкт-Петербурге. На практике было бы правильно установить причину пропусков и восстановить данные, однако такой возможности нет, поэтому придётся:

- заполнить и эти пропуски явными обозначениями,
- оценить, насколько они повредят расчётам.

Заменим пропущенные значения в столбцах `track`, `artist` и `genre` на строку `'unknown'`.

```
In [9]: # перебор названий столбцов в цикле и замена пропущенных значений на 'unknown'
columns_to_replace = ['track', 'artist', 'genre']

for col in columns_to_replace:
    df[col] = df[col].fillna('unknown')
```

Убедимся, что в таблице не осталось пропусков.

```
In [10]: # подсчёт пропусков
df.isna().sum()

Out[10]: user_id      0
track        0
artist        0
genre         0
city          0
time          0
day           0
dtype: int64
```

Дубликаты

Посчитаем явные дубликаты в таблице:

```
In [11]: # подсчёт явных дубликатов
df.duplicated().sum()

Out[11]: 3826
```

Удалим явные дубликаты:

```
In [12]: # удаление явных дубликатов
df = df.drop_duplicates()
```

Ещё раз посчитаем явные дубликаты, чтобы убедиться, что полностью от них избавились:

```
In [13]: # проверка на отсутствие дубликатов
df.duplicated().sum()

Out[13]: 0
```

Теперь избавимся от неявных дубликатов в колонке `genre`. Например, название одного и того же жанра может быть записано немного по-разному. Такие ошибки тоже повлияют на результат исследования.

Выведем на экран список уникальных названий жанров, отсортированный в алфавитном порядке. Для этого:

- извлеките нужный столбец датафрейма,
- примените к нему метод сортировки,
- для отсортированного столбца вызовите метод, который вернёт уникальные значения из столбца.

```
In [14]: # Просмотр уникальных названий жанров
df['genre'].sort_values().unique()
```

```
Out[14]: array(['acid', 'acoustic', 'action', 'adult', 'africa', 'afrikaans',
               'alternative', 'alternativepunk', 'ambient', 'americana',
               'animated', 'anime', 'arabesque', 'arabic', 'arena',
               'argentinatango', 'art', 'audiobook', 'author', 'avantgarde',
               'axe', 'baile', 'balkan', 'beats', 'bigroom', 'black', 'bluegrass',
               'blues', 'bollywood', 'bossa', 'brazilian', 'breakbeat', 'breaks',
               'broadway', 'cantautori', 'cantopop', 'canzone', 'caribbean',
               'caucasian', 'celtic', 'chamber', 'chanson', 'children', 'chill',
               'chinese', 'choral', 'christian', 'christmas', 'classical',
               'classicmetal', 'club', 'colombian', 'comedy', 'conjazz',
               'contemporary', 'country', 'cuban', 'dance', 'dancehall',
               'dancepop', 'dark', 'death', 'deep', 'deutschrock', 'deutschspr',
               'dirty', 'disco', 'dnb', 'documentary', 'downbeat', 'downtempo',
               'drum', 'dub', 'dubstep', 'eastern', 'easy', 'electronic',
               'electropop', 'emo', 'entehno', 'epicmetal', 'estrada', 'ethnic',
               'eurofolk', 'european', 'experimental', 'extrememetal', 'fado',
               'fairytail', 'film', 'fitness', 'flamenco', 'folk', 'folklore',
               'folkmatal', 'folkrock', 'folktronica', 'forró', 'frankreich',
               'französisch', 'french', 'funk', 'future', 'gangsta', 'garage',
               'german', 'ghazal', 'gitarre', 'glitch', 'gospel', 'gothic',
               'grime', 'grunge', 'gypsy', 'handsup', 'hard'n'heavy', 'hardcore',
               'hardstyle', 'hardtechno', 'hip', 'hip-hop', 'hiphop',
               'historisch', 'holiday', 'hop', 'horror', 'house', 'hymn', 'idm',
               'independent', 'indian', 'indie', 'indipop', 'industrial',
               'inspirational', 'instrumental', 'international', 'irish', 'jam',
               'japanese', 'jazz', 'jewish', 'jpop', 'jungle', 'k-pop',
               'karadeniz', 'karaoke', 'kayokyoku', 'korean', 'laiko', 'latin',
               'latino', 'leftfield', 'local', 'lounge', 'loungeelectronic',
               'lovers', 'malaysian', 'mandopop', 'marschmusik', 'meditative',
               'mediterranean', 'melodic', 'metal', 'metalcore', 'mexican',
               'middle', 'minimal', 'miscellaneous', 'modern', 'mood', 'mpb',
               'muslim', 'native', 'neoklassik', 'neue', 'new', 'newage',
               'newwave', 'nu', 'nujazz', 'numetal', 'oceania', 'old', 'opera',
               'orchestral', 'other', 'piano', 'podcasts', 'pop', 'popdance',
               'popelectronic', 'popeurodance', 'poprussian', 'post',
               'posthardcore', 'postrock', 'power', 'progmetal', 'progressive',
               'psychedelic', 'punjabi', 'punk', 'quebecois', 'ragga', 'ram',
               'rancheras', 'rap', 'rave', 'reggae', 'reggaeton', 'regional',
               'relax', 'religious', 'retro', 'rhythm', 'rnb', 'rnn', 'rock',
               'rockabilly', 'rockalternative', 'rockindie', 'rockother',
               'romance', 'roots', 'ruspop', 'rusrap', 'rusrock', 'russian',
               'salsa', 'samba', 'scenic', 'schlager', 'self', 'sertanejo',
               'shanson', 'shoegazing', 'showtunes', 'singer', 'ska', 'skarock',
               'slow', 'smooth', 'soft', 'soul', 'soulful', 'sound', 'soundtrack',
               'southern', 'specialty', 'speech', 'spiritual', 'sport',
               'stonerrock', 'surf', 'swing', 'synthpop', 'synthrock',
               'sängerportrait', 'tango', 'tanzorchester', 'taraftar', 'tatar',
               'tech', 'techno', 'teen', 'thrash', 'top', 'traditional',
               'tradjazz', 'trance', 'tribal', 'trip', 'triphop', 'tropical',
               'türk', 'türkçe', 'ukrrock', 'unknown', 'urban', 'uzbek',
               'variété', 'vi', 'videogame', 'vocal', 'western', 'world',
               'worldbeat', 'iii', 'электроника'], dtype=object)
```

В списке выявлены следующие неявные дубликаты:

- Альтернативные названия и названия с ошибками одного и того же жанра:
 - * hip*,
 - * hop*,
 - * hiphop*,
 - * hip-hop*.
- Названия одного и того же жанра на разных языках:
 - * электроника*
 - * electronic*

Чтобы очистить от них таблицу, напомним функцию `replace_wrong_genres()` с двумя параметрами:

- `wrong_genres` — список дубликатов,
- `correct_genre` — строка с правильным значением.

Функция исправляет колонку `genre` в таблице `df`: заменяет каждое значение из списка `wrong_genres` на значение из `correct_genre`.

```
In [15]: # функция для замены неявных дубликатов
def replace_wrong_genres(wrong_genres, correct_genre):
    for wrong_genre in wrong_genres:
        df['genre'] = df['genre'].replace(wrong_genre, correct_genre)
```

Устраним неявные дубликаты: вместо `hip`, `hop` и `hip-hop` в таблице должно быть значение `hiphop`; вместо `электроника` - `electronic`:

```
In [16]: # Устранение неявных дубликатов

genres_to_replace = {
    'hiphop': ['hip', 'hop', 'hip-hop'],
    'electronic': ['электроника']
}

for correct_value in genres_to_replace:
    replace_wrong_genres(genres_to_replace[correct_value], correct_value)
```

Проверим корректность устранения неявных дубликатов. Выведем отсортированный список уникальных значений столбца `genre`:

```
In [17]: # Проверка на неявные дубликаты
df['genre'].sort_values().unique()
```

```
Out[17]: array(['acid', 'acoustic', 'action', 'adult', 'africa', 'afrikaans',
        'alternative', 'alternativepunk', 'ambient', 'americana',
        'animated', 'anime', 'arabesk', 'arabic', 'arena',
        'argentinetango', 'art', 'audiobook', 'author', 'avantgarde',
        'axe', 'baile', 'balkan', 'beats', 'bigroom', 'black', 'bluegrass',
        'blues', 'bollywood', 'bossa', 'brazilian', 'breakbeat', 'breaks',
        'broadway', 'cantautori', 'cantopop', 'canzone', 'caribbean',
        'caucasian', 'celtic', 'chamber', 'chanson', 'children', 'chill',
        'chinese', 'choral', 'christian', 'christmas', 'classical',
        'classicmetal', 'club', 'colombian', 'comedy', 'conjazz',
        'contemporary', 'country', 'cuban', 'dance', 'dancehall',
        'dancepop', 'dark', 'death', 'deep', 'deutschrack', 'deutschspr',
        'dirty', 'disco', 'dnb', 'documentary', 'downbeat', 'downtempo',
        'drum', 'dub', 'dubstep', 'eastern', 'easy', 'electronic',
        'electropop', 'emo', 'entehno', 'epicmetal', 'estrada', 'ethnic',
        'eurofolk', 'european', 'experimental', 'extrememetal', 'fado',
        'fairytail', 'film', 'fitness', 'flamenco', 'folk', 'folklore',
        'folkmatal', 'folkrock', 'folktronica', 'forró', 'frankreich',
        'französisch', 'french', 'funk', 'future', 'gangsta', 'garage',
        'german', 'ghazal', 'gitarre', 'glitch', 'gospel', 'gothic',
        'grime', 'grunge', 'gypsy', 'handsup', "hard'n'heavy", 'hardcore',
        'hardstyle', 'hardtechno', 'hiphop', 'historisch', 'holiday',
        'horror', 'house', 'hymn', 'idm', 'independent', 'indian', 'indie',
        'indipop', 'industrial', 'inspirational', 'instrumental',
        'international', 'irish', 'jam', 'japanese', 'jazz', 'jewish',
        'jpop', 'jungle', 'k-pop', 'karadeniz', 'karaoke', 'kayokyoku',
        'korean', 'laiko', 'latin', 'latino', 'leftfield', 'local',
        'lounge', 'loungeselectronic', 'lovers', 'malaysian', 'mandopop',
        'marschmusik', 'meditative', 'mediterranean', 'melodic', 'metal',
        'metacore', 'mexican', 'middle', 'minimal', 'miscellaneous',
        'modern', 'mood', 'mpb', 'muslim', 'native', 'neoklassik', 'neue',
        'new', 'newage', 'newwave', 'nu', 'nujazz', 'numetal', 'oceania',
        'old', 'opera', 'orchestral', 'other', 'piano', 'podcasts', 'pop',
        'popdance', 'popelectronic', 'popeurodance', 'poprussian', 'post',
        'posthardcore', 'postrock', 'power', 'progmetal', 'progressive',
        'psychedelic', 'punjabi', 'punk', 'quebecois', 'ragga', 'ram',
        'rancheras', 'rap', 'rave', 'reggae', 'reggaeton', 'regional',
        'relax', 'religious', 'retro', 'rhythm', 'rnb', 'rnn', 'rock',
        'rockabilly', 'rockalternative', 'rockindie', 'rockother',
        'romance', 'roots', 'ruspop', 'rusrap', 'rusrock', 'russian',
        'salsa', 'samba', 'scenic', 'schlager', 'self', 'sertanejo',
        'shanson', 'shoegazing', 'showtunes', 'singer', 'ska', 'skarock',
        'slow', 'smooth', 'soft', 'soul', 'soulful', 'sound', 'soundtrack',
        'southern', 'specialty', 'speech', 'spiritual', 'sport',
        'stonerrock', 'surf', 'swing', 'synthpop', 'synthrock',
        'sängerportrait', 'tango', 'tanzorchester', 'taraftar', 'tatar',
        'tech', 'techno', 'teen', 'thrash', 'top', 'traditional',
        'tradjazz', 'trance', 'tribal', 'trip', 'triphop', 'tropical',
        'türk', 'türkçe', 'ukrrock', 'unknown', 'urban', 'uzbek',
        'variété', 'vi', 'videogame', 'vocal', 'western', 'world',
        'worldbeat', 'iii'], dtype=object)
```

Выводы

Предобработка обнаружила три проблемы в данных:

- нарушения в стиле заголовков,
- пропущенные значения,
- дубликаты — явные и неявные.

Мы исправили заголовки, чтобы упростить работу с таблицей. Без дубликатов исследование станет более точным.

Пропущенные значения мы заменили на `'unknown'`. Ещё предстоит увидеть, не повредят ли исследованию пропуски в колонке `genre`.

Теперь можно перейти к проверке гипотез.

Проверка гипотез

Сравнение поведения пользователей двух столиц

Первая гипотеза утверждает, что пользователи по-разному слушают музыку в Москве и Санкт-Петербурге. Проверим это предположение по данным о трёх днях недели — понедельник, среде и пятницу. Для этого:

- Разделим пользователей Москвы и Санкт-Петербурга
- Сравним, сколько треков послушала каждая группа пользователей в понедельник, среду и пятницу.

Оценим активность пользователей в каждом городе. Сгруппируем данные по городу и посчитаем прослушивания в каждой группе.

```
In [18]: # Подсчёт прослушиваний в каждом городе
df.groupby('city')['user_id'].count()
```

```
Out[18]: city
Moscow      42741
Saint-Petersburg  18512
Name: user_id, dtype: int64
```

В Москве прослушиваний больше, чем в Петербурге. Из этого не следует, что московские пользователи чаще слушают музыку. Просто самих пользователей в Москве больше.

Теперь сгруппируем данные по дню недели и подсчитаем прослушивания в понедельник, среду и пятницу.

```
In [19]: # Подсчёт прослушиваний в каждый из трёх дней
df.groupby('day')['user_id'].count()
```

```
Out[19]: day
Friday      21840
Monday       21354
Wednesday   18059
Name: user_id, dtype: int64
```

В среднем пользователи из двух городов менее активны по средам. Но картина может измениться, если рассмотреть каждый город в отдельности.

Объединим эти два расчёта. Создадим функцию `number_tracks()`, которая посчитает прослушивания для заданного дня и города. Ей понадобятся два параметра:

- день недели (`day`),
- название города (`city`).

Функция возвращает количество значений в столбце `user_id` для заданных `day` и `city`.

```
In [20]: # Функция для подсчёта прослушиваний для конкретного города и дня.
def number_tracks(day, city):
    track_list = df[df['day'] == day]
    track_list = track_list[track_list['city'] == city]

    track_list_count = track_list['user_id'].count()
    return track_list_count
```

Вызовем `number_tracks()`, меняя значение параметров так, чтобы получить данные для каждого города в каждый из трёх дней.

```
In [21]: # количество прослушиваний в Москве по понедельникам
number_tracks('Monday', 'Moscow')

Out[21]: 15740
```

```
In [22]: # количество прослушиваний в Санкт-Петербурге по понедельникам
number_tracks('Monday', 'Saint-Petersburg')

Out[22]: 5614
```

```
In [23]: # количество прослушиваний в Москве по средам
number_tracks('Wednesday', 'Moscow')

Out[23]: 11056
```

```
In [24]: # количество прослушиваний в Санкт-Петербурге по средам
number_tracks('Wednesday', 'Saint-Petersburg')

Out[24]: 7003
```

```
In [25]: # количество прослушиваний в Москве по пятницам
number_tracks('Friday', 'Moscow')

Out[25]: 15945
```

```
In [26]: # количество прослушиваний в Санкт-Петербурге по пятницам
number_tracks('Friday', 'Saint-Petersburg')

Out[26]: 5895
```

Создадим с помощью конструктора `pd.DataFrame` таблицу, где

- названия колонок — `['city', 'monday', 'wednesday', 'friday']`;
- данные — результаты, которые вы получили с помощью `number_tracks`.

```
In [27]: # Таблица с результатами
data = [['Moscow', 15740, 11056, 15945],
        ['Saint-Petersburg', 5614, 7003, 5895]]
col = ['city', 'monday', 'wednesday', 'friday']

pd.DataFrame(data, columns=col)
```

Out[27]:

	city	monday	wednesday	friday
0	Moscow	15740	11056	15945
1	Saint-Petersburg	5614	7003	5895

Выводы

Данные показывают разницу поведения пользователей:

- В Москве пик прослушиваний приходится на понедельник и пятницу, а в среду замечен спад.
- В Петербурге, наоборот, больше слушают музыку по средам. Активность в понедельник и пятницу здесь почти в равной мере уступает среде.

Значит, данные говорят в пользу первой гипотезы.

Музыка в начале и в конце недели

Согласно второй гипотезе, утром в понедельник в Москве преобладают одни жанры, а в Петербурге — другие. Так же и вечером пятницы преобладают разные жанры — в зависимости от города.

Сохраним таблицы с данными в две переменные:

- по Москве — в `moscow_general`;
- по Санкт-Петербургу — в `spb_general`.

```
In [28]: # получение таблицы moscow_general из тех строк таблицы df,
# для которых значение в столбце 'city' равно 'Moscow'
moscow_general = df[df['city'] == 'Moscow']
```

```
In [29]: # получение таблицы spb_general из тех строк таблицы df,
# для которых значение в столбце 'city' равно 'Saint-Petersburg'
spb_general = df[df['city'] == 'Saint-Petersburg']
```

Создадим функцию `genre_weekday()` с четырьмя параметрами:

- таблица (датафрейм) с данными,
- день недели,
- начальная временная метка в формате 'hh:mm',

- последняя временная метка в формате 'hh:mm'.

Функция должна вернуть информацию о топ-10 жанров тех треков, которые прослушивали в указанный день, в промежутке между двумя отметками времени.

```
In [30]: # Объявление функции которая возвращает информацию о самых популярных жанрах
# в указанный день в заданное время
def genre_weekday(table, day, time1, time2):
    genre_df = table[table['day'] == day]
    genre_df = genre_df[genre_df['time'] > time1]
    genre_df = genre_df[genre_df['time'] < time2]

    genre_df_count = genre_df.groupby('genre')['genre'].count()
    genre_df_sorted = genre_df_count.sort_values(ascending=False)
    return genre_df_sorted.head(10)
```

Сравним результаты функции `genre_weekday()` для Москвы и Санкт-Петербурга в понедельник утром (с 7:00 до 11:00) и в пятницу вечером (с 17:00 до 23:00):

```
In [31]: # вызов функции для утра понедельника в Москве
genre_weekday(moscow_general, 'Monday', '07:00', '11:00')
```

```
Out[31]: genre
pop      781
dance    549
electronic 480
rock     474
hiphop   286
ruspop   186
world    181
rusrap   175
alternative 164
unknown  161
Name: genre, dtype: int64
```

```
In [32]: # вызов функции для утра понедельника в Петербурге (вместо df – таблица spb_general)
genre_weekday(spb_general, 'Monday', '07:00', '11:00')
```

```
Out[32]: genre
pop      218
dance    182
rock     162
electronic 147
hiphop    80
ruspop    64
alternative 58
rusrap    55
jazz      44
classical  40
Name: genre, dtype: int64
```

```
In [33]: # вызов функции для вечера пятницы в Москве
genre_weekday(moscow_general, 'Friday', '17:00', '23:00')
```

```
Out[33]: genre
pop      713
rock     517
dance    495
electronic 482
hiphop   273
world    208
ruspop   170
alternative 163
classical 163
rusrap   142
Name: genre, dtype: int64
```

```
In [34]: # вызов функции для вечера пятницы в Петербурге
genre_weekday(spb_general, 'Friday', '17:00', '23:00')
```

```
Out[34]: genre
pop      256
electronic 216
rock     216
dance    210
hiphop    97
alternative 63
jazz      61
classical  60
rusrap    59
world     54
Name: genre, dtype: int64
```

Выводы

Если сравнить топ-10 жанров в понедельник утром, можно сделать такие выводы:

1. В Москве и Петербурге слушают похожую музыку. Единственное отличие — в московский рейтинг вошёл жанр “world”, а в петербургский — джаз и классика.
2. В Москве пропущенных значений оказалось так много, что значение `'unknown'` заняло десятое место среди самых популярных жанров. Значит, пропущенные значения занимают существенную долю в данных и угрожают достоверности исследования.

Вечер пятницы не меняет эту картину. Некоторые жанры поднимаются немного выше, другие спускаются, но в целом топ-10 остаётся тем же самым.

Таким образом, вторая гипотеза подтвердилась лишь частично:

- Пользователи слушают похожую музыку в начале недели и в конце.
- Разница между Москвой и Петербургом не слишком выражена. В Москве чаще слушают русскую популярную музыку, в Петербурге — джаз.

Однако пропуски в данных ставят под сомнение этот результат. В Москве их так много, что рейтинг топ-10 мог бы выглядеть иначе, если бы не утерянные данные о жанрах.

Жанровые предпочтения в Москве и Петербурге

Гипотеза: Петербург — столица рэпа, музыку этого жанра там слушают чаще, чем в Москве. А Москва — город контрастов, в котором, тем не менее, преобладает поп-музыка.

Сгруппируем таблицу `moscow_general` по жанру и посчитаем прослушивания треков каждого жанра методом `count()`. Затем отсортируем результат в порядке убывания и сохраним его в таблице `moscow_genres`.

```
In [35]: # группировка таблицы moscow_general
moscow_genres = moscow_general.groupby('genre')['genre'].count().sort_values(ascending=False)
```

Выведем на экран первые десять строк `moscow_genres`:

```
In [36]: # просмотр первых 10 строк moscow_genres
moscow_genres.head(10)
```

```
Out[36]: genre
pop      5892
dance    4435
rock     3965
electronic 3786
hiphop    2096
classical 1616
world     1432
alternative 1379
ruspop    1372
rusrap     1161
Name: genre, dtype: int64
```

Теперь повторим то же для Петербурга.

Сгруппируем таблицу `spb_general` по жанру. Посчитаем прослушивания треков каждого жанра. Результат отсортируем в порядке убывания и сохраним в таблице `spb_genres`:

```
In [37]: # группировка таблицы spb_general
spb_genres = spb_general.groupby('genre')['genre'].count().sort_values(ascending=False)
```

Выведем на экран первые десять строк `spb_genres`:

```
In [38]: # просмотр первых 10 строк spb_genres
spb_genres.head(10)
```

```
Out[38]: genre
pop      2431
dance    1932
rock     1879
electronic 1737
hiphop     960
alternative 649
classical  646
rusrap     564
ruspop     538
world      515
Name: genre, dtype: int64
```

Выводы

Гипотеза частично подтвердилась:

- Поп-музыка — самый популярный жанр в Москве, как и предполагала гипотеза. Более того, в топ-10 жанров встречается близкий жанр — русская популярная музыка.
- Вопреки ожиданиям, рэп приблизительно одинаково популярен в Москве и Петербурге. В Москве данный жанр замыкает ТОП-10, в Петербурге - входит в его вторую половину, т.е. незначительно популярнее.

Итоги исследования

Мы проверили три гипотезы и установили:

1. День недели по-разному влияет на активность пользователей в Москве и Петербурге.

Первая гипотеза полностью подтвердилась.

1. Музыкальные предпочтения не сильно меняются в течение недели — будь то Москва или Петербург. Небольшие различия заметны в начале недели, по понедельникам:
 - в Москве слушают музыку жанра “world”,
 - в Петербурге — джаз и классику.

Таким образом, вторая гипотеза подтвердилась лишь отчасти. Этот результат мог оказаться иным, если бы не пропуски в данных.

1. Во вкусах пользователей Москвы и Петербурга больше общего чем различий. Вопреки ожиданиям, предпочтения жанров в Петербурге напоминают московские.

Третья гипотеза не подтвердилась. Если различия в предпочтениях и существуют, на основной массе пользователей они незаметны.