

Команда № 25.

Участники команды:

Деданов Алексей

Конева Анастасия

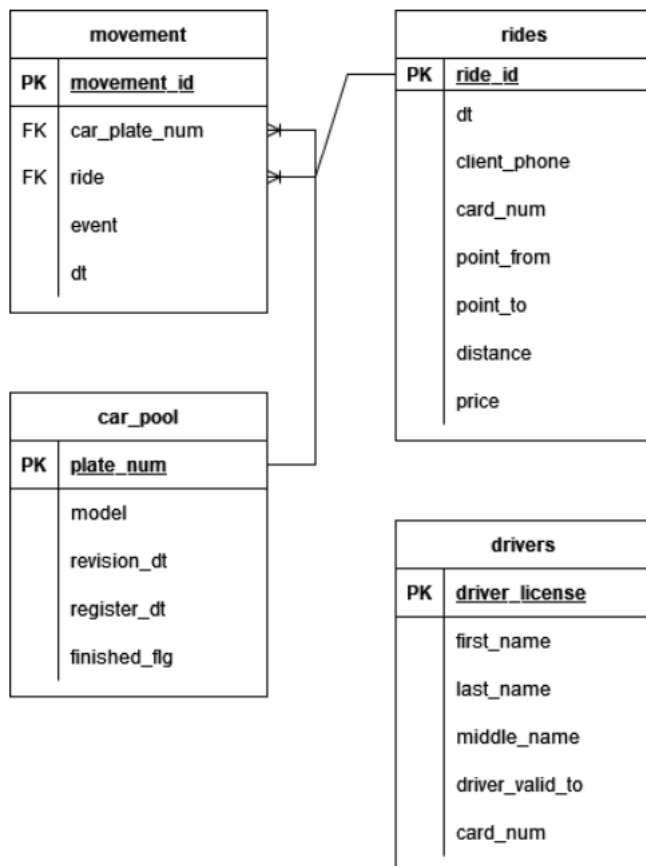
Конева Анна

Есаков Вячеслав

Техническое задание на разработку хранилища данных. Этап 1.

1. Настроить ETL процесс, который будет забирать данные из источника и раскладывать их по целевым таблицам в описанную в документации структуру в хранилище

Первым шагом для извлечения данных с сервера необходимо создать базу данных на сервере PostgreSQL (таблицы work_movement, work_rides, work_car_pool, work_drivers) при помощи базовых команд SQL (используя данные из таблиц сущностей)



С помощью этих таблиц будет извлекаться информация с исходного сервера в нашу базу данных (через команды питона)

После создания таблиц подключимся к серверу по предоставленным нам реквизитам

Реквизиты подключения к источнику техническим пользователем:

- host: de-edu-db.chronosavant.ru
- port: 5432
- database: taxi
- schema: main
- user: etl_tech_user
- password: etl_tech_user_password

Для того, чтобы подключиться к серверу с данными необходимо подключить и импортировать библиотеку для работы с *psycopg2*, которая предназначена для адаптеров баз данных (субд PostgreSQL)

```
1 #Сначала подключаемся к их бд и извлекаем нужные данные
2 conn = psycopg2.connect(dbname='taxi', user='etl_tech_user',
3                           password='etl_tech_user_password', host='de-edu-db.chronosavant.ru')
4 cursor = conn.cursor()
5
```

Вводим необходимые реквизиты и извлекаем данные, которые постоянно обновляются (добавляются)

Поскольку изначально (после создания пустых таблиц) в нашей таблице отсутствуют данные используем счетчик для каждой из таблиц (чтобы подгружать только что добавившиеся данные)

```
In [21]: 1 cur_count_movement = 0 # Изначально кол-во данных в нашей бд равно нулю
2 cur_count_rides = 0
3 cur_count_car_pool = 0
4 cur_count_drivers = 0
```

Для каждой из таблиц извлекаем данные с сервера, чтобы потом добавить их на другой сервер

Демонстрация извлечения данных для каждой из 4ех таблиц:

```

6 # SQL запрос извлекающий новые данные
7 query = f'''
8     select * from main.movement
9     limit 100000 offset {cur_count_movement};
10 '''
11 cursor.execute(query)
12 data_movement = list(cursor.fetchall()) # Формируем список новых только что полученных данных
13
14 query = f'''
15     select * from main.rides
16     limit 100000 offset {cur_count_rides};
17 '''
18 cursor.execute(query)
19 data_rides = list(cursor.fetchall())
20
21 query = f'''
22     select * from main.car_pool
23     limit 100000 offset {cur_count_car_pool};
24 '''
25 cursor.execute(query)
26 data_car_pool = list(cursor.fetchall())
27
28 query = f'''
29     select * from main.drivers
30     limit 100000 offset {cur_count_drivers};
31 '''
32 cursor.execute(query)
33
34 data_drivers = list(cursor.fetchall())
35

```

После того, как данные извлеклись, их необходимо добавить на другой сервер и загрузить в пустые таблицы (при этом новые данные также должны добавляться)

Для того, чтобы загрузить извлечённые данные на другой сервер, нужно подключиться к этому серверу (аналогично с извлечением данных)

```

2 conn = psycopg2.connect(dbname='taxi', user='etl_tech_user',
3                          password='etl_tech_user_password', host='de-edu-db.chronosavant.ru')
4 cursor = conn.cursor()
5

```

Далее идет процесс загрузки данных:

```

42 # movement
43
44 query = 'INSERT INTO work_movement (movement_id, car_plate_num, ride, event, dt) VALUES (%s, %s, %s, %s, %s)'
45 for val in data_movement:
46     cursor.execute(query, val)
47 conn.commit() #фиксируем изменения в бд
48
49 cursor.execute("select count(movement_id) from work_movement;")
50
51 cur_count_movement = cursor.fetchall()[0][0] #обновляем количество данных в нашей бд.
52 print(cur_count_movement)
53 data_movement[:5]
54
55 # rides
56
57 query = 'INSERT INTO work_rides (ride_id, dt, client_phone, card_num, point_from, point_to, distance, price) VALUES (%s, %s, %s, %s, %s, %s, %s, %s)'
58 for val in data_rides:
59     cursor.execute(query, val)
60 conn.commit() #фиксируем изменения в бд
61
62 cursor.execute("select count(ride_id) from work_rides;")
63
64 cur_count_rides = cursor.fetchall()[0][0] #обновляем количество данных в нашей бд.
65 print(cur_count_rides)
66 data_rides[:5]
67
68
69
70 ## drivers
71
72 query = 'INSERT INTO work_drivers (driver_license, first_name, last_name, middle_name, driver_valid_to, card_num, update_dt)'
73 for val in data_drivers:
74     cursor.execute(query, val[:7])
75 conn.commit() #фиксируем изменения в бд
76
77 cursor.execute("select count(driver_license) from work_drivers;")
78
79 cur_count_rides = cursor.fetchall()[0][0] #обновляем количество данных в нашей бд.
80 print(cur_count_drivers)
81 data_drivers[:5]
82
83
84
85 ## car_pool
86
87 query = 'INSERT INTO work_car_pool (plate_num, model, revision_dt, register_dt, finished_flg) VALUES (%s, %s, %s, %s, %s)'
88 for val in data_car_pool:
89     cursor.execute(query, val[:5])
90 conn.commit() #фиксируем изменения в бд
91
92 cursor.execute("select count(plate_num) from work_car_pool;")
93
94 cur_count_rides = cursor.fetchall()[0][0] #обновляем количество данных в нашей бд.
95 print(cur_count_car_pool)
96 data_car_pool[:5]
97

```

6

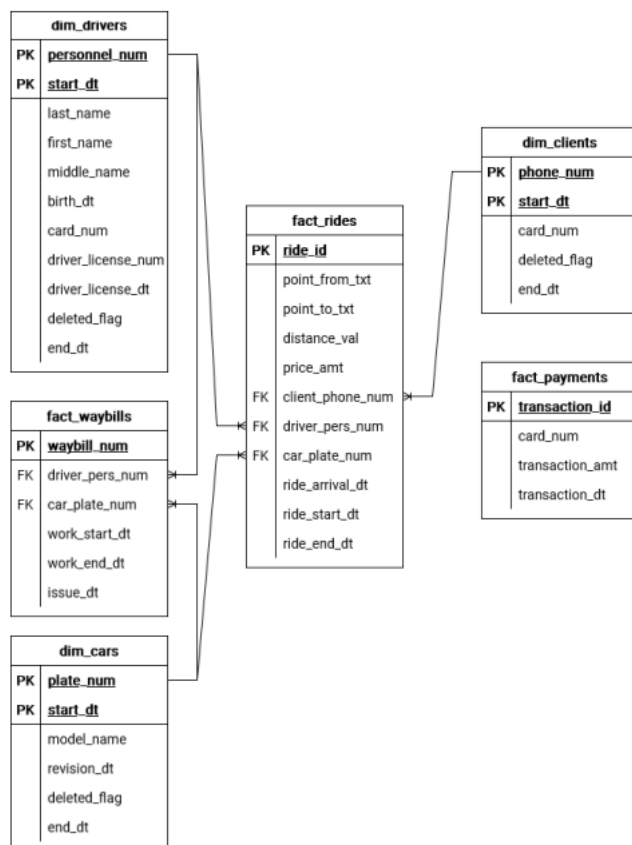
Демонстрация извлечения и загрузки данных на сервер (данные извлеклись одного сервера и загрузились на другой). Аналогично для остальных таблиц.

driver_license_dt	first_name	last_name	middle_name	driver_valid_to	card_num	update
13 51 130647 2:47:03+00	рэшһыһ	һхыһхыһ	һхыһхыһхыһ	2023-04-12	4205 4648 9442 7155	2022-10-12 1
38 10 977977 2:50:03+00	һюһырт	һюһыртһюһ	һюһыртһюһ	2023-04-12	7719 1583 7663 4759	2022-10-12 1
34 48 464847 2:52:03+00	һырэ	һрыһыһ	һыһыһыһ	2023-04-12	1875 8923 3744 8928	2022-10-12 1
42 30 001788 2:54:03+00	һюһыһ	һрыһыһ	һыһыһыһ	2023-04-12	2186 7201 8599 1460	2022-10-12 1
31 39 276519 2:54:03+00	һюһыһ	һрыһыһ	һыһыһыһ	2023-04-12	9590 3819 7544 6557	2022-10-12 1
91 95 521602 2:55:03+00	һюһыһ	һрыһыһ	һыһыһыһ	2023-04-12	5198 9779 9847 9503	2022-10-12 1
68 87 611785 2:55:03+00	һюһыһ	һрыһыһ	һыһыһыһ	2023-04-12	1486 9214 3856 5510	2022-10-12 1
99 00 949942 2:57:03+00	һрыһыһ	һюһыһыһ	һхыһыһыһ	2023-04-12	6222 8600 4484 3992	2022-10-12 1
22 79 062689 2:57:03+00	һыһыһыһ	һхыһыһыһ	һыһыһыһ	2023-04-12	5247 4495 9014 2398	2022-10-12 1
96 34 555121 2:58:03+00	һшүһыһ	һрыһыһ	һшүһыһыһ	2023-04-12	3503 3829 4377 8481	2022-10-12 1
39 24 868812 3:02:03+00	һртхыһыһ	һрыһыһыһ	һшүһыһыһ	2023-04-12	4059 9235 5280 7658	2022-10-12 1
57 45 478087 3:05:03+00	һсесрьһ	һрыһыһыһ	һшүһыһыһ	2023-04-12	4478 1758 5828 2283	2022-10-12 1
60 53 123506 -- Далее --	һшүһыһыһ	һрыһыһыһ	һшүһыһыһ	2023-04-12	5136 9184 5614 7650	2022-10-12 1

Работоспособность кода проверяли на локальном сервере PjsgreSQL, после чего подключились к нужному серверу.

Данные извлеклись и загрузились.

Следующим шагом мы создаем в субд PosgreSQL данные таблицы (с помощью базовых команд языка SQL).



С помощью библиотеки *ftp lib* мы можем подключаться к серверу для выгрузки путевых листов, требуется подключиться к *ftp* серверу. Подключиться к серверу нам удалось, но извлечь данные не получилось.

```

: 1 #from ftplib import FTP_TLS
  2 import ftplib
  3
  4 import ssl
  5 import sys, os
  
```



```

1 host = 'de-edu-db.chronosavant.ru'
2 ftp_user = 'etl_tech_user'
3 ftp_password = 'etl_tech_user_password'
4 port = 21
5 ftps = ftplib.FTP_TLS(host)
6
7 ftps.auth()
8 ftps.prot_p()
9 ftps.login(ftp_user, ftp_password)
10
11 ftps.cwd('/waybills')
12 ftps.cwd('/payments')
13
14 print(ftps.pwd())
15 #ftps.retrlines('LIST')
16 #ftps.nlst()
17 print(ftps.sendcmd('PASV'))
18 #ftps.set_pasv(21)
19 print(ftps.sendcmd('pwd'))
20
21 ftps.quit()
22 ftps.close()
23

```

/payments

227 Entering Passive Mode (10,128,0,31,171,46).

257 "/payments" is the current directory

ИТОГ:

Была проделана работа по извлечению и переносу данных с одного сервера на другой в необходимую базу данных

Данные были извлечены и загружены. Реализовано динамическое обновление данных.