

بسمه تعالی

پایتون و استخراج، تبدیل و بارگذاری داده‌ها

اغلب کاربران نیاز به «استخراج» (Extract)، «تبدیل» (Transform) و «بارگذاری» (Load) داده‌ها از «پایگاه داده‌های» (databases) گوناگون به «انبارهای داده» (data warehouse) به منظور تجمع داده‌ها جهت تحلیل‌های «هوش تجاری» (BI | Business Intelligence) و «تحلیل‌های داده» (Data Analysis) دارند. نرم‌افزارهای گوناگونی برای انجام فرآیند ETL تولید شده‌اند، ولی استفاده از آن‌ها برای کاربردهای ساده اغلب زیاده‌روی به نظر می‌رسد. در این مطلب چگونگی استخراج داده‌ها از MySQL، SQL-Server و Firebird، تبدیل داده‌ها و بارگذاری آن‌ها در SQL-Server (به عنوان انبار داده) با استفاده از پایتون ۳.۶ آموزش داده شده است. اکنون، کار با راه‌اندازی دایرکتوری پروژه آغاز می‌شود:

```
1 python_etl
2 |__main__.py
3 |__db_credentials.py
4 |__variables.py
5 |__sql_queries.py
6 |__etl.py
```

علاوه بر نصب MySQL، SQL-Server و Firebird، نیاز به سه ماژول پایتونی که در زیر بیان شده و اجرای pip install [module_name] وجود دارد.

- mysql-connector-python : اتصال به MySQL
- pyodbc : اتصال به sql-server
- fdb : اتصال به Firebird

راه‌اندازی مجوزها و متغیرهای پایگاه داده

راه‌اندازی یک «متغیر» برای ذخیره‌سازی نام پایگاه داده انبار داده در variables.py با دستور زیر انجام می‌شود.

```
datawarehouse_name = 'your_datawarehouse_name'
```

نصب همه پایگاه داده‌های منبع و هدف‌گیری رشته‌های اتصال پایگاه داده و مجوزها در db_credentials.py به صورتی که در زیر نمایش داده شده قابل انجام است. سپس، کانفیگ باید به صورت list ذخیره شود تا بتوان بعداً آن را در پایگاه داده‌های زیادی تکرار کرد.

```

1 from variables import datawarehouse_name
2
3 # sql-server (target db, datawarehouse)
4 datawarehouse_db_config = {
5     'Trusted_Connection': 'yes',
6     'driver': '{SQL Server}',
7     'server': 'datawarehouse_sql_server',
8     'database': '{}'.format(datawarehouse_name)
9     'user': 'your_db_username',
10    'password': 'your_db_password',
11    'autocommit': True,
12 }
13
14 # sql-server (source db)
15 sqlserver_db_config = [
16     {
17         'Trusted_Connection': 'yes',
18         'driver': '{SQL Server}',
19         'server': 'your_sql_server',
20         'database': 'db1'
21         'user': 'your_db_username',
22         'password': 'your_db_password',
23         'autocommit': True,
24     }
25 ]
26
27 # mysql
28 mysql_db_config = [
29     {
30         'user': 'your_user_1',
31         'password': 'your_password_1',
32         'host': 'db_connection_string_1',
33         'database': 'db_1',
34     },
35     {
36         'user': 'your_user_2',
37         'password': 'your_password_2',
38         'host': 'db_connection_string_2',
39         'database': 'db_2',
40     },
41 ]
42
43 # firebird
44 fdb_db_config = [
45     {
46         'dsn': "/your/path/to/source.db",
47         'user': "your_username",
48         'password': "your_password",
49     }
50 ]

```

کوئری‌های SQL

sql_queries.py جایی است که همه کوئری‌های SQL برای استخراج از پایگاه‌داده‌های منبع و بارگذاری در پایگاه‌داده هدف (انبار داده) ذخیره می‌شوند. با توجه به اینکه در این فرآیند با پلتفرم‌های گوناگونی کار می‌شود، می‌توان از «نحو» (syntax) گوناگونی برای هر سکوی داده با جداسازی کوئری‌ها مطابق با نوع پایگاه داده استفاده کرد.

```

1 # example queries, will be different across different db platform
2 firebird_extract = '''
3     SELECT fbd_column_1, fbd_column_2, fbd_column_3
4     FROM fbd_table;
5 '''
6
7 firebird_insert = '''
8     INSERT INTO table (column_1, column_2, column_3)
9     VALUES (?, ?, ?)
10 '''
11
12 firebird_extract_2 = '''
13     SELECT fbd_column_1, fbd_column_2, fbd_column_3
14     FROM fbd_table_2;
15 '''
16
17 firebird_insert_2 = '''
18     INSERT INTO table_2 (column_1, column_2, column_3)
19     VALUES (?, ?, ?)
20 '''
21
22 sqlserver_extract = '''
23     SELECT sqlserver_column_1, sqlserver_column_2, sqlserver_column_3
24     FROM sqlserver_table
25 '''
26
27 sqlserver_insert = '''
28     INSERT INTO table (column_1, column_2, column_3)
29     VALUES (?, ?, ?)
30 '''
31
32 mysql_extract = '''
33     SELECT mysql_column_1, mysql_column_2, mysql_column_3
34     FROM mysql_table
35 '''
36
37 mysql_insert = '''
38     INSERT INTO table (column_1, column_2, column_3)
39     VALUES (?, ?, ?)
40 '''
41
42 # exporting queries
43 class SqlQuery:
44     def __init__(self, extract_query, load_query):
45         self.extract_query = extract_query
46         self.load_query = load_query
47
48 # create instances for SqlQuery class
49 fbd_query = SqlQuery(firebird_extract, firebird_insert)
50 fbd_query_2 = SqlQuery(firebird_extract_2, firebird_insert_2)
51 sqlserver_query = SqlQuery(sqlserver_extract, sqlserver_insert)
52 mysql_query = SqlQuery(mysql_extract, mysql_insert)
53
54 # store as list for iteration
55 fbd_queries = [fbd_query, fbd_query_2]
56 sqlserver_queries = [sqlserver_query]
57 mysql_queries = [mysql_query]

```

استخراج، تبدیل، بارگذاری

در etl.py، ماژول‌ها و متغیرهای پایتون زیر باید بارگذاری شوند تا کار آغاز شود.

```

1 # python modules
2 import mysql.connector
3 import pyodbc
4 import fdb
5 # variables
6 from variables import datawarehouse_name

```

در اینجا دو متد `etl()` و `etl_process()` وجود دارند. `etl_process()` متدی برای راه‌اندازی اتصال منبع پایگاه داده مطابق با پلتفرم پایگاه‌داده و فراخوانی متد `etl()` است. در متد `etl()`، ابتدا کوئری استخراج اجرا، داده‌های SQL در متغیر Data ذخیره و این متغیر در پایگاه‌داده هدف که در واقع انبار داده کاربر است «درج» می‌شود. تبدیل داده‌ها با دستکاری متغیر data که از نوع «تاپل» است صورت می‌پذیرد.

```
1 def etl(query, source_cnx, target_cnx):
2     # extract data from source db
3     source_cursor = source_cnx.cursor()
4     source_cursor.execute(query.extract_query)
5     data = source_cursor.fetchall()
6     source_cursor.close()
7
8     # load data into warehouse db
9     if data:
10        target_cursor = target_cnx.cursor()
11        target_cursor.execute("USE {}".format(datawarehouse_name))
12        target_cursor.executemany(query.load_query, data)
13        print('data loaded to warehouse db')
14        target_cursor.close()
15    else:
16        print('data is empty')
17
18 def etl_process(queries, target_cnx, source_db_config, db_platform):
19     # establish source db connection
20     if db_platform == 'mysql':
21         source_cnx = mysql.connector.connect(**source_db_config)
22     elif db_platform == 'sqlserver':
23         source_cnx = pyodbc.connect(**source_db_config)
24     elif db_platform == 'firebird':
25         source_cnx = fdb.connect(**source_db_config)
26     else:
27         return 'Error! unrecognised db platform'
28
29     # loop through sql queries
30     for query in queries:
31         etl(query, source_cnx, target_cnx)
32
33     # close the source db connection
34     source_cnx.close()
```

قرار دادن کلیه موارد در کنار هم

در `main.py` می‌توان از طریق مجوزها و انجام ETL برای همه پایگاه‌های داده حلقه زد. کد زیر همه متغیرها و متدهای مرتبط را ایمپورت می‌کند.

```

1 # variables
2 from db_credentials import datawarehouse_db_config, sqlserver_db_config, mysql_db_config, fbd_db_config
3 from sql_queries import fbd_queries, sqlserver_queries, mysql_queries
4 from variables import *
5 # methods
6 from etl import etl_process

```

```

1 def main():
2     print('starting etl')
3
4     # establish connection for target database (sql-server)
5     target_cnx = pyodbc.connect(*datawarehouse_db_config)
6
7     # loop through credentials
8
9     # mysql
10    for config in mysql_db_config:
11        try:
12            print("loading db: " + config['database'])
13            etl_process(mysql_queries, target_cnx, config, 'mysql')
14        except Exception as error:
15            print("etl for {} has error".format(config['database']))
16            print('error message: {}'.format(error))
17            continue
18
19    # sql-server
20    for config in sqlserver_db_config:
21        try:
22            print("loading db: " + config['database'])
23            etl_process(sqlserver_queries, target_cnx, config, 'sqlserver')
24        except Exception as error:
25            print("etl for {} has error".format(config['database']))
26            print('error message: {}'.format(error))
27            continue
28
29    # firebird
30    for config in fbd_db_config:
31        try:
32            print("loading db: " + config['database'])
33            etl_process(fbd_queries, target_cnx, config, 'firebird')
34        except Exception as error:
35            print("etl for {} has error".format(config['database']))
36            print('error message: {}'.format(error))
37            continue
38
39    target_cnx.close()
40
41 if __name__ == "__main__":
42     main()

```

اکنون کافیسٲ کاربر در ترمینال python main.py را وارد کند و از اسکرپت پایتونی که برای انجام فرآیند ETL نوشته شده استفاده کند.

منبع: <https://codeburst.io/using-python-script-for-data-etl-53138c567906>