

Группа DE\_622, Студент: Виталий Зайцев

Домашнее задание по уроку 5

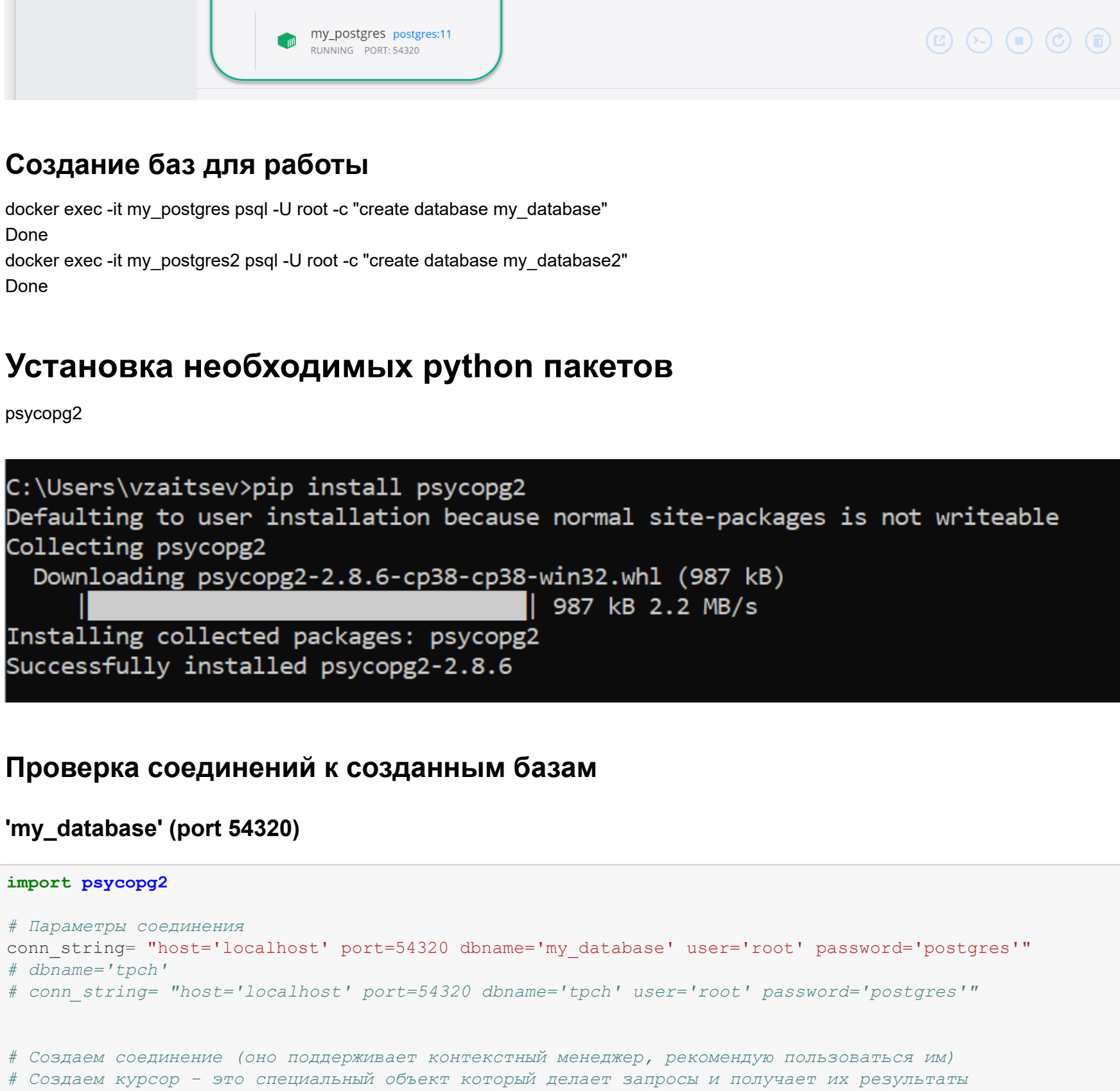
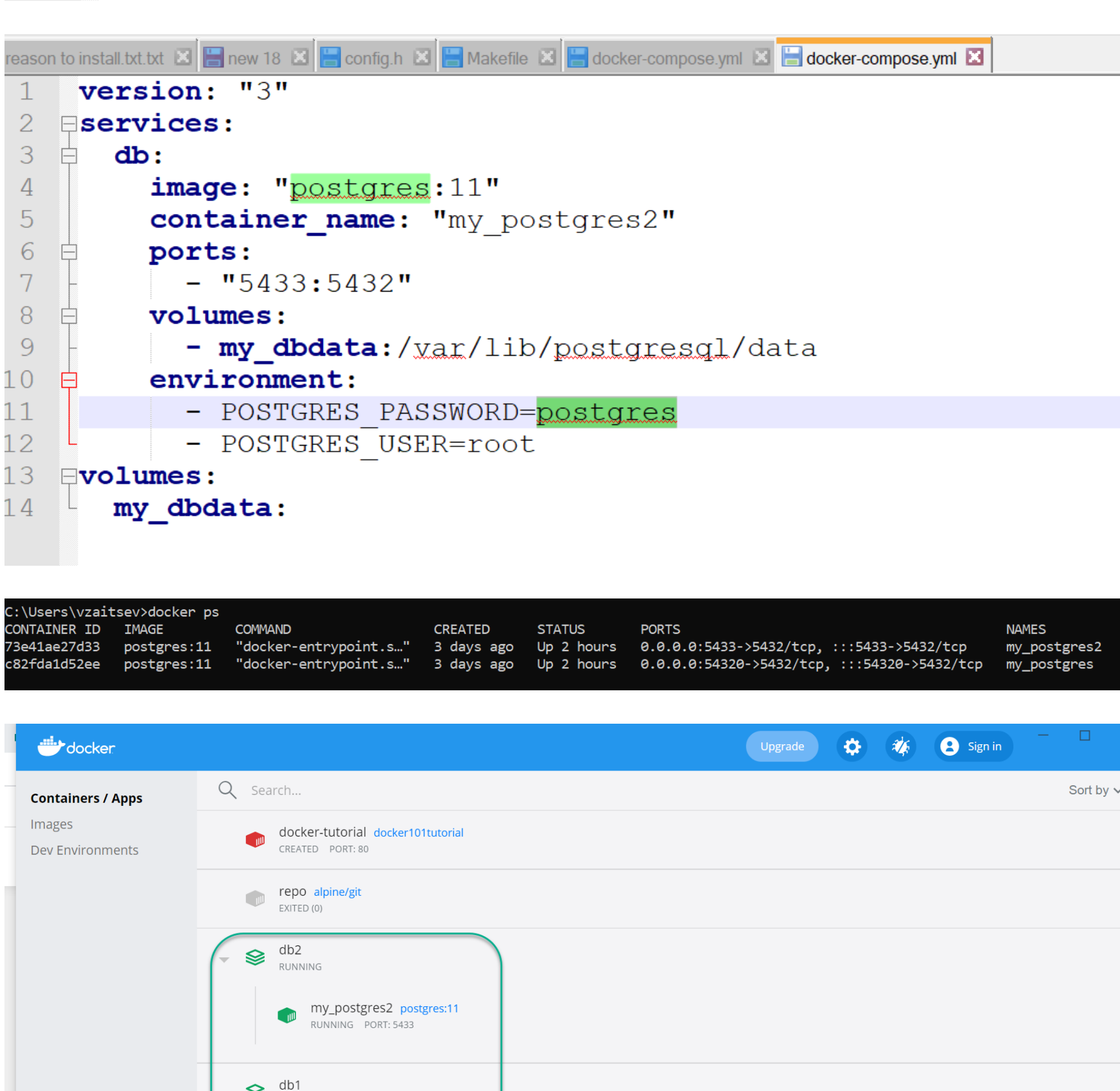
- 1) Развернуть всю архитектуру у себя
- 2) Написать ETL процесс для загрузки ВСЕХ таблиц из postgres-источника в postgres-приемник

Выполнение домашнего задания

0. Подготовительная часть

Подготовка докер образов

docker ps



Создание баз для работы

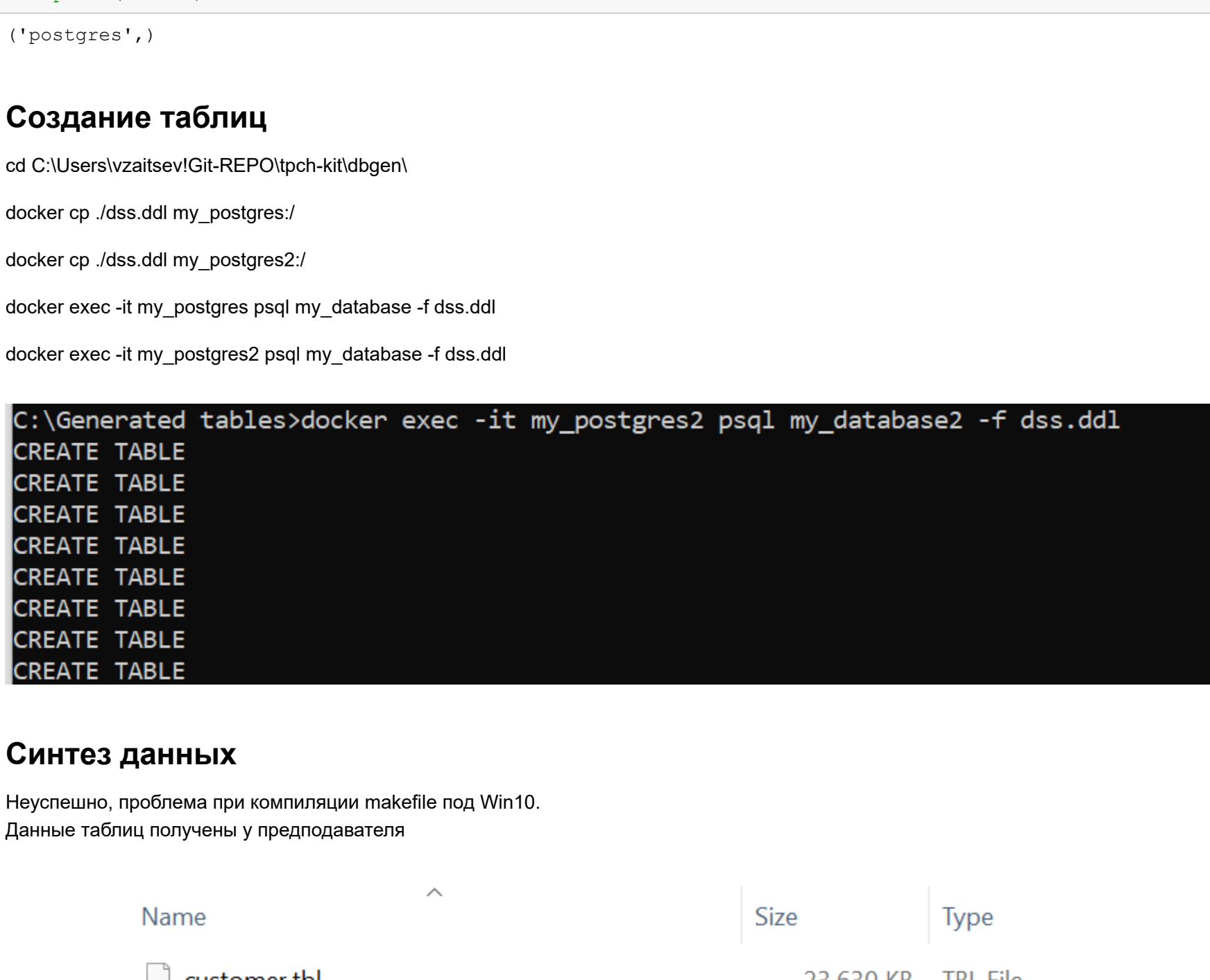
docker exec -it my\_postgres psql -U root -c "create database my\_database"

Done

docker exec -it my\_postgres2 psql -U root -c "create database my\_database2"

Done

Установка необходимых python пакетов



Проверка соединений к созданным базам

'my\_database' (port 54320)

```
In [6]: import psycopg2

# Параметры соединения
conn_string = "host='localhost' port=54320 dbname='my_database' user='root' password='postgres'"
# dbname='tpch'
# conn_string = "host='localhost' port=54320 dbname='tpch' user='root' password='postgres'"

# Создаем соединение (оно поддерживает контекстный менеджер, рекомендую пользоваться им)
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    # query = "select * from customer limit 1" # запрос к БД
    query = "SELECT database FROM pg_database;" # запрос к БД
    cursor.execute(query) # выполнение запроса
    result = cursor.fetchone() # получение результата
    print(result)

('postgres',)
```

'my\_database2' (port 5433)

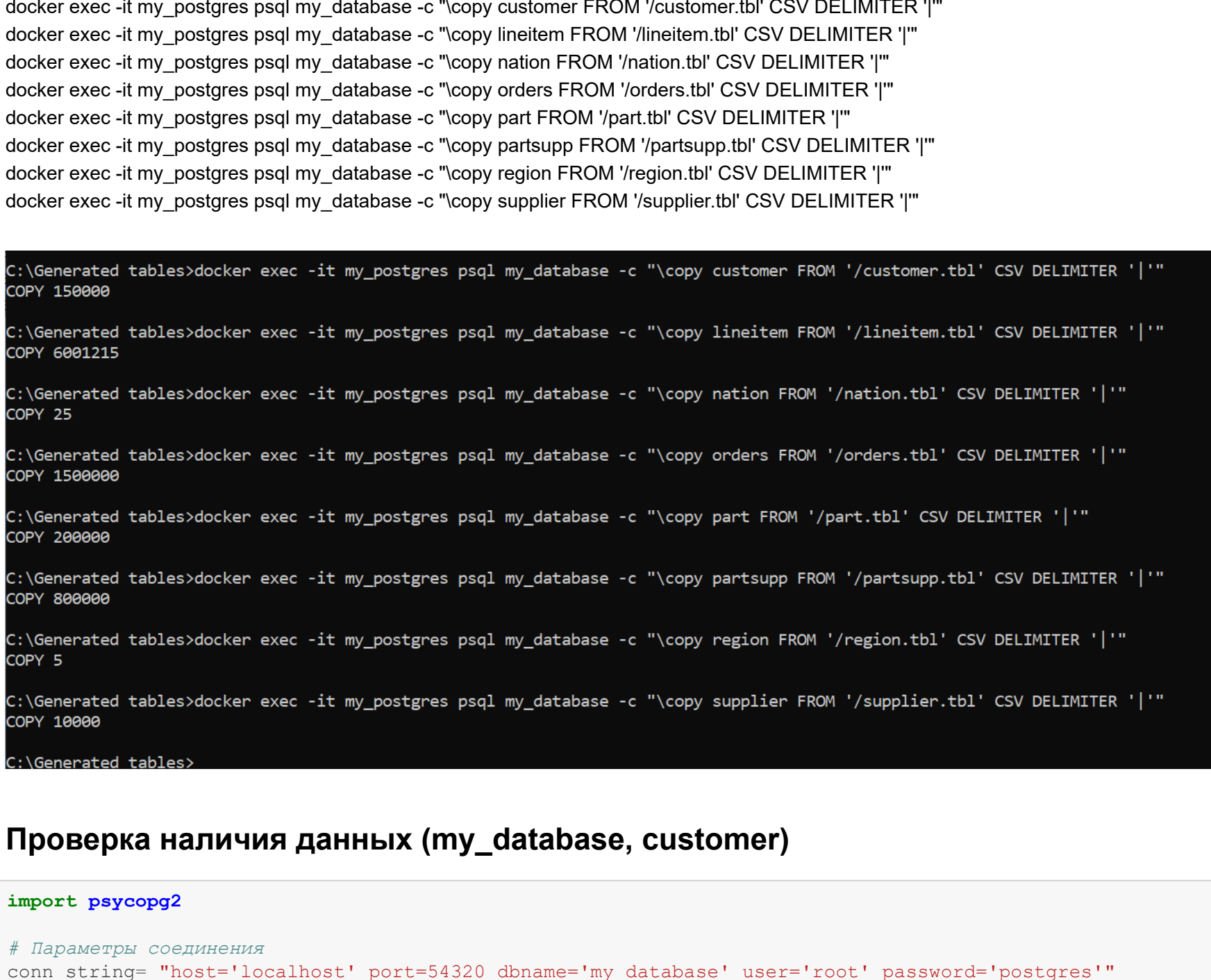
```
In [5]: import psycopg2

# Параметры соединения
conn_string = "host='localhost' port=5433 dbname='my_database2' user='root' password='postgres'"
# dbname='tpch'
# conn_string = "host='localhost' port=5433 dbname='tpch' user='root' password='postgres'"

# Создаем соединение (оно поддерживает контекстный менеджер, рекомендую пользоваться им)
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    # query = "select * from customer limit 1" # запрос к БД
    query = "SELECT database FROM pg_database;" # запрос к БД
    cursor.execute(query) # выполнение запроса
    result = cursor.fetchone() # получение результата
    print(result)

('postgres',)
```

Создание таблиц



Синтез данных

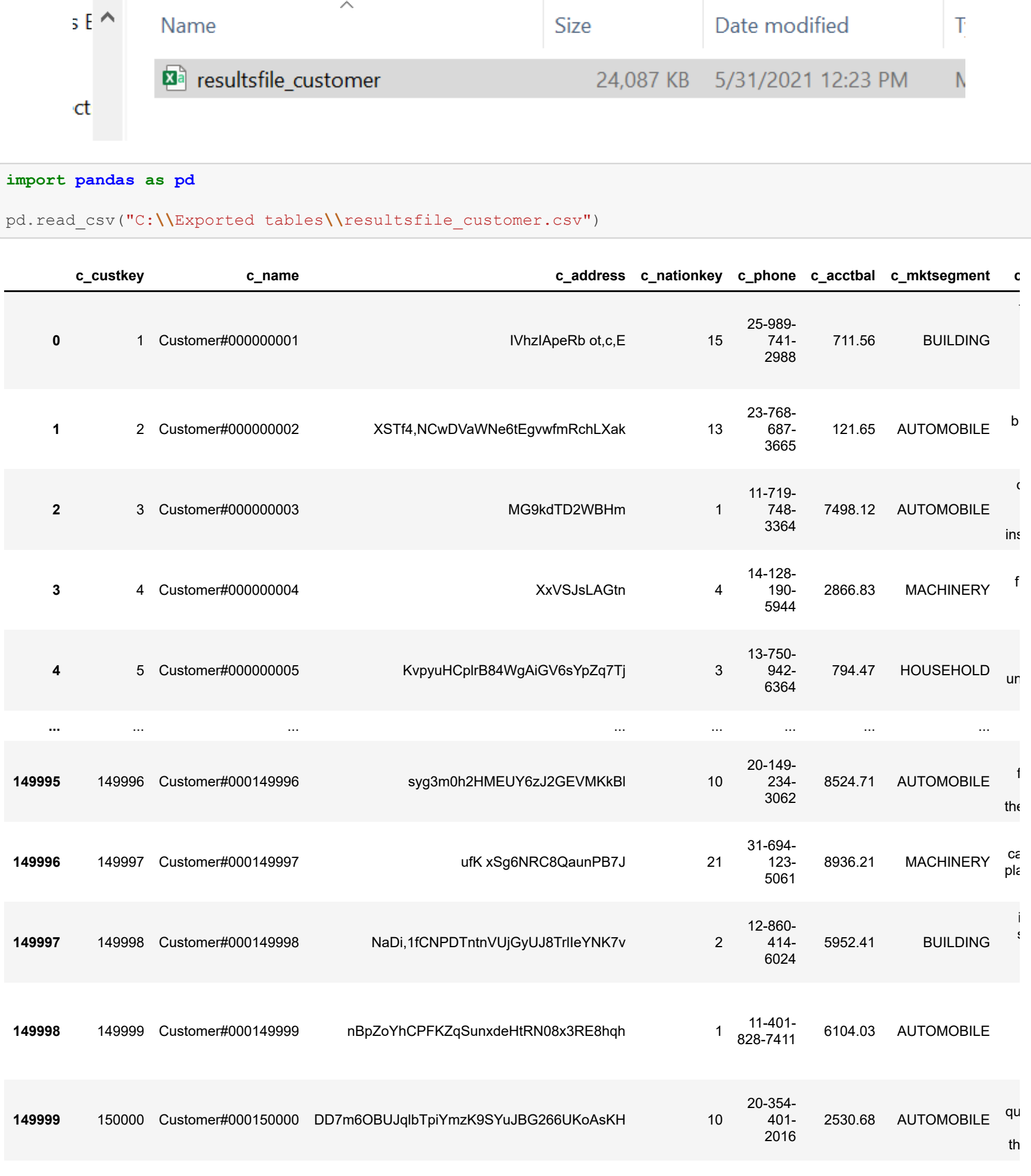
Неуспешно, проблема при компиляции makefile под Win10.

Данные таблиц получены у преподавателя

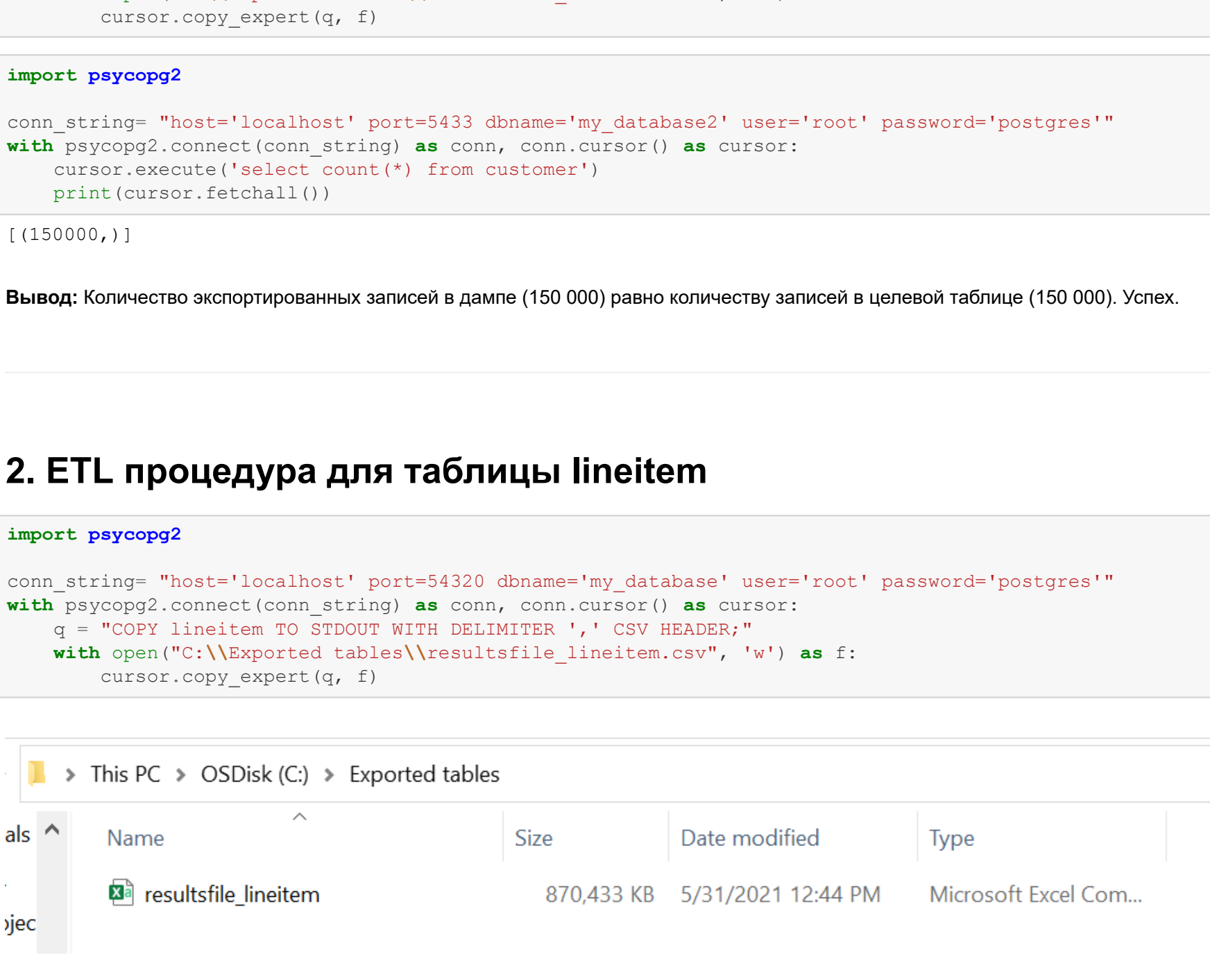
Name	Size	Type
customer.tbl	23,630 KB	TBL File
lineitem.tbl	736,194 KB	TBL File
nation.tbl	3 KB	TBL File
orders.tbl	166,458 KB	TBL File
part.tbl	23,375 KB	TBL File
partsupp.tbl	115,415 KB	TBL File
region.tbl	1 KB	TBL File
supplier.tbl	1,367 KB	TBL File

my\_database

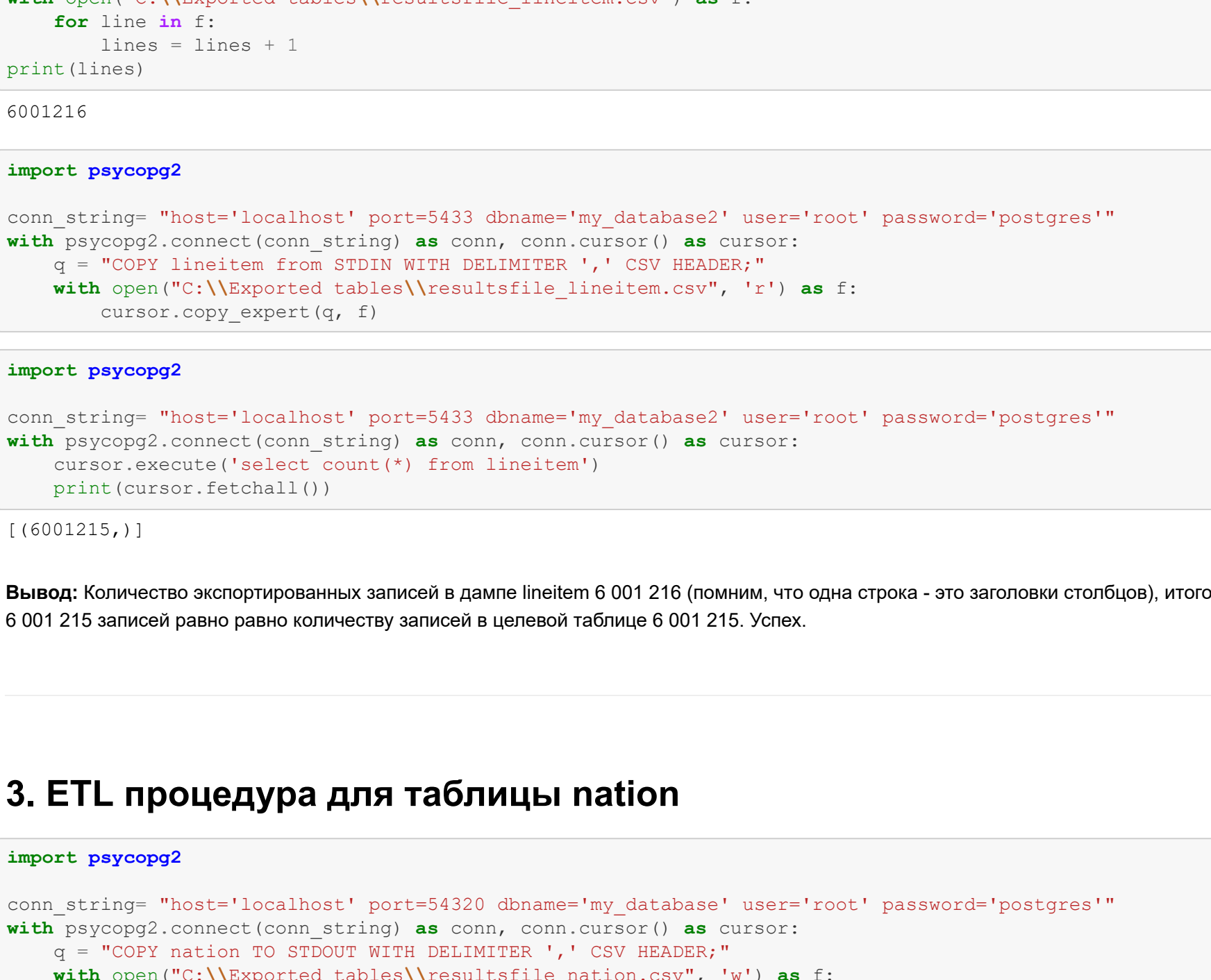
Копирование данных в контейнер



Загрузка данных в базу



Проверка наличия данных (my\_database, customer)



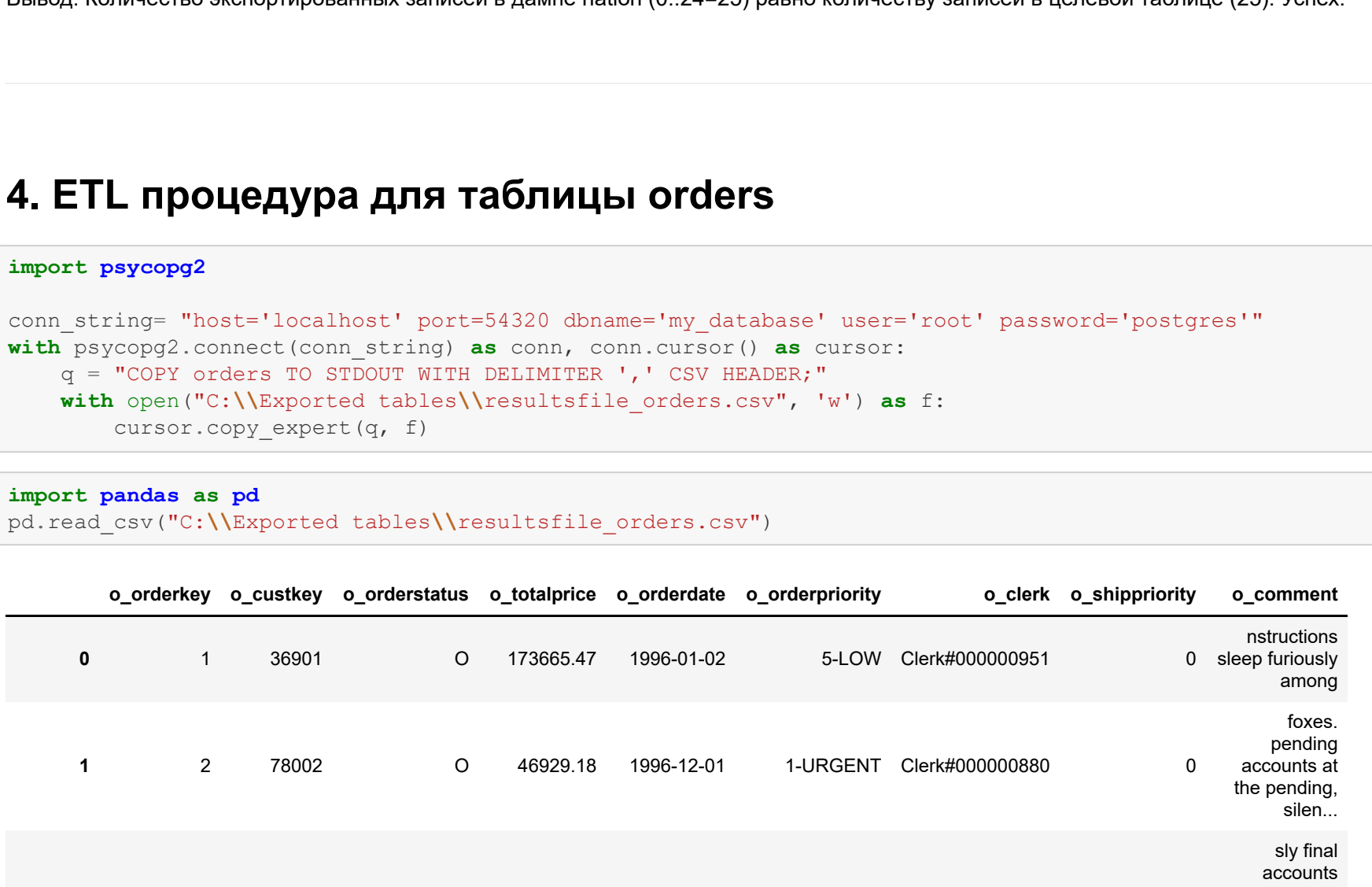
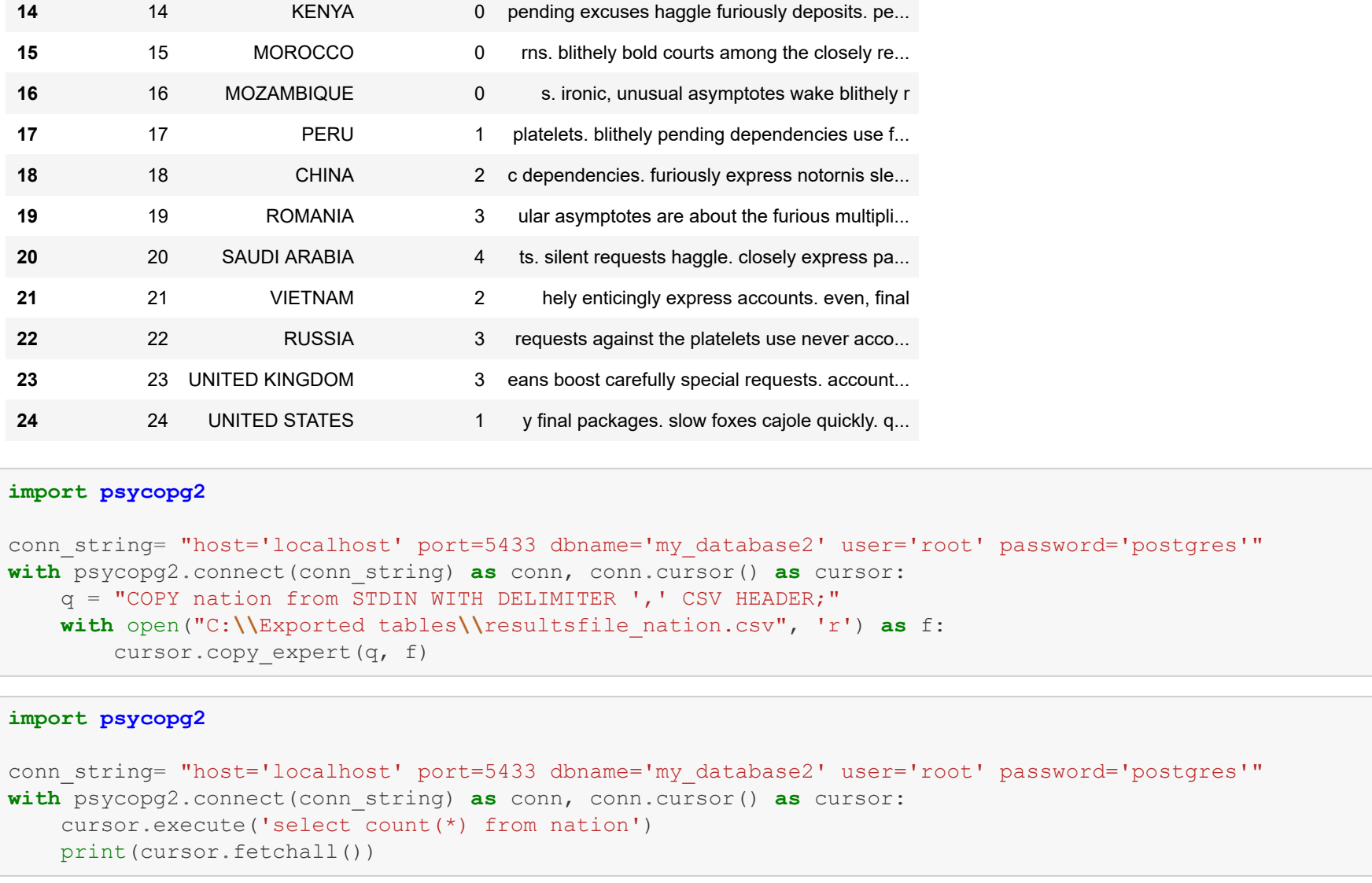
1. ETL процедура для таблицы customer

Для каждой таблицы:

- дампа в файл из исходной таблицы
- проверка выгруженного дампа в pandas
- загрузка из дампа в целевую таблицу
- select count(\*) по целевой таблице

Перечень таблиц:

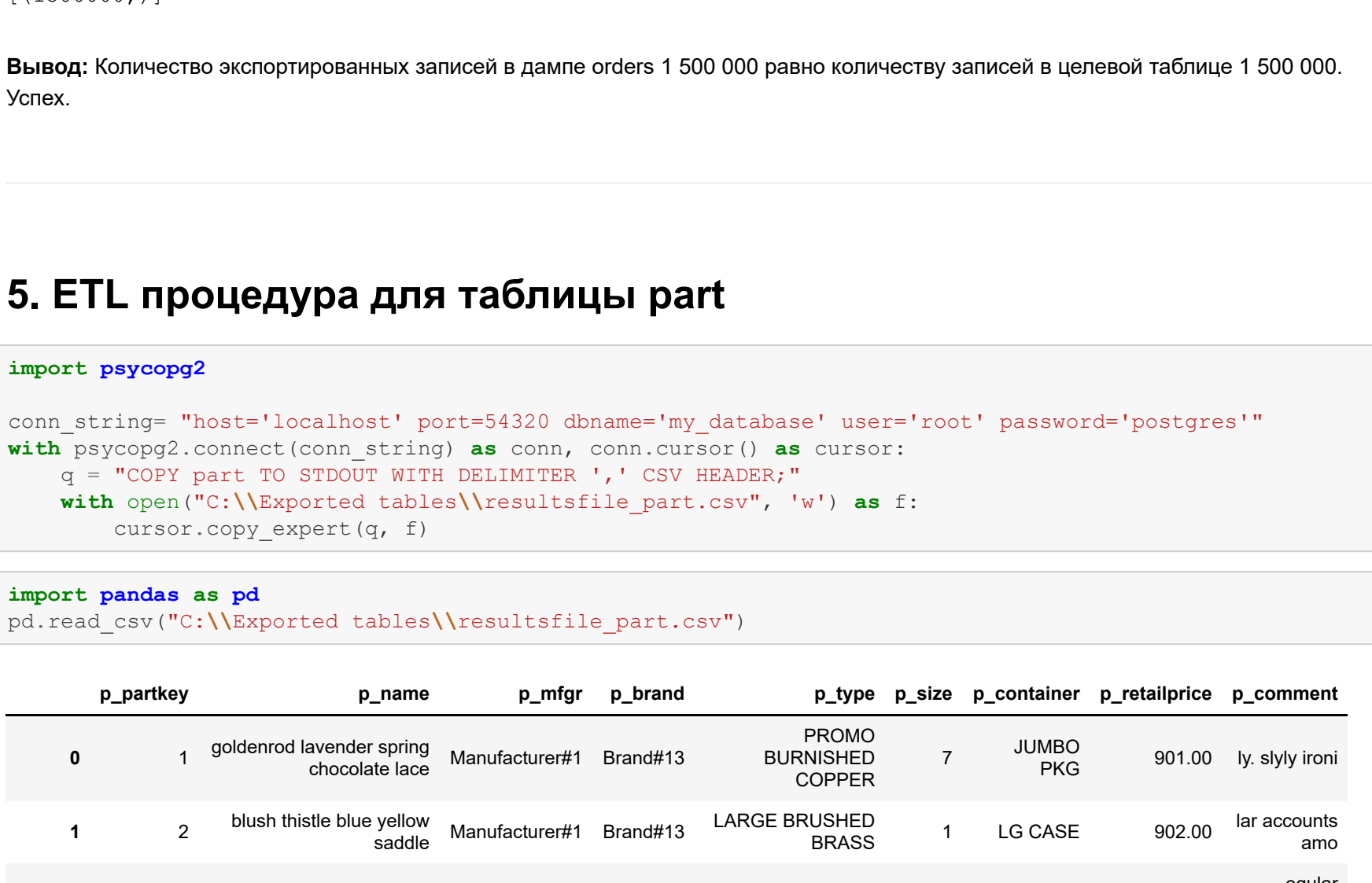
- 1. customer
- 2. lineitem
- 3. nation
- 4. orders
- 5. part
- 6. partsupp
- 7. region
- 8. supplier



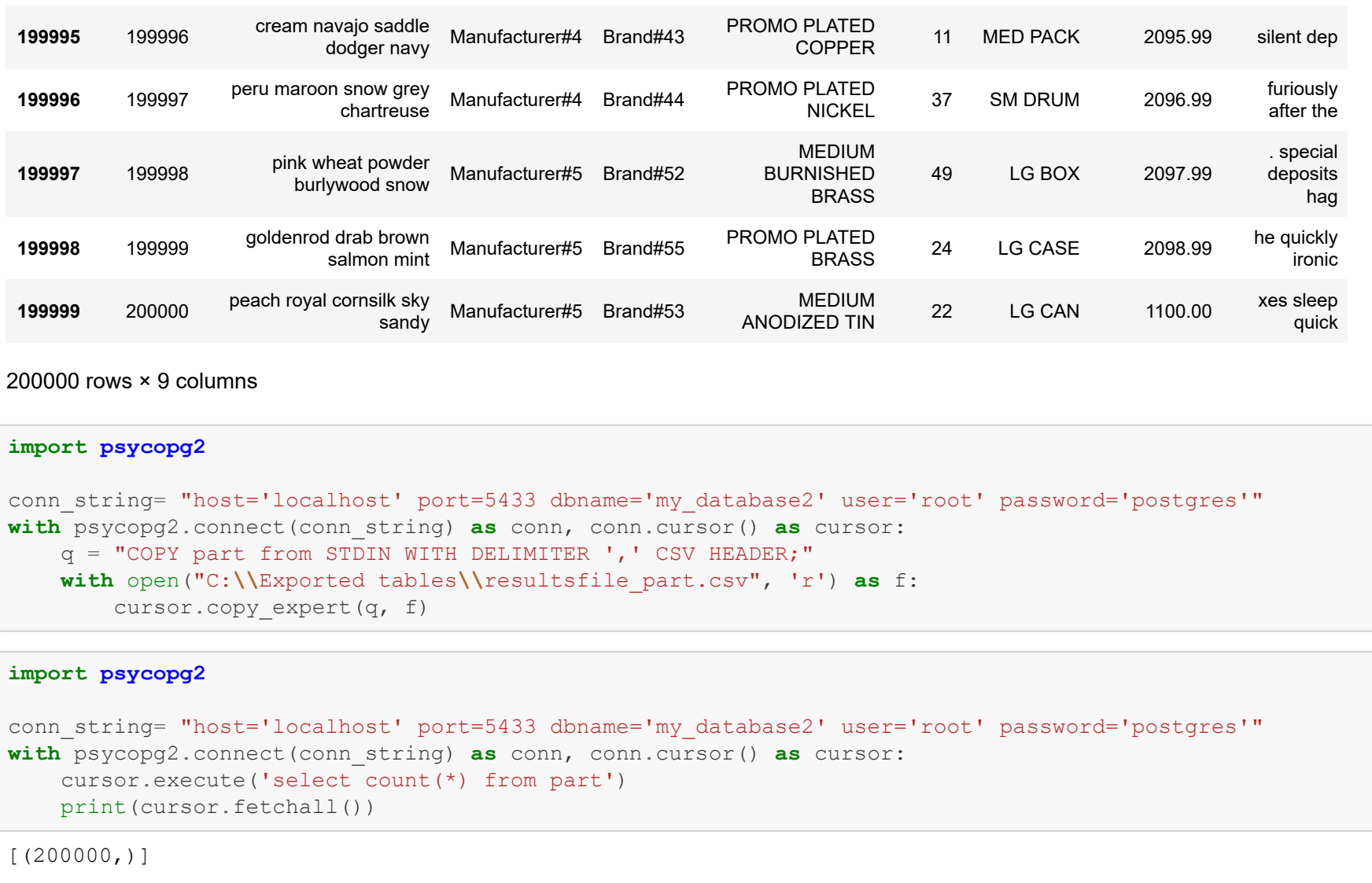
2. ETL процедура для таблицы lineitem



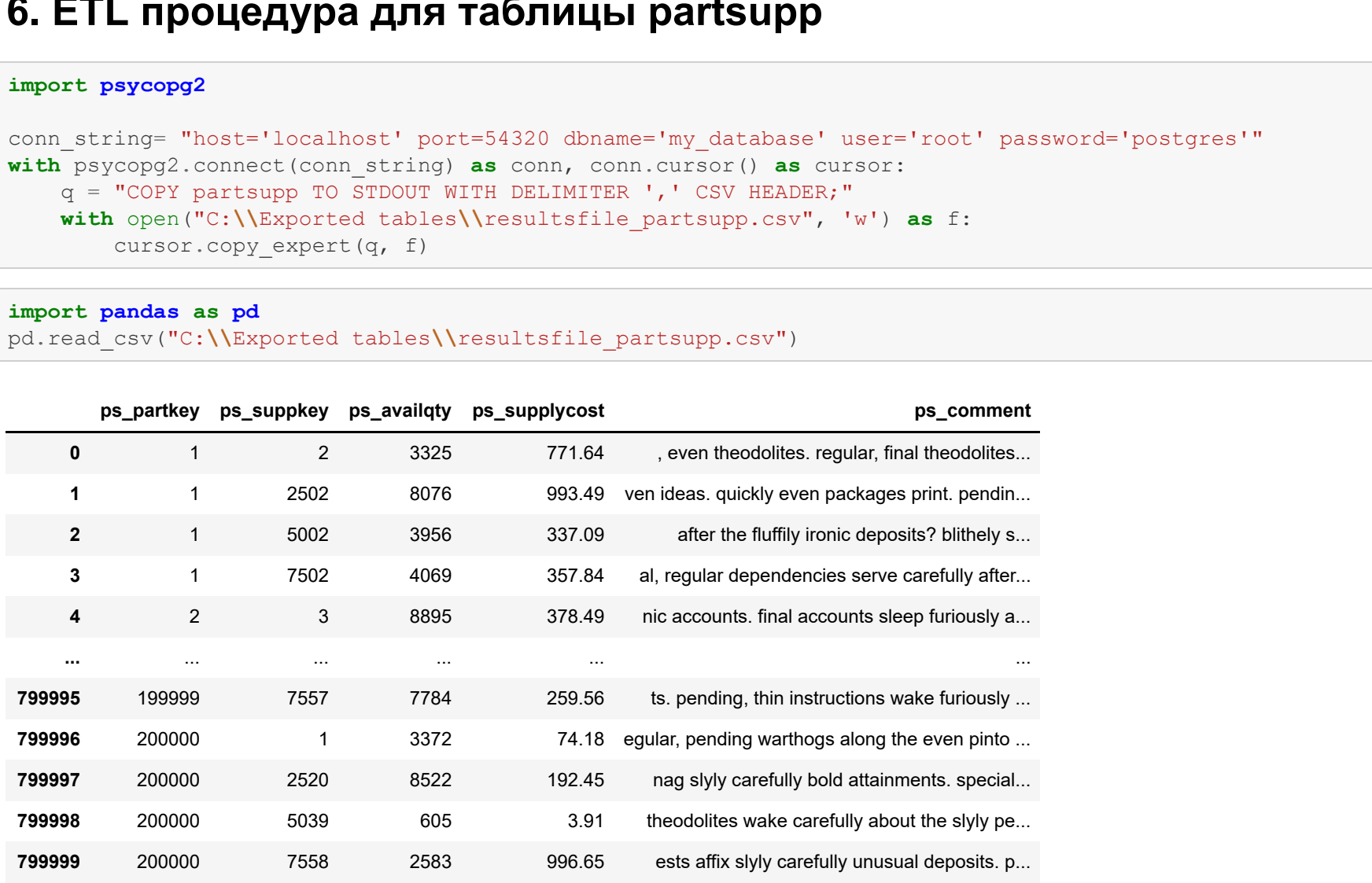
3. ETL процедура для таблицы nation



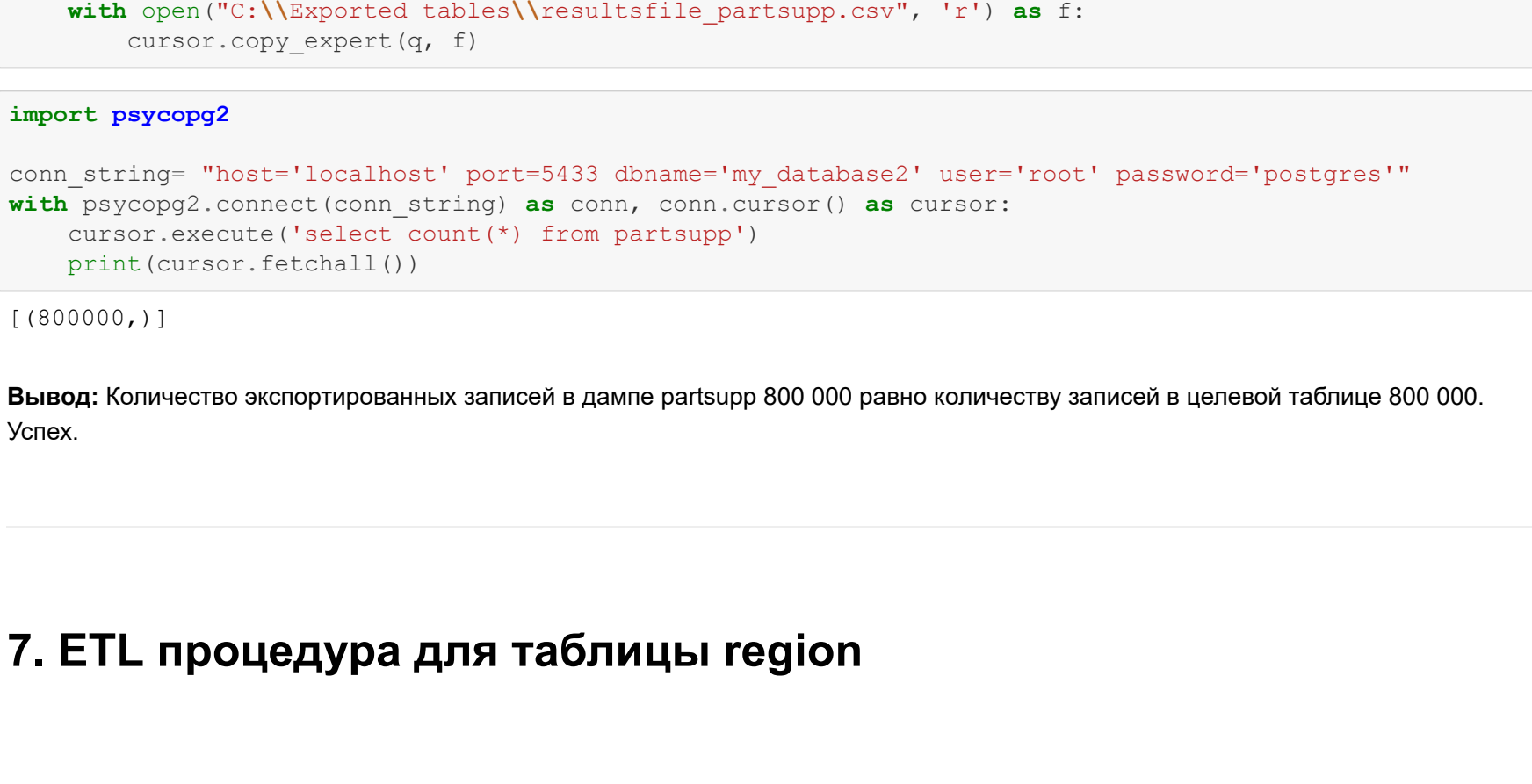
4. ETL процедура для таблицы orders



5. ETL процедура для таблицы part



6. ETL процедура для таблицы partsupp



7. ETL процедура для таблицы region





```
In [20]: import psycopg2

conn_string= "host='localhost' port=54320 dbname='my_database' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    q = "COPY region TO STDOUT WITH DELIMITER ',' CSV HEADER;"
    with open("C:\\Exported tables\\resultsfile_region.csv", 'w') as f:
        cursor.copy_expert(q, f)
```

```
In [21]: import pandas as pd

pd.read_csv("C:\\Exported tables\\resultsfile_region.csv")
```

Out [21]:

	r_regionkey	r_name	r_comment
0	0	AFRICA	lar deposits. blithely final packages cajole...
1	1	AMERICA	hs use ironic, even requests. s
2	2	ASIA	ges. thinly even pinto beans ca
3	3	EUROPE	ly final courts cajole furiously final excuse
4	4	MIDDLE EAST	uckly special accounts cajole carefully blith...

```
In [22]: import psycopg2

conn_string= "host='localhost' port=54330 dbname='my_database2' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    q = "COPY region TO STDOUT WITH DELIMITER ',' CSV HEADER;"
    with open("C:\\Exported tables\\resultsfile_region.csv", 'w') as f:
        cursor.copy_expert(q, f)
```

```
In [23]: import psycopg2

conn_string= "host='localhost' port=5433 dbname='my_database2' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    cursor.execute('select count(*) from region')
    print(cursor.fetchall())

[(5,)]
```

**Вывод:** Количество экспортированных записей в дампе region 5 равно количеству записей в целевой таблице 5. Успех.

## 8. ETL процедура для таблицы supplier

```
In [24]: import psycopg2

conn_string= "host='localhost' port=54320 dbname='my_database' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    q = "COPY supplier TO STDOUT WITH DELIMITER ',' CSV HEADER;"
    with open("C:\\Exported tables\\resultsfile_supplier.csv", 'w') as f:
        cursor.copy_expert(q, f)
```

```
In [25]: import pandas as pd

pd.read_csv("C:\\Exported tables\\resultsfile_supplier.csv")
```

Out [25]:

	s_supplierkey	s_name	s_address	s_nationkey	s_phone	s_acctbal	s_comment
0	1	Supplier#000000001	N KD4onr9OM lpw3.gfuBoQDd7IgrzndZ	17	27-818-335-1736	5755.94	each shyly above the careful
1	2	Supplier#000000002	89eJ5ksX3ImuJQbvOvCvC	5	15-479-861-2259	4032.68	stilyly bold instructions. idle dependen
2	3	Supplier#000000003	ct.G3Pj9CjUuYUlcH18BFTKPSaU9SEV3	1	11-383-516-1199	4192.40	blithely silent requests after the express dep...
3	4	Supplier#000000004	Bk7ahwCKBSYQTerEmvMkkgMwg	15	25-843-787-7479	4641.08	iously even requests above the exp
4	5	Supplier#000000005	Gcdm2uJRz5qTVzr	11	21-151-690-3663	-263.84	shly regular pinto bea
...	...	...	...	...	...	...	...
9995	9996	Supplier#000009996	a4eQd7Sz2NRrcCwyM5ley	10	20-998-443-4436	6209.67	s above the blithely even deposits play carefu...
9996	9997	Supplier#000009997	VcQgaTCWQYMS	15	25-177-344-7203	7011.83	ve the furiously ironic plateaus, evenly
9997	9998	Supplier#000009998	lRTcQwCJzbx7GAJcLajdt.8K	1	11-122-533-7674	2801.35	e regular excuses, blithely final pinto beans ...
9998	9999	Supplier#000009999	mX37oAqzqBPfNfLlWdz p	9	19-773-990-9603	231.69	ounts cajole fluffily among the quickly ironic...
9999	10000	Supplier#000010000	aTGLEwaCL4F PDBdv665XBJhPyCOBO	19	26-576-432-2146	8968.42	ly regular boxes boost shyly, quickly special ...

10000 rows x 7 columns

```
In [26]: import psycopg2

conn_string= "host='localhost' port=5433 dbname='my_database2' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    q = "COPY supplier from STDIN WITH DELIMITER ',' CSV HEADER;"
    with open("C:\\Exported tables\\resultsfile_supplier.csv", 'r') as f:
        cursor.copy_expert(q, f)
```

```
In [27]: import psycopg2

conn_string= "host='localhost' port=5433 dbname='my_database2' user='root' password='postgres'"
with psycopg2.connect(conn_string) as conn, conn.cursor() as cursor:
    cursor.execute('select count(*) from supplier')
    print(cursor.fetchall())

[(10000,)]
```

**Вывод:** Количество экспортированных записей в дампе supplier 10 000 равно количеству записей в целевой таблице 10 000. Успех.

In [ ] :