

UNIVERSITÀ DEGLI STUDI DI VERONA

Big Data

REPORT DEL PROGETTO

Mattia Zorzan - VR464472

9 dicembre 2021

Indice

1	Descrizione del Progetto	2
2	Descrizione delle Interrogazioni	3
2.1	Query 1	3
2.2	Query 2	4

1 Descrizione del Progetto

Il progetto si propone di eseguire alcune interrogazioni sul dataset MovieLens. Nello specifico queste vengono eseguite sulla variante **MovieLens 1M**, di dimensione minore.

Di seguito un elenco delle interrogazioni richieste:

- **EXPLORATORY ANALYSIS**

- *Number of ratings* for each **movie** (and its *distribution*)
- *Number of ratings* for each **user** (and its *distribution*)
- *Average* score received by each **movie**
- *Average* score given by each **user**
- Top **K** *movies* with at least **R** *ratings*

- **ADVANCED QUERIES**

- Find if there is a *correlation* between the *standard deviation* of the ratings a movie has received, and the *number of ratings*
- Find the *evolution over time* (with a granularity of **N** months) of the *number of ratings* and the *average rating*: do high rated movies maintain their ratings? Are low rated movies “abandoned” after a while?
- Find how the *text* of each movie changes as we progressively **remove** the ratings from users that rated more and more movies. For instance, we can identify *different groups of users* (who rated less than 10 movies, who rated between 11 and 30 movies, ...) and we can compute the *average rating* considering all the groups, then only the groups of users with at least 11 ratings, and so forth
- Is it possible to identify *groups of similar movies* based on the ratings they received from the users? For instance, if movies **m1** and **m2** have both obtained 5 stars from users **u1** and **u2**, they may be considered similar

Si poteva scegliere se eseguire tutte le interrogazioni oppure solo un sottinsieme di esse.

2 Descrizione delle Interrogazioni

Per comodità, in questo report è stata omessa la descrizione delle soluzioni per l'*exploratory analysis*, verranno trattate solo le *advanced queries*.

2.1 Query 1

Per dimostrare la correlazione tra *standard deviation* ed il numero dei *ratings* si è in primo luogo ottenuto il numero di valutazioni per ogni film presente nel dataset. Limitando (in 4 *DataFrame* diversi) il numero di film presi in considerazione si può a questo punto vedere come la diminuzione dei campioni presi in considerazione "*ammorbidisca*" la curva. I 4 *DataFrame* sono:

1. **1M Samples:** Prende in considerazione l'intero dataset
2. **100K Samples:** Prende in considerazione i primi 100.000 film del risultato del conteggio
3. **100K Samples:** Prende in considerazione i primi 1.000 film del risultato del conteggio
4. **100K Samples:** Prende in considerazione i primi 100 film del risultato del conteggio

Il campione può essere considerato "*randomico*" in quanto non viene fatto alcun tipo di ordinamento sui risultati della query, che vengono già restituiti in ordine sparso.

Dal plot dei dati la forma della distribuzione potrebbe sembrare **gaussiana**, per verificare ciò è stato eseguito uno **skewness** test su tutti e 4 i *DataFrame*, dando risultati nulli o negativi. Non è quindi una distribuzione normale.

Osservando l'andamento del grafico è possibile riconoscere un andamento iperbolico, riconducibile ad una distribuzione **paretiana**. Non avendo trovato alcun test di *paretianità* già implementato è stata generata una distribuzione paretiana generica tramite il metodo **genpareto** di *scipy*. Accostando i due grafici la somiglianza risulta evidente.

2.2 Query 2

Per comodità nell'analisi è stato deciso di riferirsi ai quartili, quindi $N=4$, per il calcolo delle medie. Come primo passo è stata creata una **UDF** (*User Defined Function*) che andasse ad etichettare ogni riga del dataset come segue:

- **Q1:** Etichetta di tutte le *Row* aventi campo **Date** compreso tra *Aprile* e *Giugno 2000*
- **Q2:** Etichetta di tutte le *Row* aventi campo **Date** compreso tra *Aprile* e *Giugno 2000*
- **Q3:** Etichetta di tutte le *Row* aventi campo **Date** compreso tra *Aprile* e *Giugno 2000*
- **Q4:** Etichetta di tutte le *Row* aventi campo **Date** compreso tra *Gennaio* e *Marzo 2001*
- **Remaining:** Etichetta di tutte le *Row* aventi campo **Date** successivo a *Aprile 2001* (compreso)

In seguito all'etichettatura delle *Row* si è implementata un'ulteriore **UDF**, che andasse a valutare la "frequenza" nei quartili di appartenenza alla categoria **High Rated** o alla categoria **Low Rated**. Per differenziare le due classi si è deciso di utilizzare il voto a *3 stelle* come threshold, se per la maggior parte dei quartili il film è etichettato come **High Rated** allora la sua etichetta "globale" sarà quella, **Low Rated** altrimenti.

In seguito a quest'operazione, possiamo vedere che in media i film etichettati come **High Rated** sono ottengono in media il doppio delle recensioni dopo l'anno rispetto agli altri.

Con la seconda cella Jupyter invece si va ad analizzare la stabilità nel numero di rating dopo l'anno.

Dallo *Scatter* possiamo vedere che per entrambe le categorie c'è una tendenza all'abbandono, quasi totale per i film **Low Rated** mentre più mitigata per i film **High Rated**.

Segno evidente di questo è la grande presenza di film "apprezzati" nel range [40, 100] (pallini *blu*), sono invece quasi del tutto assenti nello stesso range i film "poco apprezzati".