

Task3

```
1 //task3
2 // Write DataFrame-based Spark code to find the number of distinct movies in the file movies.csv
3 df1.select("title").distinct().count() //selecting the distinct movie title and counting them

+ (3) Spark jobs
res24: Long = 1499
Command took 2.05 seconds --- by sc2002@pghn.wiscouristate.edu at 10/29/2023, 12:40:40 PM on UW
Cmd 3
```

Task 4

```
1 //task4
2 //Write DataFrame-based Spark code to find the titles of the movies that appear in the file movies.csv but do not have a rating in the file movie_ratings.csv.
3 val df2join = df1.withColumnRenamed("title", "join_title") // renaming the title in movie_rating dataframe
4 df1.join(df2join, "title"=== "join_title", "left_anti").select("title").show(false) // left anti join on title column with movie and movie_rating dataframe, showing the title as output

+ (3) Spark jobs
+ df2join: org.apache.spark.sql.DataFrame = [rating: double, join_title: string ... 1 more field]
+-----+
+ [title]
+-----+
+-----+

df2join: org.apache.spark.sql.DataFrame = [rating: double, join_title: string ... 1 more field]
Command took 3.25 seconds --- by sc2002@pghn.wiscouristate.edu at 10/29/2023, 12:40:50 PM on UW
Cmd 4
```

Task 5

```
1 //task5
2 //Write DataFrame-based Spark code to find the number of movies that appear in the ratings file (i.e., movie_ratings.csv) but not in the movies file (i.e., movies.csv).
3 val df1join = df1.withColumnRenamed("title", "join_title") // renaming the title of the movie dataframe for join
4 df2.join(df1join, "title"=== "join_title", "left_anti").select("title").count() // left anti join on title column with movie_rating and movie dataframe, showing the distinct title count as output

+ (6) Spark jobs
+ df1join: org.apache.spark.sql.DataFrame = [actor: string, join_title: string ... 1 more field]
df1join: org.apache.spark.sql.DataFrame = [actor: string, join_title: string ... 1 more field]
res52: Long = 2327
Command took 2.77 seconds --- by sc2002@pghn.wiscouristate.edu at 10/29/2023, 12:50:30 PM on UW
```

Task 6

```
1 //task6
2 //Write DataFrame-based Spark code to find the total number of distinct movies that appear in either movies.csv, or movie_ratings.csv, or both.
3 df1.join(df2,eq("title"), "outer").select("title").distinct().count() // Outer joining the two dataframe based on similar title and showing the distinct title count as output

+ (3) Spark jobs
res33: Long = 3536
Command took 2.43 seconds --- by sc2002@pghn.wiscouristate.edu at 10/29/2023, 12:50:40 PM on UW
Cmd 6
```

Task 7

```
1 //task7
2 //Write DataFrame-based Spark code to find the title and year for movies that were remade. These movies appear more than once in the ratings file with the same title but different years. Sort the output by title.
3 val df2join=df2.groupBy("title").count().where("count">1).withColumnRenamed("title", "join_title") // Grouping the movie_rating dataframe with title and counting them. Filtering where count1 to get the remake, renaming the column title of the movie_rating dataframe for join
4 df1.join(df2join, "title"=== "join_title").select("title", "year").sort(desc("year")).show(false) // inner join with movie_rating dataframe, selecting year and title showing the output sorted descending order with title

+ (3) Spark jobs
+ df2join: org.apache.spark.sql.DataFrame = [join_title: string, count: long]
+-----+
+ [title] [year]
+-----+
+-----+
[A Nightmare on Elm Street][1984]
[A Nightmare on Elm Street][1984]
[Casino Royale [2006]
[Casino Royale [2007]
[Conan the Barbarian [1982]
[Conan the Barbarian [2012]
[Death at a Funeral [2008]
[Death at a Funeral [2007]
[Dracula [1992]
[Dracula [1979]
[Footloose [1984]
[Footloose [2011]
[Fright Night [2012]
[Fright Night [1985]
[Halifax [2007]
[Halifax [1988]
[Halifax [1978]
[Halifax [2007]
Command took 4.03 seconds --- by sc2002@pghn.wiscouristate.edu at 10/30/2023, 10:14:51 PM on UW
Cmd 7
```

Task 8

```
//task8
//Write DataFrame-based Spark code to find the rating for every movie that the actor "Brennon, Richard" appeared in. Schema of the output should be (title, year, rating)
val dfJoin = df1.withColumnRenamed("title", "join_title").withColumnRenamed("year", "join_year").select("join_title","actor", "join_year") //renaming the column title,year of the movie_rating dataframe for join
dfJoin(dfJoin, "title"="join_title") // joining the dataframes based on similar title
.where(col("actor")=equalTo("Brennon, Richard")) // filtering actor equal "Brennon, Richard"
.where(dfJoin.col("join_year")=equalTo(dfJoin.col("year")) // filtering the results actor did not worked
.select("title","year","rating").show(false) // selecting title, year, rating as output

// (2) Spark jobs
dfJoin: org.apache.spark.sql.DataFrame = [join_title: string, actor: string ... 1 more field]

+-----+-----+
|title      |year|rating|
+-----+-----+
|Casino Royale      |2006|9.2878|
|Around the World in 80 Days|2004|1.8623|
|Superman Returns    |2006|9.1888|
+-----+-----+

dfJoin: org.apache.spark.sql.DataFrame = [join_title: string, actor: string ... 1 more field]

Command took 3.08 seconds -- by sc2020light@mscsource1data-ssh at 10/20/2023, 12:05:04 PM on HW
```

Task 9

```
//task9
//Write DataFrame-based Spark code to find the highest-rated movie per year and include all the actors in that movie. The output should have only one movie per year, and it should contain four columns: year, movie title, rating, and a list of actor names. Sort the output by year.
val windowSpec = Window.partitionBy("year").orderBy(col("rating").desc) // creating window partition by year and descending order by rating to get the top rated movies
val rankedMovie = df1.withColumn("rank", rank().over(windowSpec)).where("rank" == 1) // using rank function and filtering based on rank = 1 to get the top movie of each year
rankedMovie.join(dfJoin, "title"="year", "year"="year", "left_actor") // left actor join two dataframes based on year and title
.drop(dfJoin.col("year")) // dropping duplicate column year
.groupBy("year", "title", "rating") // grouping the dataframe based on "year", "title", "rating"
.aggregate_list("actor", as="actor_list") // aggregating the actors as a list and named the column as actor_list
.select("year", "title", "rating", "actor_list").show(false) // selecting showing the "year", "title", "rating", "actor_list" as output

// (3) Spark jobs
rankedMovie: org.apache.spark.sql.DataFrame[org.apache.spark.sql.Row] = [rating: double, title: string ... 2 more field]

+-----+-----+-----+
|year|title      |rating|actor_list|
+-----+-----+-----+
|1937|Snow White and the Seven Dwarfs|2.2287|[ ]|
|1939|The Wizard of Oz      |7.9215|[ ]|
|1940|Pinocchio      |7.8557|[ ]|
|1942|Bambi      |1.5853|[ ]|
|1944|Song of the South      |7.462|[ ]|
|1950|Cinderella      |9.4226|[ ]|
|1953|Peter Pan      |5.4756|[ ]|
|1954|Bear Witness      |18.7625|[ ]|
|1955|Lady and the Tramp      |5.1258|[ ]|
|1956|Around the World in Eighty Days|14.9687|[ ]|
|1959|Sleeping Beauty      |6.3919|[ ]|
|1960|Psycho      |18.6375|[ ]|
|1961|One Hundred and One Dalmatians|18.6726|[Wright, Ben (I), Wickes, Mary]|
|1962|The Longest Day      |12.8866|[ ]|
|1963|It's a Mad Mad Mad Mad World      |6.426|[ ]|
|1964|My Fair Lady      |7.587|[ ]|
|1965|Doctor Zhivago      |4.9384|[ ]|
|1966|Who's Afraid of Virginia Woolf?|13.1311|[ ]|
+-----+-----+-----+

Command took 3.28 seconds -- by sc2020light@mscsource1data-ssh at 10/20/2023, 12:04:08 PM on HW
```

Task 10

```
//task10
//Write DataFrame-based Spark code to determine which pair of actors worked together most. Working together is defined as appearing in the same movie. The output should have three columns: actor 1, actor 2, and count. The output should be sorted by the count in descending order.
val windowSpec = Window.orderBy(col("actor").asc) // creating window ascending order by rating to get the top rated movies
val rankedActor = df1.withColumnRenamed("rank", "rank1").over(windowSpec2).withColumnRenamed("actor", "actor 1") // using rank function to rank all the actor and renaming the actor column as actor 1
val rankedMovie = rankedMovie.withColumnRenamed("rank", "rank2").withColumnRenamed("title", "movie_title").withColumnRenamed("year", "year_new") // creating another dataframe with renaming rank as rank2, title as movie_title, actor 1 as actor 2 and year as year_new
val rankedMovie2 = df1.withColumnRenamed("title", "movie_title").where(rankedMovie2("title")=rankedMovie("movie_title")).where(rankedMovie2("rank")=rankedMovie("rank2")) // joining the two dataframes based on similar title and actor 1 rank greater than actor 2 rank to avoid similar actor.
.groupBy("actor 1", "actor 2") // Group by the dataframe with actor 1 and actor 2
.count() // counting the movies based on actor pair and showing the output as descending order

// (3) Spark jobs
rankedMovie: org.apache.spark.sql.DataFrame = [actor 1: string, title: string ... 2 more field]
rankedMovie2: org.apache.spark.sql.DataFrame = [actor 2: string, movie_title: string ... 2 more field]

+-----+-----+
|actor 1|actor 2|count|
+-----+-----+
|McGowan, Mickie|Lynn, Sherry (I)|23|
|Lynn, Sherry (I)|Bergen, Bob (I)|19|
|McGowan, Mickie|Bergen, Bob (I)|19|
|Lynn, Sherry (I)|Angel, Jack (I)|17|
|McGowan, Mickie|Angel, Jack (I)|17|
|Rabson, Jan|McGowan, Mickie|16|
|Rabson, Jan|Lynn, Sherry (I)|16|
|McGowan, Mickie|Darling, Jennifer|15|
|Harnell, Jessi|Bergen, Bob (I)|14|
|Lynn, Sherry (I)|Darling, Jennifer|14|
|McGowan, Mickie|Harnell, Jessi|14|
|McGowan, Mickie|Farmer, Bill (I)|14|
|Rabson, Jan|Bergen, Bob (I)|14|
|Schneider, Bob (I)|Jennifer, Adam (I)|14|
|Bergen, Bob (I)|Angel, Jack (I)|13|
|Bumpass, Rodger|Bergen, Bob (I)|13|
|Lynn, Sherry (I)|Harnell, Jessi|13|
|Lynn, Sherry (I)|Farmer, Bill (I)|13|
+-----+-----+

Command took 4.08 seconds -- by sc2020light@mscsource1data-ssh at 10/20/2023, 12:05:09 PM on HW
```