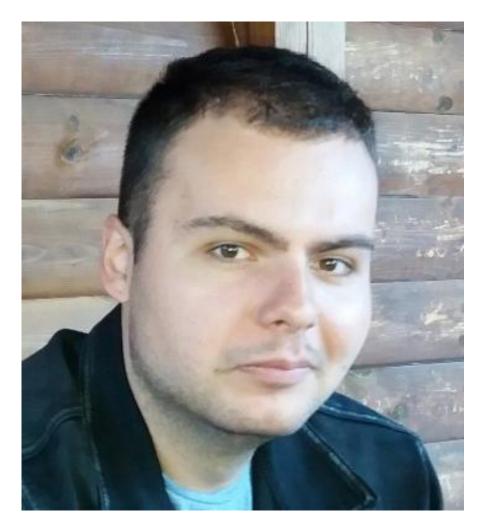
Име на проекта: *Анализ на потребителски данни за препоръчителна система*



Автор: Николай Вецов – Data Analysis & AI Junior Freelancer

Задача: Анализ на потребителски данни за препоръчителна система

📽 Цел

Представи си, че работиш в екип, който изгражда препоръчителна система за онлайн платформа. Твоята задача е да анализираш данните за потребителска активност и да извлечеш основни статистики и зависимости.

Данни

Ще работиш със следната примерна таблица:

	User_ID	Product_ID	Rating	Timestamp
	101	A234	4	06/01/2025 14:25
	102	B678	5	06/01/2025 14:27
	101	B678	3	06/01/2025 14:30
•	103	A234	2	06/01/2025 14:33

Таблицата може да бъде записана като CSV файл: user ratings.csv

🛠 Задачи

1. **C SQL**:

- о Извлечи уникалните потребители и броя оценки, които всеки е дал.
- о Намери средната оценка за всеки продукт.

2. C Excel:

- Зареди файла и използвай Pivot таблица, за да покажеш средната оценка по продукт.
- Филтрирай по оценки ≥ 4.

3. C Python (c Pandas):

python

```
import pandas as pd

df = pd.read_csv("user_ratings.csv")
print(df.groupby("Product_ID")["Rating"].mean())
print(df.groupby("User_ID").size())
```

Така, аз, Николай Вецов съм на ход:

1. В SQL качих файла, първоначално файл XLSX (Екселски файл), после конвертиран в SQL, след качването в системата на SQL workbench, установих следното:

	User_ID	Product_ID	Rating	Timestamp
	101	A234	4	06/01/2025 14:25
	102	B678	5	06/01/2025 14:27
	101	B678	3	06/01/2025 14:30
•	103	A234	2	06/01/2025 14:33

При установяване на средната оценка, установих следното:

Имаме общо 3 уникални потребители, които са 101 (гласувал 2 пъти), 102 и 103. Продуктите са A234 и B678.

За продукт А234, средната оценка е 3. Защо? Защото имаме две оценки до тук и те са 4 и 2 и цифрата между тях е 3 и затова.

За продукт В678, средната оценка е 4. Тъй като има 5 и 3 и така средната оценка е 4.

Извод:

- Средната оценка на продукт А234 е 3.
- Средната оценка на продукт В678 е 4.
 - 2.След качване на данните в Microsoft Excel, след надстойване се появи средната оценка и филтрирах по оценка ≥ 4 .

A	D		U	E	Г	G	П	'	,		L	IVI	IN	0	r	
ow Labels 🔻	Average of Rating		Row Labels 🎜	Average of Rating		Based on	what we se	ee for the	roduct wi	th ID A234	the avera	ge rating is	3, becaus	e there are	two	
234	3		B678	4		ratings: 4	from uniq	ue user 10	1) and 2 (f	rom uniqu	e user 103)	. Therefor	e, by this l	ogic, the av	/erage	
78	4		Grand Total	4		rating is 3. The same applies to product B678, which has two ratings: 5 (from user 102) and 3 (from user										
rand Total	3,5					101). Thus, the average rating is 4.										
						In the upp	er-left par	t of the tal	ole, the av	erage ratir	g for both	products i	s displaye	d. In the tal	ble to the	
User_ID	Product_ID	Rating	Timestamp			right, the average rating for product B678 is shown, because the task I was given requires filtering the										
101	A234	4	1.6.2025 14:25			data by ra	tings≥4 (g	reater tha	n or equal	to 4).						
102	B678	5	1.6.2025 14:27													
101	B678	3	1.6.2025 14:30													
103	A234	2	1.6.2025 14:33													
отребител 1 педователн	което виждаме за п 03). Следователно п о средната оценка е о задачата, която ми	о тази логика, и 4. Горе вляво на	маме средна оце гтаблицата е пок	нка = 3. Същото се азано за всеки еди	отнася и за ин от двата п	продукт с і родукта ср	номер В67 едната оц	8. Общо д енка. А в т	ве оценки	- 5 (потре	бител с но	мер 102) і	13 (потре	бител с но	мер 101).	
∕ъ, гъи като																_
го, тый като																
/ъ, гъи като																
га, гый като																

3. След стартиране на кода на Python, получих следния резултат:

	• •	<u> </u>
4	А	В
	Product_ID	Average_Rating
	B678	4
	A234	3

📆 Среден рейтинг по продукт. (4 и 3)

А	В			
User_ID	Ratings_Count			
101	2			
102	1			
103	1			

Потребителска активност. (потребител 101 - 2 пъти гласувал, потребител 102 - 1 път и потребител 103 - 1 път.

A	В	С	D
User_ID	Product_ID	Rating	Timestamp
101	A234	4	1.6.2025 14:25
102	B678	5	1.6.2025 14:27

Продукти с рейтинг ≥4 – т.е да показва продукти с рейтинг по-голям или = на 4. Оказва се, че и двата продукта имат такъв рейтинг.

Α Α	В	С
User_ID	A234	B678
101	4	3
102		5
103	2	

Рейтинг матрица (таблица, която показва **оценките на потребители към различни обекти).** Според нея – потребител 101 е оценил продукт A234 и продукт B678. Потребител 102 е оценил само и единствено продукт B678, а потребител 103 е оценил само и единствено продукт A234.

ВСИЧКО ТОВА, КОЕТО ПОКАЗАХ В ТЕЗИ 4 ФИГУРИ, ПОКАЗВА АНАЛИЗА НА ПОТРЕБИТЕЛИТЕ ПРИ РАЗЛИЧНИТЕ ПРОДУКТИ. ВСЕКИ ОТ ТЯХ ОЦЕНЯВА РАЗЛИЧНИТЕ ПРОДУКТИ ПО РАЗЛИЧЕН НАЧИН. ОКАЗВА СЕ, ЧЕ ПРОДУКТ В678 Е ПОВЕЧЕ ПРЕДПОЧИТАН, ЗА РАЗЛИКЕ ОТ ПРОДУКТ А234. И ВСИЧКО ТОВА Е КРАЙНИЯТ РЕЗУЛТАТ НА ТОЗИ ПОТРЕБИТЕЛСКИ АНАЛИЗ. ПО-НАДОЛУ ЩЕ ВИДИТЕ И КОДА, КОЙТО ИЗПОЛЗВАХ НА РҮТНОN.

Кодът, който използвах (Python):

```
import pandas as pd
# Зареждане на CSV с правилен разделител
df = pd.read_csv("user_ratings.csv", sep=";", encoding="utf-8")
# Анализ 1: Средна оценка по продукт
product_avg =
df.groupby("Product_ID")["Rating"].mean().sort_values(ascending=False)
# Анализ 2: Активност на потребители
user_activity = df["User_ID"].value_counts()
# Анализ 3: Продукти с рейтинг ≥ 4
top_rated = df[df["Rating"] >= 4]
# Анализ 4: User-Product рейтинг матрица
rating_matrix = df.pivot_table(index="User_ID", columns="Product_ID",
values="Rating")
# Записване във Excel файл
with pd.ExcelWriter("user_analysis.xlsx") as writer:
    product_avg.to_frame(name="Average_Rating").to_excel(writer,
sheet_name="Среден рейтинг по продукт")
    user_activity.to_frame(name="Ratings_Count").to_excel(writer,
sheet_name="Активност на потребители")
    top_rated.to_excel(writer, sheet_name="Продукти с рейтинг ≥4",
index=False)
    rating_matrix.to_excel(writer, sheet_name="Рейтинг матрица")
```

С този код, генерирах самия краен резултат в Excel файл, които в момента e user_analysis.xlsx. Точно там е целият анализ.