

Üzleti intelligencia házi feladat dokumentáció

Autók árának előrejelzése, Netflix és Youtube nézettségek, vélemények

Inez Anna Papp

May 2025

1 Forráskód, videó

Link a forráskód GitHub repojához: https://github.com/pappinez/BI_hw
A videó a repoban található.

2 Bemutató

A feladatom folyamán 3 adathalmazt használok, az egyikben autók eladási árai szerepelnek, a másikban országonkénti youtube nézettségek két különböző hónapban, a harmadikban pedig Netflixes sorozat értékelések. Az autós adathalmaz idősoros adatokat tartalmaz, így a LSTM alapú predikcióhoz fogom alkalmazni, a másik két adathalmazon pedig ETL jobokat fogok végrehajtani és reportokat készíteni belőlük. A Netflixesnél a KPI, hogy melyiknek a legjobb az értékelése, a youtubeosnál, hogy melyik ország nézi a legtöbbet, az autósánál pedig, hogy az ár emelkedik vagy csökken különböző tényezők figyelembevételével. Az ETL jobokkal megtisztítom az adatokat, majd átalakítom őket az elvárt formátumokba. A reportoknál különböző megjelenítéseket készítek, melyik interaktív, és az adatokról mutatnak be különböző nézeteket, ezzel valamilyen új információt bemutatva.

3 Főbb funkciók összefoglaló

3.1 Általam választott adatforrások

- Autók eladási árai
- Országonkénti youtube nézettség
- Netflix sorozatok és értékelések

3.2 Adattárolás

- Raw réteg: az adatok .csv fileokban vannak tárolva
- Stage réteg: megtörténik az adattisztítás, hibás adatok kiszűrése
- Data Warehouse réteg: dimenzió bevezetése

3.3 Megvalósítandó ETL jobok

- Transzformáció 1: Hibás adatokat keres és kilogolja azokat a youtube-os adatforráson
- Job 1: Ha a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot
- Transzformáció 2: Adattisztítás, hibás értékek kiszűrése a Netflixes adatforráson
- Job 2: Ha a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot
- Transzformáció 3: Netflix értékelések alapján kategorizálja, hogy ez egy Jó, Közepes, vagy Rossz film, ehhez sávok vannak megadva, az összegzést kimentí egy JSON fileba
- Job 3: a jó filmek feltöltése webszerverre
- Transzformáció 4: Youtube nézettség csoportosítás kontinensenként, majd aggregáció: átlagosnézettség a megadott két hónapban
- Job 4: Ha a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot

Az adatokból olyan adattípus készül, ami alkalmas arra, hogy a PowerBI fel tudja használni

3.4 Reportok

- Report 1: Kördiagrammon ábrázolom a 3 legtöbb youtube felhasználóval rendelkező kontinenst, mellette egy táblázatban a kontinensek nevét és az átlagos nézettségüket a két megadott hónapban. A másik oldalon clustered bar charttal ábrázolom az országokat, és nézettségüket, és lefűrást használok.
- Report 2: Csatlakoztatom a youtube és netflix táblákat a Netflix nézettség ország és a filmet készítő ország alapján. Egy táblázatban megjelenítem azáltal a top 2 ország által készített filmeket, amely két országban a legmagasabb a youtube nézettség, az értékelésük szerint csökkenő sorrendben.

- Report 3: 3 táblázatban ábrázolom a Jó, Rossz és Közepes értékelésű Netflix filmeket, majd az x tengelyen évek szerint megjelenítem a kategóriákat egy scatter charton
- Report 4: Kördiagramon ábrázolom a Netflixes filmekben a műfajok arányát, egy másikon a film felfedezésének arányát

3.5 Data Science

LSTM hálózatot használva az autó árakról szóló idősoros adatok alapján készítek predikciót a jövő adatokról

4 Választott technológiák

- Pentaho az ETL transzformációkhoz és jobokhoz
- Power BI a reportokhoz
- Python a predikcióhoz

5 Transzformációk, Jobok

5.1 Transzformáció 1

Hibás adatokat keres és kilogolja azokat a youtube-os adatforráson.

Ehhez szükség volt a forrás mappában lévő fájlnevek összegyűjtésére, majd CSV-ből az adatok beolvasására. Utána a nézettségszám adatokból kifiltereltem a 0-nál kisebb értékeket, végül adatvalidációval kiszűrtem a nem 2 nagy betűből álló zászló kódokat, és egy .txt-be kimentettem őket. A tisztított adatokat egy .csv-be, és a PowerBI-hoz Excelként is mentettem. *Figure 1 Figure 2*

5.2 Job 1

Ha a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot

Megnéztem, hogy üres-e a mappa, ha igen, akkor kilogoljuk ezt, ha nem, akkor viszont elvégzi rajta a transtformációt. *Figure 3*

5.3 Transzformáció 2

Adattisztítás, hibás értékek kiszűrése a Netflixes adatforráson

Ehhez szükség volt a forrás mappában lévő fájlnevek összegyűjtésére, majd CSV-ből az adatok beolvasására. Utána kifiltereltem azokat az adatokat, ahol nem 0 és 10 közötti az értékelés, végül adatvalidációval kiszűrtem azokat az évszámokat, ami nem az 1900-as, vagy 2000-es években keletkezett, és egy .txt-be kimentettem őket. A tisztított adatokat egy .csv-be, és a PowerBI-hoz Excelként is mentettem.

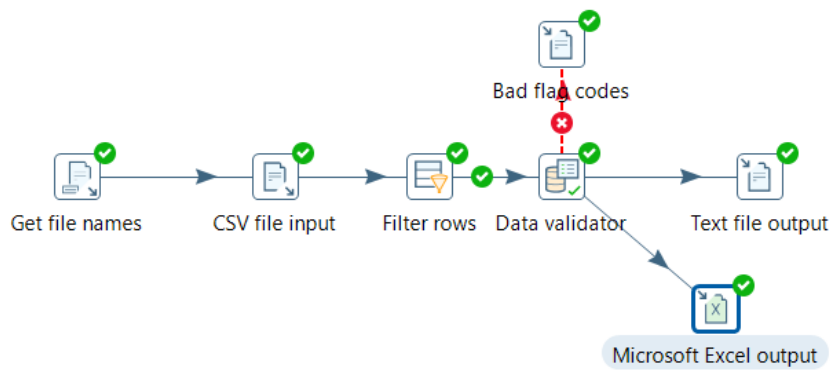


Figure 1: Transformation 1

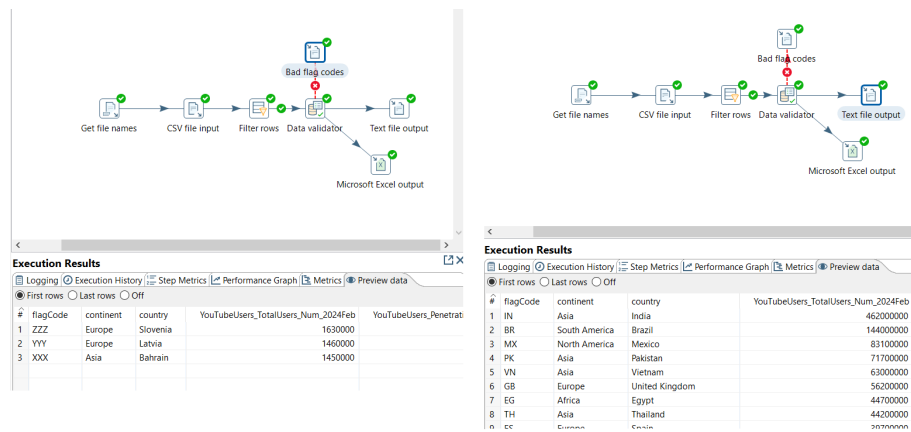


Figure 2: Result of Transformation 1

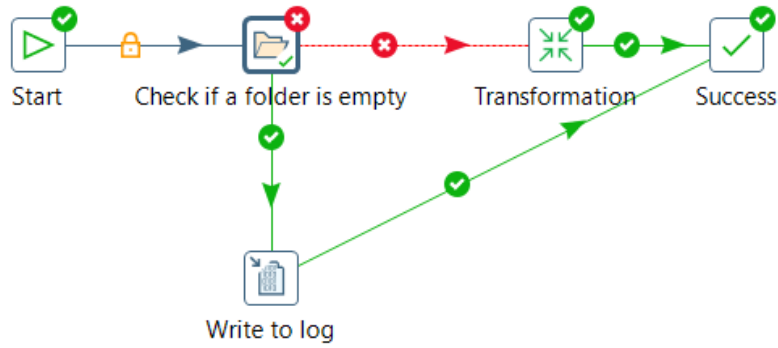


Figure 3: Job 1

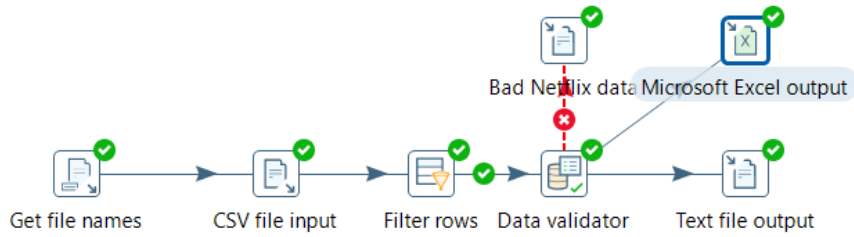


Figure 4: Transformation 2

Figure 4 Figure 5

5.4 Job 2

A a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot

Megnéztem, hogy üres-e a mappa, ha igen, akkor kilogoljuk ezt, ha nem, akkor viszont elvégzi rajta a transtformációt. *Figure 6*

5.5 Transzformáció 3

Netflix értékelések alapján kategorizálja, hogy ez egy Jó, Közepes, vagy Rossz film, ehhez sávok vannak megadva, a kategóriánként összegzést kimentti egy

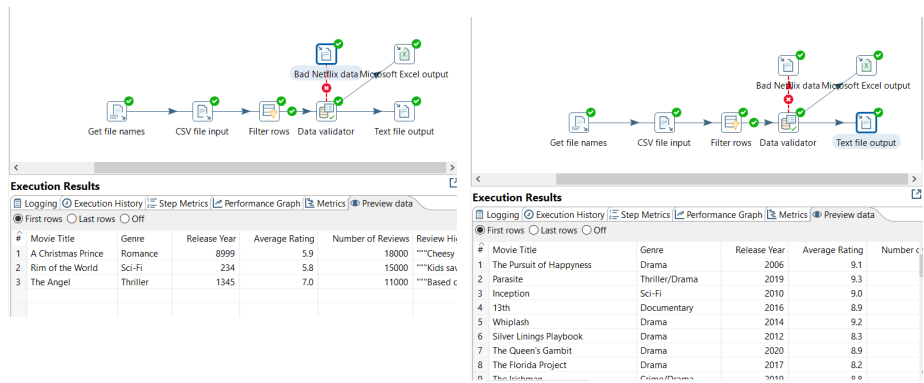


Figure 5: Result of Transformation 2

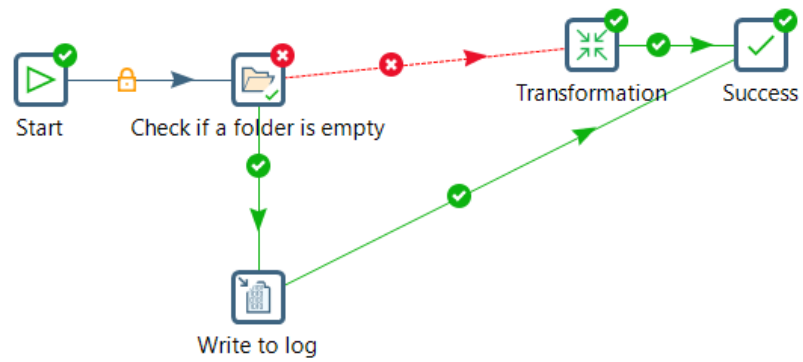


Figure 6: Job 2

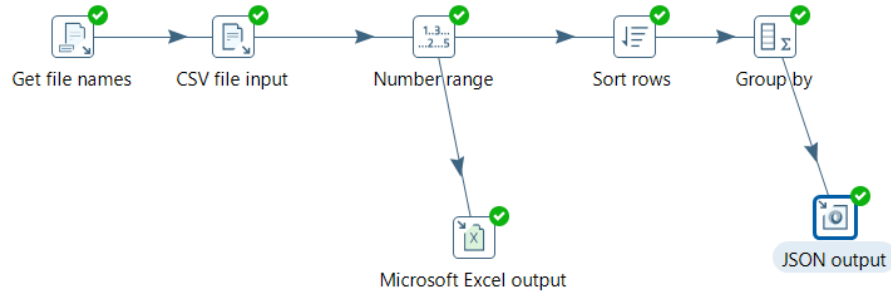


Figure 7: Transformation 3

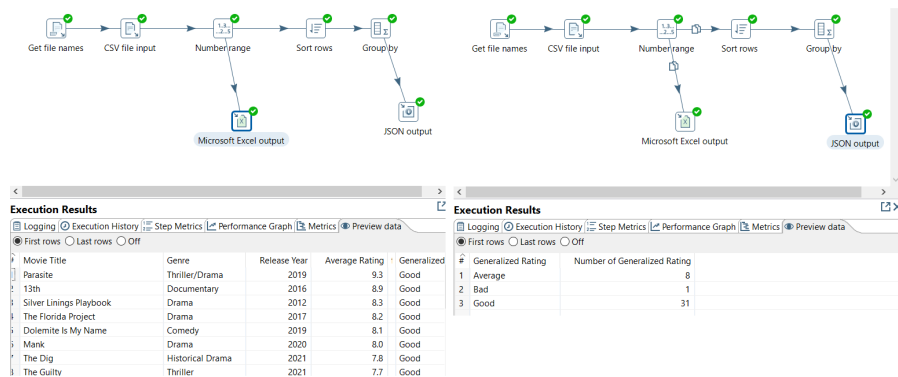


Figure 8: Result of Transformation 3

JSON fileba.

Ehhez szükség volt a forrás mappában lévő fájlnevek összegyűjtésére, majd CSV-ből az adatok beolvasására. Ezután 1-5-ig 'Bad', 5-7-ig 'Average' 7-10-ig 'Good' értékelést kaptak, majd értékelés szerinti rendezés után JSON-ban kimentettem az adatokat. Egy Excel formátumú mentés is készült a PowerBI-hoz. *Figure 7*

5.6 Job 3

A jó filmek feltöltése webszerverre.

Egy HTTP requesttel feltöltöttem webszerverre, ha pedig nem sikerült, azt kilogoltam. *Figure 9*

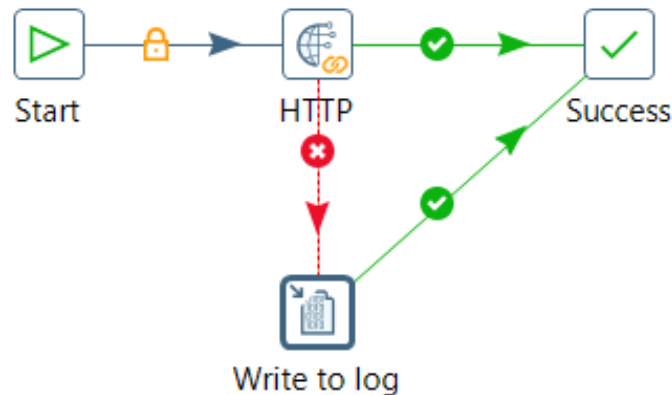


Figure 9: Job 3

5.7 Transzformáció 4

Youtube nézettség csoportosítás kontinensenként, majd aggregáció: átlagosnézettség a megadott két hónapban

Ehhez szükség volt a forrás mappában lévő fájlnevek összegyűjtésére, majd CSV-ből az adatok beolvasására. Ezután rendeztem, és az átlagos nézettséget kiszámoltam a két adott hónapra és .csv-be mentettem. *Figure 10 Figure 11*

5.8 Job 4

Ha a bemenetnek megadott mappába került egy .csv file, akkor elvégzi a transzformációt, ha nem, akkor erről készít logot

Megnéztem, hogy üres-e a mappa, ha igen, akkor kilogoljuk ezt, ha nem, akkor viszont elvégzi rajta a transzformációt. *Figure 12*

6 Reportok

6.1 Report 1

Kördiagrammon ábrázolom a 3 legtöbb youtube felhasználóval rendelkező kontinentst, mellette egy táblázatban a kontinensek nevét és az átlagos nézettségüket a két megadott hónapban. A másik oldalon clustered bar charttal ábrázolom

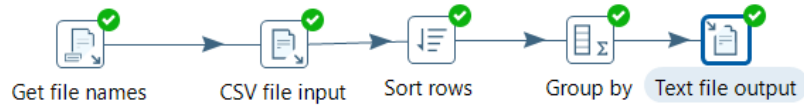


Figure 10: Transformation 4

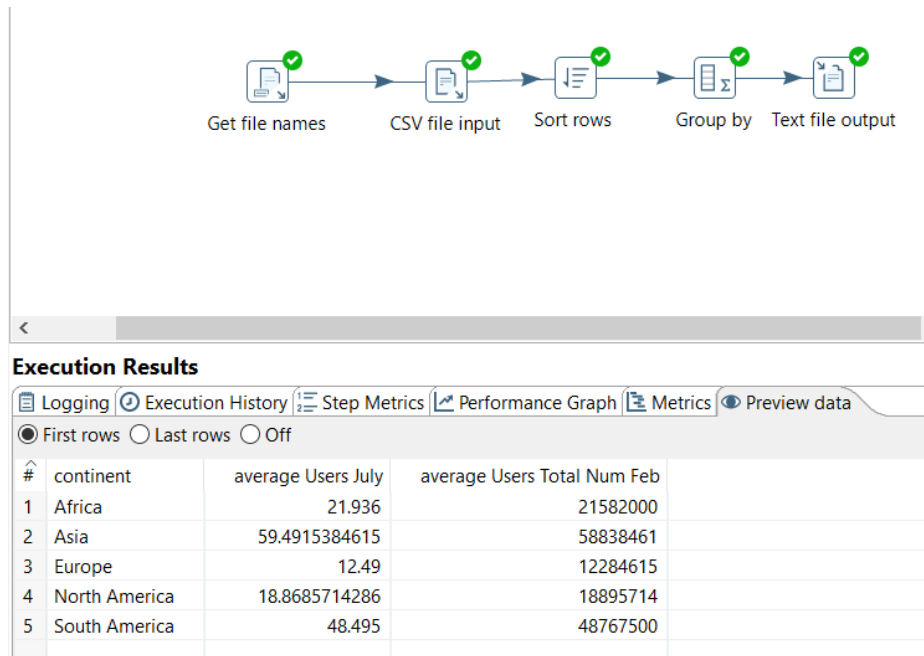


Figure 11: Result of Transformation 4

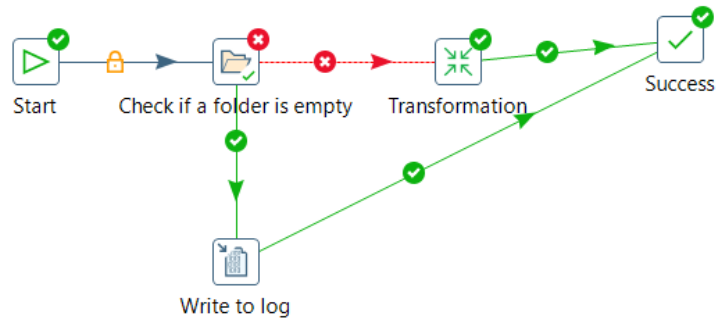


Figure 12: Job 4

az országokat, és nézettségüket, és lefűrást használlok. *Figure 13 Figure 14 Figure 15*

6.2 Report 2

Csatlakoztatom a youtube és netflix táblákat a Netflix nézettség ország és a filmet készítő ország alapján. Egy táblázatban megjelenítem azáltal a top 2 ország által készített filmeket, amely két országban a legmagasabb a youtube nézettség, az értékelésük szerint csökkenő sorrendben. *Figure 16 Figure 17*

6.3 Report 3

Táblázatban ábrázolom a Jó, Rossz és Közepes értékelésű Netflix filmeket, majd az x tengelyen évek szerint megjelenítem a kategóriákat egy scatter charton *Figure 18 Figure 19*

6.4 Report 4

Kördiagramon ábrázolom a Netflixes filmekben a műfajok arányát, egy másikon a film felfedezésének arányát *Figure 20 Figure 21*

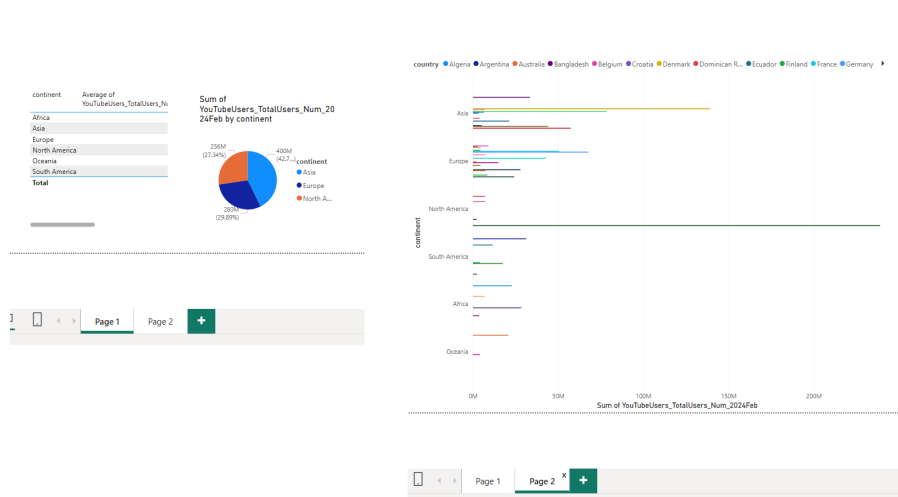


Figure 13: Report 1

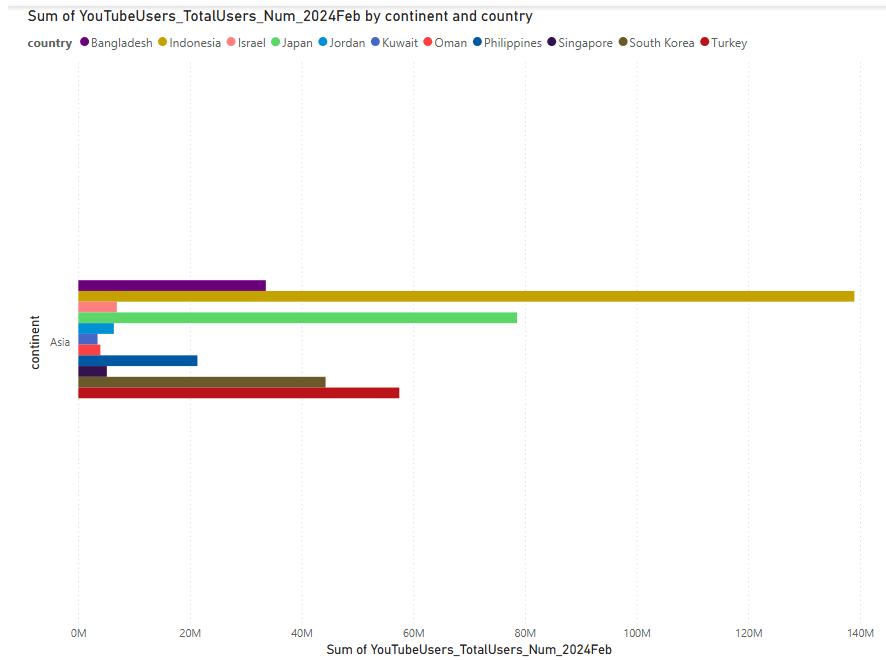


Figure 14: Report 1 Drill down

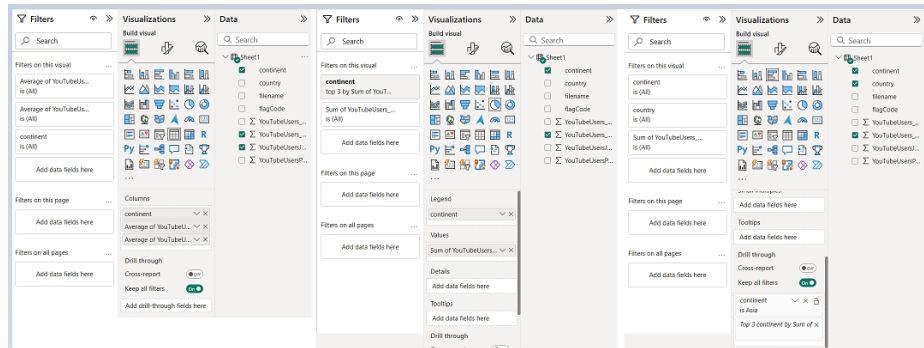


Figure 15: Settings for Report 1

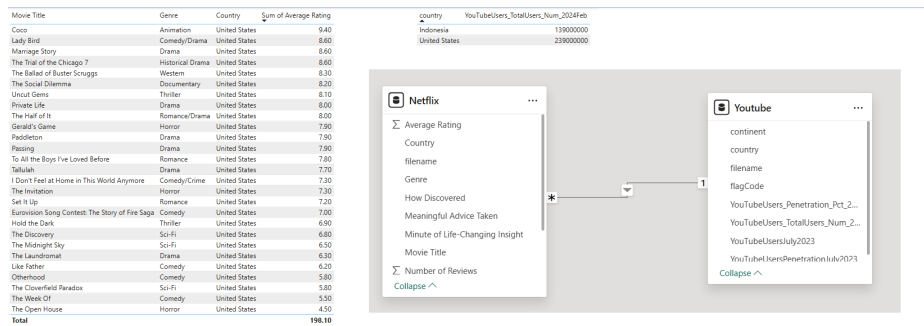


Figure 16: Report 2

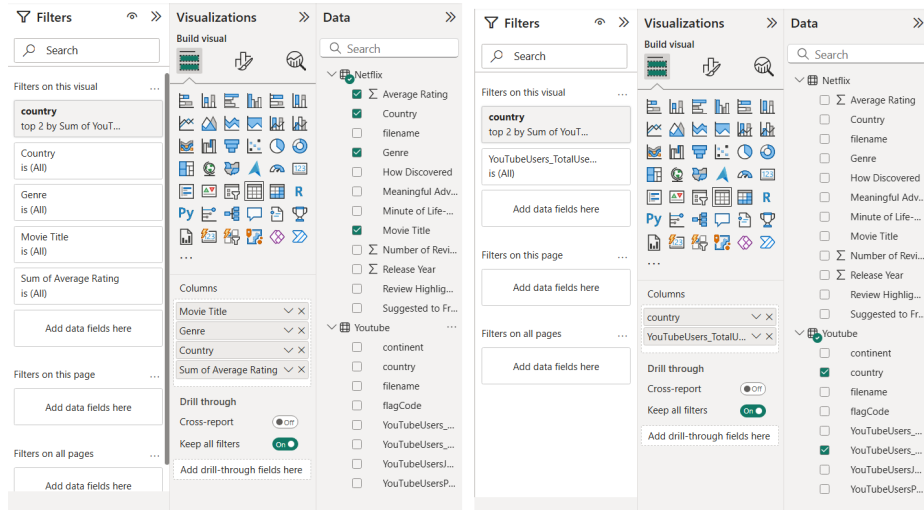


Figure 17: Settings for Report 2

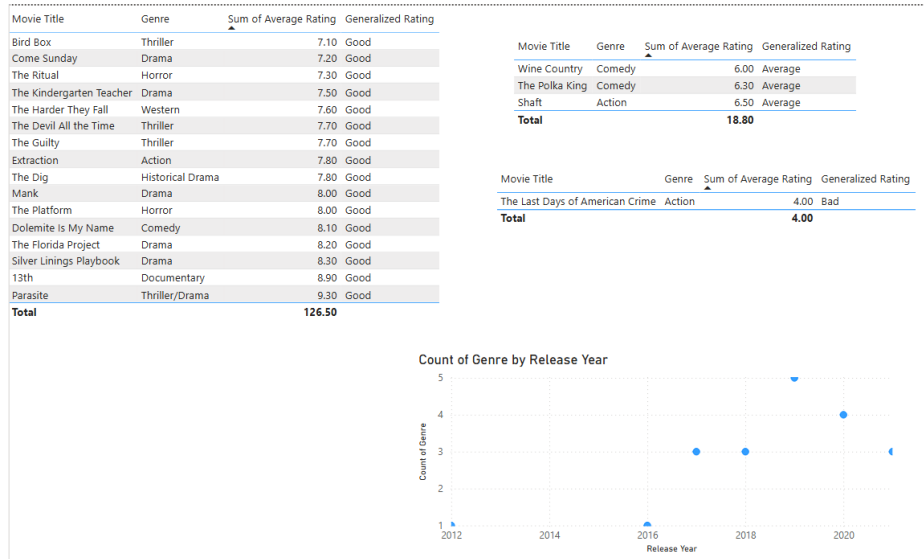


Figure 18: Report 3

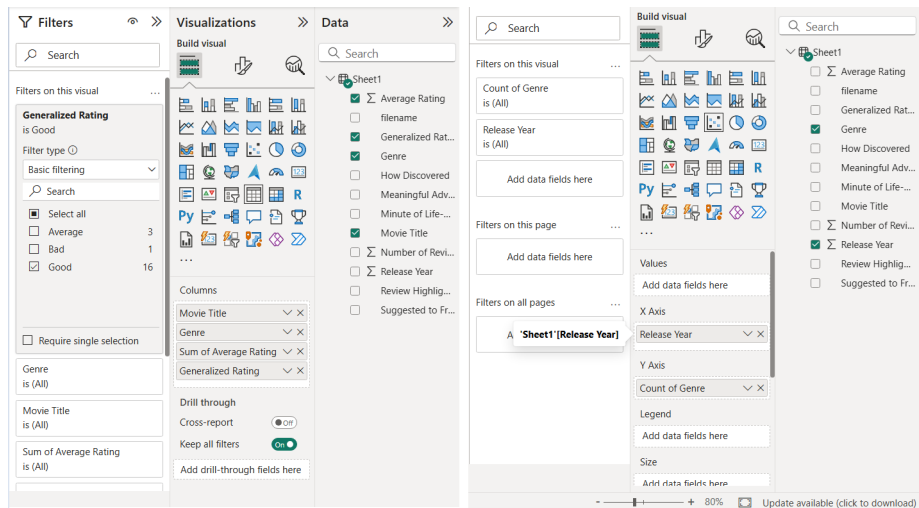
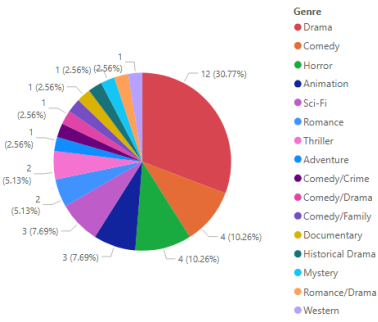


Figure 19: Setting for Report 3

Count of Genre by Genre



Count of How Discovered by How Discovered

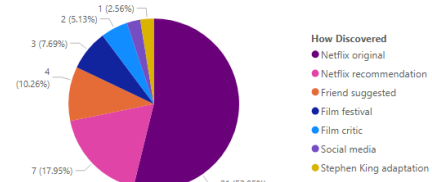


Figure 20: Report 4

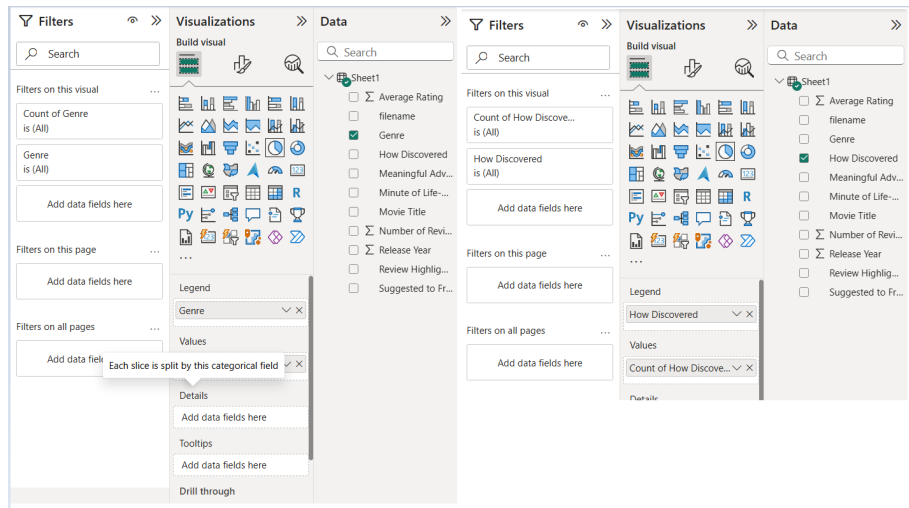


Figure 21: Settings for Report 4

7 Idősoros előrejelzés Pythonban

Az alábbi linken megtalálható a forráskód és kommentek formájában a dokumentáció:

<https://colab.research.google.com/drive/1gUd3XhPzx1xIK0XbaRlgnwsPhP54fL6K?usp=sharing>

Szükség lehet a sample data mappába feltölteni a következő fület, amit a házi feladat GitHub repojában lehet elérni: `automobile_prices_economics_2019_2023`