



GeoTEDx

Matteo Cesari (MAT. 1073570)

Davide Girolamo (MAT. 1073645)

GitHub: <https://github.com/TurnTheBait/GeoTEDx>

Job AWS Glue

Job Create_Data_Lake

Lo scopo principale dello script è elaborare, aggregare e archiviare i dati dei dataset TEDx su MongoDB, consentendo un'analisi e una consultazione più efficienti.

- Data cleaning: rimozione dei dati nulli e duplicati per tutti i dataset.
- PySpark Job: creazione delle collections per identificare ogni video.

```
##### READ TAGS DATASET
tags_dataset_path = "s3://bucket-dati-test1-2023/tags_dataset.csv"
tags_dataset = spark.read.option("header", "true").csv(tags_dataset_path)
```

```
##### CLEAN TAGS DATASET
print(f"TOTAL TAGS DATASET: {tags_dataset.count()}")
tags_dataset = tags_dataset.dropDuplicates()
#### REMOVE DUPLICATES
print(f"TAGS DATASET without DUPLICATES {tags_dataset.count()}")
```

```
##### CLEAN TEDX DATASET
print(f"TOTAL DATASET IDX: {tedx_dataset.count()}")
tedx_dataset = tedx_dataset.filter(length("idx") == 32) ##### REMOVE INVALID IDX
print(f"DATASET without INVALID IDX: {tedx_dataset.count()}")

tedx_dataset = tedx_dataset.dropDuplicates() ##### REMOVE DUPLICATES
print(f"DATASET without DUPLICATES: {tedx_dataset.count()}")

##### FILTER ITEMS WITH NULL POSTING KEY
count_items = tedx_dataset.count()
count_items_null = tedx_dataset.filter("idx is not null").count()

print(f"Number of items from RAW DATA {count_items}")
print(f"Number of items from RAW DATA with NOT NULL KEY {count_items_null}")
```

Job AWS Glue

Job Watch_Next

Lo scopo principale dello script è quello di incorporare all'interno della documentazione su MongoDB di ogni talk, l'array watch_next contenente gli id dei video consigliati.

- Data cleaning: rimozione dei dati nulli e duplicati per tutti i dataset.

```
##### READ WATCH NEXT DATASET
watch_next_dataset_path = "s3://bucket-dati-test1-2023/watch_next_dataset.csv"
watch_next_dataset = spark.read.option("header", "true").csv(watch_next_dataset_path)
||
##### CLEAN WATCH NEXT DATASET
print(f"TOTAL WATCH NEXT DATASET: {watch_next_dataset.count()}")
watch_next_dataset = watch_next_dataset.dropDuplicates() ##### REMOVE DUPLICATES
print(f"WATCH NEXT DATASET without DUPLICATES {watch_next_dataset.count()}")
```

- Aggregate model: Creazione del modello aggregato aggiungendo i dati "watch_next" al dataset aggregato TEDx e Tags.

```
# CREATE THE AGGREGATE MODEL, ADD TAGS TO TEDX_DATASET
tags_dataset_agg = tags_dataset.groupBy(col("idx").alias("idx_ref")).agg(collect_list("tag").alias("tags"))
tags_dataset_agg.printSchema()
tedx_dataset_agg = tedx_dataset.join(tags_dataset_agg, tedx_dataset.idx == tags_dataset_agg.idx_ref, "left") \
    .drop("idx_ref") \
    .select(col("idx").alias("_id"), col("*")) \
    .drop("idx") \
tedx_dataset_agg.printSchema()
```

Collection su MongoDB

Collection: **tedx_data**

- Ogni documento della collezione rappresenta un talk. Ogni elemento è identificato da un id univoco, vengono inoltre forniti varie informazioni riguardo il contenuto e le visualizzazioni del video.
- I **tag** associati al talk vengono inseriti in un array.
- L'array **watch_next** contiene gli id dei talk indicati come "watch_next" del talk considerato.

```
1  _id: ObjectId('64662e740bd25e13055800cd')      ObjectId
2  main_speak... : "TED Audio Collective"         String
3  title: "Introducing Body Stuff with Dr. Jen Gunter" String
4  detai... : "Should you do a juice cleanse? Is it ac" String
5  post... : "Posted May 2021"                    String
6  u... : "https://www.ted.com/talks/ted_audio_collect" String
7  num_vie... : "0"                               String
8  durati... : "2:10"                             String
9  ▼ tags: Array                                  Array
10     0: "human body"                             String
11     1: "science"                                String
12     2: "society"                                String
13     3: "TED"                                    String
14     4: "health"                                 String
15     5: "talks"                                  String
10  ▼ watch_next: Array                          Array
17     0: "a2a717eaff1c4f504dc36dc908285207"      String
18     1: "9f7b1654e792011b7e1c6f4288520226"      String
19     2: "3d09fa5a4a64a82554252244c9420355"      String
20     3: "c4b0ade4a4862ecfc44b19adb7db339d"      String
```

Sviluppi futuri

- Inserire una maggiore quantità di tag all'interno dei dati di ciascun talk in modo più attinente per migliorare la qualità dei video suggeriti
- Prendere in considerazione la data di pubblicazione dei talk per consigliare video più recenti

