

ЗАСТОСУВАННЯ ПЕРЕТВОРЕННЯ SQL

Мета цього експеримента - продемонструвати, як можна виконувати стандартні операції SQL, такі як з'єднання, об'єднання та його узагальнення, в машинному навчанні Microsoft Azure за допомогою модуля «Застосувати перетворення SQL». Використовуючи цей модуль, можна виконувати різні перетворення даних за допомогою SQL.

В цьому експерименті ми використовуємо три набори даних (клієнт ресторану, характеристика ресторану і рейтинги ресторану). Набори даних включають як числові, так і категоріальні характеристики.

На наступній діаграмі зображено загальний робочий процес експерименту (рис. 1):

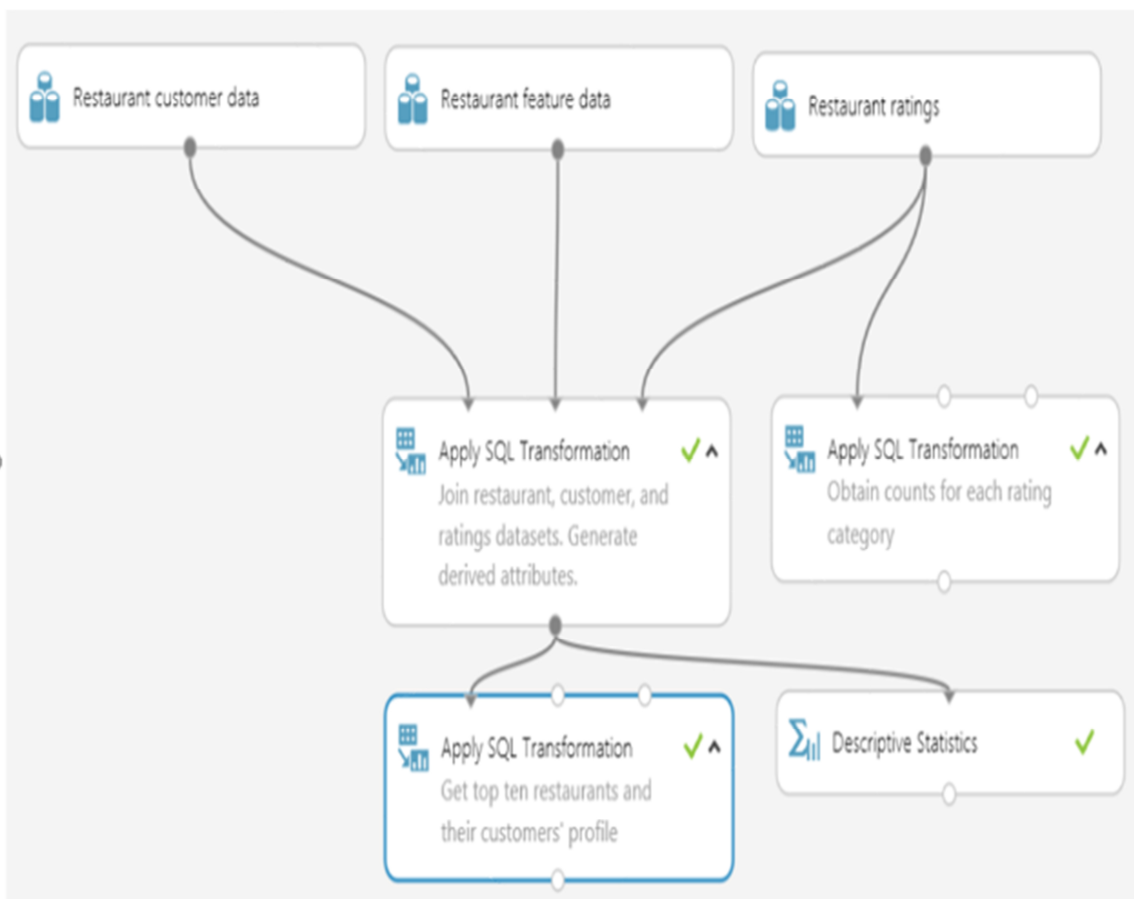


Рисунок 1 – процес експерименту

¹Студент групи М КН 2020-1, ХНУМГ ім. О.М.Бекетова

² К.т.н., доцент кафедри КНтаІТ ХНУМГ ім. О.М.Бекетова

По-перше, було проведено просте дослідження набору даних рейтингів ресторанів і отримали кількість ресторанів для кожної рейтингової категорії. В цьому випадку в модуль «Застосувати перетворення SQL» було надано тільки один вхідний набір даних (рис. 2).

```
▲ Apply SQL Transformation
SQL Query Script
1 SELECT rating, COUNT(rating) AS count
2 FROM t1
3 GROUP BY rating
4 ORDER BY rating
5 ;
```

Рисунок 2 – Дослідження набору даних рейтингів ресторану

Результат зображено нижче (рис. 3):

rating	count
0	254
1	421
2	486

Рисунок 3 – Кількість ресторанів для кожної рейтингової категорії

По-друге, було зроблено внутрішнє з'єднання всіх трьох наборів даних в SQL на основі ідентифікатора користувача і ідентифікатора місця. Крім того, ми витягли такі атрибути, як Вік (у 2015 році), на основі Birth_year з даних про відвідувачів ресторану (рис. 4).

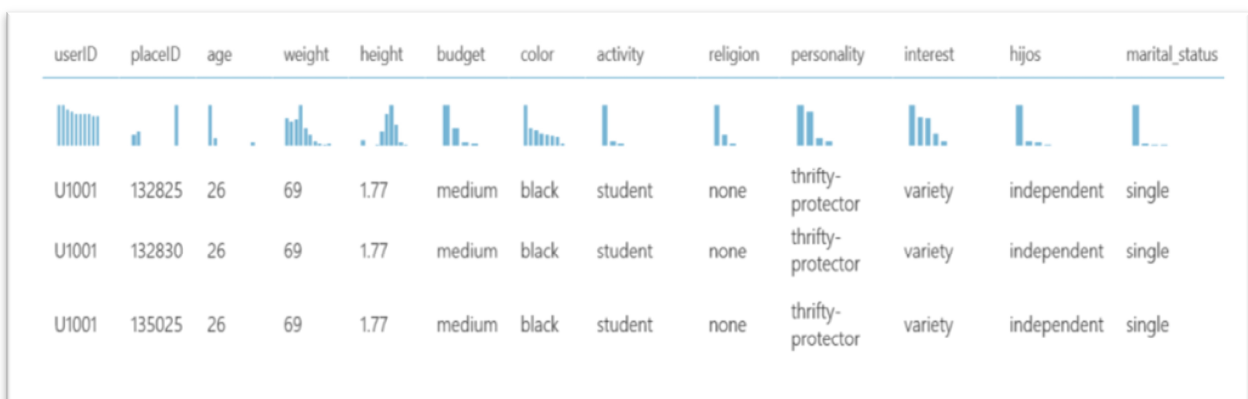
Apply SQL Transformation

SQL Query Script

```
1 SELECT |
2 t3.userID AS userID,
3 t3.placeID AS placeID,
4 (2015-t1.birth_year) AS age,
5 weight,
6 height,
7 budget,
8 color,
9 activity,
10 religion,
11 personality,
12 interest,
13 hijos,
14 marital_status,
15 transport,
16 ambience,
17 dress_preference,
18 drink_level,
19 smoker,
20 abs(t1.latitude-t2.latitude) AS latitude_diff,
21 abs(t1.longitude-t2.longitude) AS longitude_diff,
22 city,
23 state,
24 country,
25 alcohol,
26 smoking_area,
27 dress_code,
28 accessibility,
29 price,
30 Rambience,
31 franchise,
32 area,
33 other_services,
34 rating
35 FROM t1, t2, t3
36 WHERE t1.userID =t3.userID
37 AND t2.placeID=t3.placeID;
```

Рисунок 4 – Внутрішнє з'єднання трьох наборів даних

Результат зображено нижче (рис. 5):



userID	placeID	age	weight	height	budget	color	activity	religion	personality	interest	hijos	marital_status
U1001	132825	26	69	1.77	medium	black	student	none	thrifty-protector	variety	independent	single
U1001	132830	26	69	1.77	medium	black	student	none	thrifty-protector	variety	independent	single
U1001	135025	26	69	1.77	medium	black	student	none	thrifty-protector	variety	independent	single

Рисунок 5 – Результат SQL запроса

Нарешті, було використано об'єднання результатів трьох таблиць (зверху) і обчислено такі показники, як середній рейтинг і загальні оцінки, отримані для кожного ресторану (рис. 6).

```

Apply SQL Transformation
SQL Query Script
1 SELECT
2 placeID,
3 AVG(rating) AS Average_Rating,
4 COUNT(rating) AS Num_of_Ratings,
5 price,
6 Rambience,
7 smoking_area,
8 AVG(age) AS Customer_Average_Age,
9 AVG(weight) AS Customer_Average_Weight,
10 AVG(height) AS Customer_Average_Height
11 FROM t1
12 GROUP BY
13 placeID
14 ORDER BY Average_Rating DESC
15 LIMIT 10
16 :

```

Рисунок 6 – Обчислення середнього рейтингу і оцінки ресторанів

Десять кращих ресторанів, відсортованих за середньою оцінкою в порядку убубання (на виході з модуля «Застосувати перетворення SQL»), перераховані нижче (рис. 7).

placeID	Average_Rating	Num_of_Ratings	price	Rambience	smoking_area	Customer_Average_Age	Customer_Average_Weight	Customer_Average_Height
132955	2	5	low	familiar	none	35.6	47.8	1.564
134986	2	8	high	familiar	none	33.625	67.125	1.65875
135034	2	5	medium	familiar	none	28	53.6	1.67
132922	1.833333	6	medium	familiar	permitted	25.333333	57.666667	1.628333
132755	1.8	5	medium	familiar	none	27.4	59	1.69
134976	1.75	4	low	familiar	none	30.25	60.25	1.6425
135013	1.75	4	low	familiar	none	24.25	70.75	1.6275
135074	1.75	4	high	familiar	section	28.75	54.5	1.665
135055	1.714286	7	high	familiar	section	25.571429	59.285714	1.594286
135075	1.692308	13	medium	familiar	none	29.538462	60.384615	1.686154

Рисунок 7 – Відображення десяти кращих ресторанів

ПЕРЕЛІК ПОСИЛАНЬ

1. Tan P-N., Steinbach M., Karpatne A. and Kumar V. “Introduction to data mining”. Pearson; 2nd edition. 2018. 864p.
2. Ben-Gan, Itzik. “T-SQL Fundamentals”. 3rd edition. 2016. 235p.