

Table Schema Design:

- `movie_id`: Type: `int`, Description: Partition key
- `user_id`: Type: `int`, Description: Clustering key, for compound primary key
- `rating`: Type: `int`, Description: Rating value
- `ts`: Type: `timestmap`, Description: Rating timestamp (converted from UNIX timestamp)

The primary key is defined as `(movie_id, user_id)`, indicating that each user can have only one rating record per movie.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:33:38 PM.

```
%pyspark
# 2.2 PySpark 读取 + 类型转换代码
from pyspark.sql.types import StructType, IntegerType, LongType
from pyspark.sql.functions import col, from_unixtime
```

FINISHED

```
# 定义 schema
schema_rating = StructType([
    StructField("user_id", IntegerType(), True),
    StructField("movie_id", IntegerType(), True),
    StructField("rating", IntegerType(), True),
    StructField("unix_ts", LongType(), True)
])
```

```
# 加载数据
rating_raw_df = spark.read.csv(
    "hdfs://user/maria_dev/ml-100k/u.data",
    sep="\t",
    schema=schema_rating
)
```

```
# 转换 UNIX 时间戳为 timestamp
rating_df = rating_raw_df.withColumn(
    "ts", from_unixtime(col("unix_ts")).cast("timestamp")
).select("movie_id", "user_id", "rating", "ts")
```

```
rating_df.show(5, truncate=False)
```

```
+-----+|movie_id|user_id|rating|ts|-----+
|242 |196 |3 |1997-12-04 15:55:49|
|302 |186 |3 |1998-04-04 19:22:22|
|377 |22 |1 |1997-11-07 07:18:36|
|51 |244 |2 |1997-11-27 05:02:03|
|346 |166 |1 |1998-02-02 05:33:16|
+-----+
```

only showing top 5 rows

Took 2 sec. Last updated by anonymous at June 19 2025, 9:33:43 PM.

```
%pyspark
# 2.3 写入 Cassandra 表 ratings
rating_df.write.partitionBy("user_id").format("org.apache.spark.sql.cassandra") \
    .option("table", "ratings") \
    .option("keyspace", "movielens") \
    .mode("append") \
    .save()
```

FINISHED

Took 1 min 3 sec. Last updated by anonymous at June 19 2025, 9:34:45 PM.

```
%pyspark
# 2.4 验证写入成功
spark.read \
    .format("org.apache.spark.sql.cassandra") \
    .options(table="ratings", keyspace="movielens") \
    .load().show(10, truncate=False)
```

FINISHED

```
+-----+|movie_id|user_id|rating|ts|-----+
|1281 |98 |5 |1998-03-31 22:34:47|
|79 |1 |14 |1997-09-24 03:47:45|
|79 |5 |3 |1997-09-30 16:11:35|
|79 |6 |3 |1997-12-31 20:39:07|
|79 |7 |14 |1998-03-31 13:51:01|
|79 |8 |14 |1997-11-12 19:18:06|
|79 |11 |4 |1998-04-14 23:36:23|
|79 |13 |3 |1997-12-14 22:49:06|
|79 |16 |5 |1997-10-24 21:05:22|
|79 |18 |4 |1997-11-21 16:57:30|
+-----+
```

only showing top 10 rows

Took 1 min 4 sec. Last updated by anonymous at June 19 2025, 9:34:47 PM.

Primary Key Design Explanation:

- Partition Key: `movie_id`, to facilitate aggregating ratings for each movie
- Clustering Key: `user_id`, to enable quick lookup of whether a user has rated a specific movie
- Using a compound primary key `(movie_id, user_id)` prevents users from rating the same movie multiple times

Timestamp Field Conversion:

- Use `from_unixtime()` to convert `int` to Spark-supported `timestmap`
- Automatically maps to `timestmap` type when writing to Cassandra, facilitating subsequent time analysis

Took 0 sec. Last updated by anonymous at June 19 2025, 9:33:38 PM.

Section 3 | Load Movie Information and Write to Cassandra Table `movies`

FINISHED

The file `u.item` contains detailed movie information, including:

- `movie_id` (Primary Key)
- `title` (Movie Title)
- 19 boolean flags indicating whether the movie belongs to corresponding genres (e.g., Action, Comedy)

To improve readability and facilitate subsequent analysis, we will convert the genres into a `set type`, storing multiple genres for each movie.

Field Definitions:

- `movie_id`: Type: `int`, Description: Primary key
- `title`: Type: `text`, Description: Movie title
- `genres`: Type: `set<text>`, Description: Set of genres (e.g., `{"Action", "Drama"}`)

Took 0 sec. Last updated by anonymous at June 19 2025, 9:33:38 PM.

```
%pyspark
# 3.1 加载 genre 映射表
genre_map = (
    spark.read.text("hdfs://user/maria_dev/ml-100k/u.genre")
    .rdd.map(lambda r: r[0].split("|"))
    .map(lambda kv: (int(kv[1]), kv[0]))
    .collectAsMap()
)
```

FINISHED

```
print(genre_map) # 可查看 genre 索引和名称映射关系
```

```
{0: 'Unknown', 1: 'Action', 2: 'Adventure', 3: 'Animation', 4: 'Children\'s', 5: 'Comedy', 6: 'Crime', 7: 'Documentary', 8: 'Drama', 9: 'Fantasy', 10: 'Film-Noir', 11: 'Horror', 12: 'Musical', 13: 'Mystery', 14: 'Romance', 15: 'Sci-Fi', 16: 'Thriller', 17: 'War', 18: 'Western'}
```

Took 2 sec. Last updated by anonymous at June 19 2025, 9:34:47 PM.

```
%pyspark
# 3.3 加载 u.item 文件并解析 genre 位
from pyspark.sql.types import *
from pyspark.sql import DataFrame, Row, Column
```

FINISHED

```
# 注意 schema (前 5 列 + 19 个 genre 布尔列)
predefined_schema = StructField("movie_id", IntegerType(), True),
StructField("title", StringType(), True),
StructField("genres", SetStringType(), True), ...
```

```

structField("release_date", StringType(), true),
StructField("video_release_date", StringType(), true),
StructField("imdb_url", StringType(), true)
]

schema_item = StructType(prefix_fields + genre_fields)

item_raw = spark.readCSV(
    "hdfs://user/maria_dev/ml-100k/u.item",
    sep="|",
    schema=schema_item
)

item_raw.show(5, truncate=False)

+-----+-----+-----+-----+-----+
|movie_id|title |release_date|video_release_date|imdb_ur l|
+-----+-----+-----+-----+-----+
|1       |Toy Story (1995)|01-Jan-1995| null      |http://us.imdb.com//title-exact?Toy%20Story%20(1995) |
|2       |GoldenEye (1995)|01-Jan-1995| null      |http://us.imdb.com//title-exact?GoldenEye%20(1995) |
|3       |Four Rooms (1995)|01-Jan-1995| null      |http://us.imdb.com//title-exact?Four%20Rooms%20(1995) |
|4       |Get Shorty (1995)|01-Jan-1995| null      |http://us.imdb.com//title-exact?Get%20Shorty%20(1995) |
|5       |Copycat (1995) |01-Jan-1995| null      |http://us.imdb.com//title-exact?Copycat%20(1995) |
+-----+-----+-----+-----+-----+
only showing top 5 rows

```

Took 0 sec. Last updated by anonymous at June 19 2025, 9:34:47 PM.

```

%spark
# 3.4 转换 genre 位 + set<text>
# 定义转换函数
def extract_genres(row):
    genres = [genre_map[i] for i in range(19) if row["g"+str(i)] == 1]
    return (row["movie_id"], row["title"], genres) # 不建议用 set, DataFrame不支持

# RDD 转换为 (movie_id, title, [genre1, genre2, ...])
movie_rdd = item_raw.rdd.map(extract_genres)

# 将转换为 DataFrame
movie_df = spark.createDataFrame(movie_rdd, ["movie_id", "title", "genres"])

# 展示前5行
movie_df.show(5, truncate=False)

+-----+-----+-----+
|movie_id|title |genres   |
+-----+-----+-----+
|1       |Toy Story (1995)|[Animation, Children's, Comedy] |
|2       |GoldenEye (1995)|[Action, Adventure, Thriller] |
|3       |Four Rooms (1995)|[Thriller] |
|4       |Get Shorty (1995)|[Action, Comedy, Drama] |
|5       |Copycat (1995) |[Crime, Drama, Thriller] |
+-----+-----+-----+
only showing top 5 rows

```

Took 1 sec. Last updated by anonymous at June 19 2025, 9:34:49 PM.

```

%spark
# 3.5 将 movie 保存入 cassandra 表 movies
movies_df.repartition(1).write \
    .format("org.apache.spark.sql.cassandra") \
    .option("table", "movies") \
    .option("keyspace", "movielens") \
    .mode("append") \
    .save()

```

Took 2 sec. Last updated by anonymous at June 19 2025, 9:34:50 PM.

```

%spark
# 3.6 验证写入是否成功
spark.read \
    .format("org.apache.spark.sql.cassandra") \
    .options(table="movies", keyspace="movielens") \
    .load().show(10, truncate=False)

+-----+-----+
|movie_id|genres   |
+-----+-----+
|1       |[Animation, Children's, Comedy] |
|2       |[Action, Adventure, Thriller] |
|3       |[Thriller] |
|4       |[Action, Comedy, Drama] |
|5       |[Crime, Drama, Thriller] |
+-----+-----+
only showing top 10 rows

```

Took 3 sec. Last updated by anonymous at June 19 2025, 9:34:51 PM.

• Each movie can correspond to multiple genres, thus a `set<text>` type is used for storage.
 • The `genres` field parsing is based on the `u.genre` mapping, avoiding direct use of boolean bit numbering.
 • Can be JOINed with the `ratings` table via `movie_id` for convenient analysis of rating performance across different genres.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:34:41 PM.

Section 4 | Answering Five Specific Query Questions (Assignment i–v)

In this section, we will use Spark SQL to answer the following five questions by querying the three previously written tables:

- The average rating for each movie
- The top 10 movies by average rating (including titles)
- The favorite movie genre of each active user
- All users under the age of 20
- Users aged between 30 and 40 whose occupation is “scientist”

To simplify the analysis, we first register the three tables as temporary views (TempView).

Took 0 sec. Last updated by anonymous at June 19 2025, 9:34:42 PM.

```

%spark
# 4.1 注册三张 Cassandra 表为视图
spark.read \
    .format("org.apache.spark.sql.cassandra") \
    .options(table="users", keyspace="movielens") \
    .load().createOrReplaceTempView("users")

spark.read \
    .format("org.apache.spark.sql.cassandra") \
    .options(table="ratings", keyspace="movielens") \
    .load().createOrReplaceTempView("ratings")

spark.read \
    .format("org.apache.spark.sql.cassandra") \
    .options(table="movies", keyspace="movielens") \
    .load().createOrReplaceTempView("movies")

```

Took 3 sec. Last updated by anonymous at June 19 2025, 9:34:54 PM.

(i) Average Rating for Each Movie

- This query uses the `ratings` table, grouping by `movie_id` and calculating `AVG(rating)`.
- The results can be used for subsequent filtering of highly-rated movies.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:33:43 PM.

```

%spark
# (1) 每部电影的平均评分
SELECT r.movie_id, m.title, ROUND(AVG(r.rating), 2) AS avg_rating
FROM ratings r
JOIN movies m ON r.movie_id = m.movie_id
GROUP BY r.movie_id

```

FINISHED

FINISHED

FINISHED

FINISHED

FINISHED

FINISHED

FINISHED

```
ORDER BY r.movie_id
""").show(10)
```

```
+-----+-----+
|movie_id|      title|avg_rating|
+-----+-----+
|    1| Toy Story (1995)|   3.88|
|    2| GoldenEye (1995)|   3.21|
|    3| Four Rooms (1995)|   3.03|
|    4| Gee Shorty (1995)|   3.55|
|    5| Copycat (1995)|   3.31|
|    6|Shanghai Tribe (Y...)|   3.58|
|    7|Twelve Monkeys (1...)|   3.8|
|    8| Babe (1995)|   4.01|
|    9|Dead Man Walking ...|   3.91|
|   10| Richard III (1995)|   3.83|
+-----+
only showing top 10 rows
```

Took 23 sec. Last updated by anonymous at June 19 2025, 9:35:14 PM.

FINISHED ▶ ✎ ↻ ⏹

(ii) Top Ten Movies with Highest Average Rating

- Join `ratings` and `movies` tables to get movie titles by movie ID.
- Sort by `avg_rating DESC` to select the top 10 movies with the highest ratings.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:41:07 PM (initiated)

FINISHED ▶ ✎ ↻ ⏹

```
%pyspark
# (ii) 平均评分最高的前十部影片 (含标题)
spark.sql("""
SELECT r.movie_id, m.title, ROUND(AVG(r.rating), 2) AS avg_rating
FROM ratings r
JOIN movies m ON r.movie_id = m.movie_id
GROUP BY r.movie_id, m.title
ORDER BY avg_rating DESC
""").show(n = 10, truncate = False)
```

```
+-----+-----+
|movie_id|title          |avg_rating|
+-----+-----+
|1201  |Marlene Dietrich: Shadow and Light (1996)| 5.0 |
|1122  |They Made Me a Criminal (1939)| 5.0 |
|1189  |Prefontaine (1997)| 5.0 |
|814   |Great Day in Harlem, A (1994)| 5.0 |
|1293  |Star Kid (1997)| 5.0 |
|1599  |Someone Else's America (1995)| 5.0 |
|1467  |Saint of Fort Washington, The (1993)| 5.0 |
|1536  |Alijing wansui (1994)| 5.0 |
|1653  |Entertaining Angels: The Dorothy Day Story (1996)| 5.0 |
|1580  |Santa with Muscles (1996)| 5.0 |
+-----+
only showing top 10 rows
```

Took 24 sec. Last updated by anonymous at June 19 2025, 9:35:19 PM.

FINISHED ▶ ✎ ↻ ⏹

(iii) For each user with at least 50 ratings, identify their favorite movie genre

- First, filter the "active users" group (those with ≥ 50 ratings).
- Then, extract the movie genres rated by these users through a JOIN operation.
- Use `explode()` to flatten the collection and count the most frequently occurring genres.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:33:48 PM.

FINISHED ▶ ✎ ↻ ⏹

```
%pyspark
# (iii) 每个评分≥50次的用户最喜欢的电影类型
spark.sql("""
WITH active_users AS (
    SELECT user_id
    FROM ratings
    GROUP BY user_id
    HAVING COUNT(*) >= 50
),
exploded_genres AS (
    SELECT user_id, genre
    FROM ratings r
    JOIN active_users a ON r.user_id = a.user_id
    JOIN movies m ON m.movie_id = r.movie_id
    LATERAL VIEW explode(m.genres) AS genre
),
genre_counts AS (
    SELECT user_id, genre, COUNT(*) AS count
    FROM exploded_genres
    GROUP BY user_id, genre
),
ranked_genres AS (
    SELECT *,
        ROW_NUMBER() OVER (PARTITION BY user_id ORDER BY count DESC) AS rank
    FROM genre_counts
)
SELECT g.user_id, u.gender, u.occupation, g.genre, g.count
FROM ranked_genres g
JOIN users u ON g.user_id = u.user_id
WHERE g.rank = 1
ORDER BY g.count DESC
""").show(10, truncate=False)
```

```
+-----+-----+-----+-----+
|user_id|gender|occupation|genre |count|
+-----+-----+-----+-----+
|655   |F     |healthcare|Drama |410 |
|405   |F     |healthcare|Drama |309 |
|537   |M     |engineer |Drama |251 |
|450   |F     |educator |Drama |237 |
|13    |M     |educator |Drama |218 |
|234   |M     |retired  |Drama |213 |
|416   |F     |student  |Drama |212 |
|279   |M     |programmer|Comedy|211 |
|281   |M     |writer   |Drama |196 |
|393   |M     |student  |Comedy|191 |
+-----+-----+-----+-----+
only showing top 10 rows
```

Took 20 sec. Last updated by anonymous at June 19 2025, 9:35:34 PM.

FINISHED ▶ ✎ ↻ ⏹

(iv) All users under the age of 20

- Simply filter the `users` table, selecting records where the age field is less than 20.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:39:32 PM.

FINISHED ▶ ✎ ↻ ⏹

```
%pyspark
spark.sql("""
SELECT * FROM users
WHERE age < 20
""").show(n = 10, truncate = False)
```

only showing top 10 rows

(v) Users aged between 30 and 40 whose occupation is "scientist"

- Use `BETWEEN` and `=` to combine filtering conditions.

Took 0 sec. Last updated by anonymous at June 19 2025, 9:34:50 PM.

```
%pyspark
# (v) 年龄在 30-40 岁之间且职业为 scientist 的用户
spark.read(*).show(n = 10, truncate=False)
```

```
+-----+
|user_id|age|gender|occupation|zip |
+-----+
|730 |31 |F |scientist |32114|
|74 |39 |M |scientist |18111|
|183 |33 |M |scientist |27708|
|107 |39 |M |scientist |60466|
|918 |40 |M |scientist |70116|
|337 |37 |M |scientist |10522|
|272 |33 |M |scientist |53706|
|430 |38 |M |scientist |98199|
|643 |39 |M |scientist |55122|
|543 |33 |M |scientist |95123|
+-----+
```

only showing top 10 rows

Took 1 sec. Last updated by anonymous at June 19 2025, 9:35:35 PM.

Section 5 | Summary and Future Work

Project Summary

In this Assignment 3, I have completed the following key tasks:

1. Used Apache Zeppelin to build an interactive Notebook, achieving an interpretable and reproducible data processing workflow.
2. Used PySpark to load three sub-files of the MovieLens dataset (`u.user`, `u.data`, `u.item`) from HDFS.
3. Combined Schema definition and cleaning logic to structure the data into DataFrames.
4. Utilized the Spark-Cassandra Connector to write three types of data (users, ratings, movies) into Cassandra tables (users, ratings, movies).
5. Designed reasonable primary keys and data types, particularly employing a composite primary key (`(movie_id, user_id)`) in `ratings`.
6. Implemented five core query tasks, covering single-table aggregation, JOINs, multi-level filtering, set unnesting, and other complex analytical operations.
7. Supplemented each analysis step with chart displays and Markdown documentation to ensure a clear and logically complete Notebook structure.

Course Knowledge Application Points

Concept Application Location

CAP Theory Selected Cassandra (AP model) for high-availability distributed writes
Hadoop HDFS Raw data loading paths based on HDFS storage structure
ETL Process Extract (HDFS) → Transform (Spark cleansing) → Load (Cassandra)
Distributed System Challenges Used <code>repartition(1)</code> to control concurrent writes, solving the risk of Zeppelin write failures
NoSQL Data Modeling Used <code>set<text></code> collection type to store movie genres, reasonably designed primary keys for optimized query performance

Took 0 sec. Last updated by anonymous at June 19 2025, 9:40:51 PM.

%pyspark