

# Проект по SQL

Коронавирус изменил привычный порядок вещей: в свободное время жители городов стали проводить больше времени за книгами, что привело к появлению большого числа приложений для тех, кто любит читать.

Наша компания решила быть на волне и купила крупный сервис для чтения книг по подписке.

**Задача:** проанализировать базу данных.

**Цель:** сформулировать предложения для нового продукта.

## Описание данных

### Таблица books

Содержит данные о книгах:

- book\_id — идентификатор книги;
- author\_id — идентификатор автора;
- title — название книги;
- num\_pages — количество страниц;
- publication\_date — дата публикации книги;
- publisher\_id — идентификатор издателя.

### Таблица authors

Содержит данные об авторах:

- author\_id — идентификатор автора;
- author — имя автора.

### Таблица publishers

Содержит данные об издательствах:

- publisher\_id — идентификатор издательства;
- publisher — название издательства;

### Таблица ratings

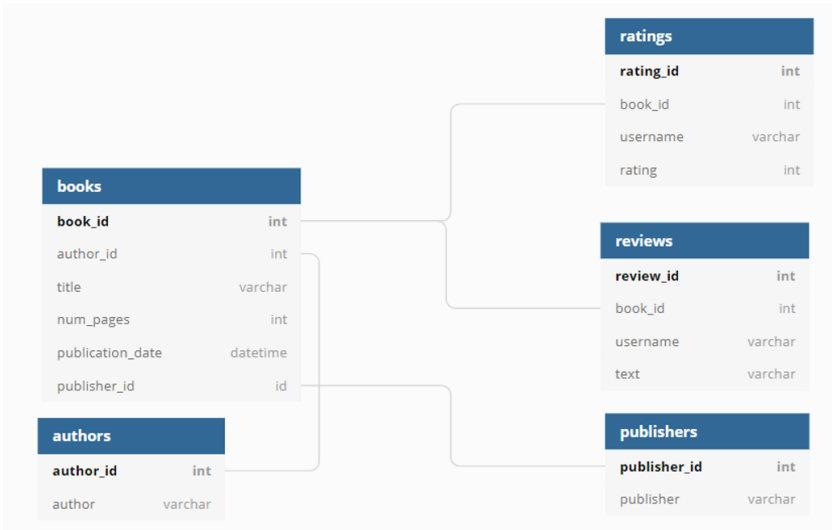
Содержит данные о пользовательских оценках книг:

- rating\_id — идентификатор оценки;
- book\_id — идентификатор книги;
- username — имя пользователя, оставившего оценку;
- rating — оценка книги.

### Таблица reviews

Содержит данные о пользовательских обзорах:

- review\_id — идентификатор обзора;
- book\_id — идентификатор книги;
- username — имя автора обзора;
- text — текст обзора.



# Содержание

[1 Подключение к базе данных](#)

[2 Исследование таблиц \(обзор\)](#)

## ▼ 3 Задания

[3.1 Количество книг вышедших после 1 января 2000 года](#)

[3.2 Количество обзоров и средняя оценка для каждой книги](#)

[3.3 Издательство, которое выпустило наибольшее число книг толще 50 страниц \(так были исключены из анализа брошюры\)](#)

[3.4 Автор с самой высокой средней оценкой книг \(книги с 50 и более оценками\)](#)

[3.5 Среднее количество обзоров от пользователей, которые поставили больше 48 оценок](#)

## 1 Подключение к базе данных

In [1]:

```
1  # Импорт библиотек
2  import pandas as pd
3  from sqlalchemy import text, create_engine
4
5  # Устанавливаем параметры подключения
6  db_config = {'user': '', # имя пользователя
7  'pwd': '', # пароль
8  'host': '',
9  'port': , # порт подключения
10 'db': '' } # название базы данных
11
12 connection_string = 'postgresql://{user}:{pwd}@{host}:{port}/{db}'.format(**db_config)
13
14 # Сохраняем коннектор
15 engine = create_engine(connection_string, connect_args={'sslmode': 'require'})
16
17 con=engine.connect()
```

## 2 Исследование таблиц (обзор)

```
In [2]: 1 # Выполняем запросы SQL и получаем информацию о датасетах
2 for value in ['books', 'authors', 'publishers', 'ratings', 'reviews']:
3     print('\033[31m' + '_____Выполняем запрос_____ ' + '\033[0m')
4     print('Название таблицы:', value)
5
6     query = '''
7     SELECT *
8     FROM {}
9     LIMIT 5;
10    '''.format(value)
11
12    df = pd.io.sql.read_sql(query, con = engine)
13    display(df)
14
15    print('\033[31m' + '_____Дубликаты_____ ' + '\033[0m')
16    print('Дубликаты:', df.duplicated().sum())
17
18    print('\033[31m' + '_____Информация о датасете_____ ' + '\033[0m')
19    query = '''
20    SELECT * FROM {};
21    '''.format(value)
22
23    df = pd.io.sql.read_sql(query, con = engine)
24
25    print(df.info())
26    print('\033[34m' + '_*50 + '\033[0m')
27    print()
```

\_\_\_\_\_Выполняем запрос\_\_\_\_\_

Название таблицы: books

	book_id	author_id		title	num_pages	publication_date	publisher_id
0	1	546		'Salem's Lot	594	2005-11-01	93
1	2	465		1 000 Places to See Before You Die	992	2003-05-22	336
2	3	407		13 Little Blue Envelopes (Little Blue Envelope...	322	2010-12-21	135
3	4	82		1491: New Revelations of the Americas Before C...	541	2006-10-10	309
4	5	125		1776	386	2006-07-04	268

\_\_\_\_\_Дубликаты\_\_\_\_\_

Дубликаты: 0

\_\_\_\_\_Информация о датасете\_\_\_\_\_

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   book_id         1000 non-null   int64
1   author_id       1000 non-null   int64
2   title           1000 non-null   object
3   num_pages       1000 non-null   int64
4   publication_date 1000 non-null   object
5   publisher_id    1000 non-null   int64
dtypes: int64(4), object(2)
memory usage: 47.0+ KB
None
```

\_\_\_\_\_Выполняем запрос\_\_\_\_\_

Название таблицы: authors

	author_id	author
0	1	A.S. Byatt
1	2	Aesop/Laura Harris/Laura Gibbs
2	3	Agatha Christie
3	4	Alan Brennert
4	5	Alan Moore/David Lloyd

\_\_\_\_\_Дубликаты\_\_\_\_\_

Дубликаты: 0

\_\_\_\_\_Информация о датасете\_\_\_\_\_

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 636 entries, 0 to 635
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   author_id       636 non-null   int64
1   author          636 non-null   object
dtypes: int64(1), object(1)
memory usage: 10.1+ KB
None
```

\_\_\_\_\_Выполняем запрос\_\_\_\_\_

Название таблицы: publishers

	<b>publisher_id</b>	<b>publisher</b>
0	1	Ace
1	2	Ace Book
2	3	Ace Books
3	4	Ace Hardcover
4	5	Addison Wesley Publishing Company

#### Дубликаты

Дубликаты: 0

#### Информация о датасете

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 340 entries, 0 to 339
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   publisher_id 340 non-null    int64
1   publisher    340 non-null    object
dtypes: int64(1), object(1)
memory usage: 5.4+ KB
None
```

#### Выполняем запрос

Название таблицы: ratings

	<b>rating_id</b>	<b>book_id</b>	<b>username</b>	<b>rating</b>
0	1	1	ryanfranco	4
1	2	1	grantpatricia	2
2	3	1	brandtandrea	5
3	4	2	lorichen	3
4	5	2	mariokeller	2

#### Дубликаты

Дубликаты: 0

#### Информация о датасете

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6456 entries, 0 to 6455
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   rating_id   6456 non-null    int64
1   book_id     6456 non-null    int64
2   username    6456 non-null    object
3   rating      6456 non-null    int64
dtypes: int64(3), object(1)
memory usage: 201.9+ KB
None
```

#### Выполняем запрос

Название таблицы: reviews

	<b>review_id</b>	<b>book_id</b>	<b>username</b>	<b>text</b>
0	1	1	brandtandrea	Mention society tell send professor analysis. ...
1	2	1	ryanfranco	Foot glass pretty audience hit themselves. Amo...
2	3	2	lorichen	Listen treat keep worry. Miss husband tax but ...
3	4	3	johnsonamanda	Finally month interesting blue could nature cu...
4	5	3	scotttamara	Nation purpose heavy give wait song will. List...

#### Дубликаты

Дубликаты: 0

#### Информация о датасете

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2793 entries, 0 to 2792
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   review_id   2793 non-null    int64
1   book_id     2793 non-null    int64
2   username    2793 non-null    object
3   text        2793 non-null    object
dtypes: int64(2), object(2)
memory usage: 87.4+ KB
None
```

### 3 Задания

#### 3.1 Количество книг вышедших после 1 января 2000 года

```
In [3]: 1 # Задание №1
2 # Выполняем запрос SQL
3 query = '''
4 SELECT COUNT(*) AS cnt_books_after_2000
5 FROM books
6 WHERE publication_date > '2000-01-01';
7
8 '''
9
10 pd.io.sql.read_sql(sql=text(query), con = con)
```

```
Out[3]: cnt_books_after_2000
0      819
```

##### Заметка:

С 1-го января 2000-го года было опубликовано 819 из 1000 книг

#### 3.2 Количество обзоров и средняя оценка для каждой книги

```
In [4]: 1 # Задание №2
2 # Выполняем запрос SQL
3 query = '''
4 SELECT
5     rat.book_id,
6     COUNT(DISTINCT rev.text) as cnt_texts,
7     AVG(rat.rating) AS avg_rating
8 FROM ratings AS rat
9 INNER JOIN reviews AS rev ON rat.book_id = rev.book_id
10 GROUP BY rat.book_id
11 ORDER BY avg_rating DESC, cnt_texts DESC
12 '''
13
14 pd.io.sql.read_sql(sql=text(query), con = con)
```

```
Out[4]: book_id  cnt_texts  avg_rating
0         17         4         5.00
1        553         3         5.00
2        444         3         5.00
3         76         2         5.00
4        993         2         5.00
...      ...      ...      ...
989       915         3         2.25
990       202         3         2.00
991       371         2         2.00
992       316         2         2.00
993       303         2         1.50
```

994 rows × 3 columns

##### Заметка:

Наибольшая оценка "5".

Отсортировав по количеству отзывов совместно с рейтингом узнаем, что на книгу с id=17 пришлось 4 отзыва

### 3.3 Издательство, которое выпустило наибольшее число книг толще 50 страниц (так были исключены из анализа брошюры)

```
In [5]: 1 # Задание №3
2 # Выполняем запрос SQL
3 query = '''
4 SELECT
5     p.publisher,
6     COUNT(DISTINCT b.book_id) AS cnt_books
7
8 FROM books AS b
9 INNER JOIN publishers AS p ON b.publisher_id = p.publisher_id
10 WHERE b.num_pages > 50
11 GROUP BY b.publisher_id, p.publisher
12 ORDER BY cnt_books DESC
13 LIMIT 1;
14 '''
15
16 pd.io.sql.read_sql(sql=text(query), con = con)
```

```
Out[5]:
```

	publisher	cnt_books
0	Penguin Books	42

#### Заметка:

Наибольшее число книг выпустило издательство Penguin Books

### 3.4 Автор с самой высокой средней оценкой книг (книги с 50 и более оценками)

```
In [6]: 1 # Задание №4↔
Old_version
```

```
In [7]: 1 # Задание №4↔
```

```
Out[7]:
```

	author_id	author	sum_cnt_ratings	avg_rating
0	236	J.K. Rowling/Mary GrandPré	310.0	4.283844

#### Заметка:

Наибольший средний рейтинг у писателя J.K. Rowling/Mary GrandPré

#### ? Идея ?

Можно выяснить также "возраст" (с первой книги до последней + год начала) этого писателя, это позволит оценить также современность стиля, но это уже другое исследование

### 3.5 Среднее количество обзоров от пользователей, которые поставили больше 48 оценок

```
In [8]: 1 # Задание №5
2 # Выполняем запрос SQL
3 query = '''
4 WITH
5 -- первый подзапрос с псевдонимом #1
6 cnt_ratings_tab AS (
7     SELECT
8         username,
9         COUNT(rating) AS cnt_ratings
10
11     FROM ratings AS rat
12
13     GROUP BY username
14 ), -- подзапросы разделяют запятыми
15
16 -- первый подзапрос с псевдонимом #2
17 cnt_reviews_tab AS (
18     SELECT
19         username,
20         COUNT(text) AS cnt_reviews
21
22     FROM reviews AS rev
23
24     GROUP BY username
25 )
26
27 -- основной запрос, в котором указаны псевдонимы для подзапросов
28 SELECT AVG(cnt_reviews_tab.cnt_reviews) AS AVG_number_of_reviews
29
30 FROM cnt_ratings_tab LEFT OUTER JOIN cnt_reviews_tab ON cnt_ratings_tab.username = cnt_reviews_tab.username
31
32 WHERE cnt_ratings_tab.cnt_ratings > 48;
33
34 '''
35
36 pd.io.sql.read_sql(sql=text(query), con = con)
```

```
Out[8]:  avg_number_of_reviews
0          24.0
```

#### Заметка:

Среди активно пользующихся платформой пользователей среднее количество оставляемых обзоров составляет 24 шт.

Активно пользующиеся согласно ТЗ - те, кто оставил более 48 оценок.

#### Итоги:

- С 1-го января 2000-го года было опубликовано 819 из 1000 книг;
- Наибольшая оценка "5";
- книга с id = 17 пришлось 4 отзыва и оценка "5" - стоит рассмотреть подробнее, узнать жанр, провести исследование касательно жанров книг;
- наибольшее число книг выпустило издательство Penguin Books - стоит рассмотреть продажи книг этого издательства, оценить прибыль относительно других издательств, выявить фаворитов;
- наибольший сердний рейтинг у писателя J.K. Rowling/Mary GrandPré - стоит присмотреться к её книгам, выяснить почему они популярны (если популярны) и почему такой рейтинг;
- среди активно пользующихся платформой пользователей среднее количество оставляемых обзоров составляет 24 шт. (активно пользующиеся согласно ТЗ - те, кто оставил более 48 оценок).