

```
1 //Task 11
2 // SQL based spark code compute the number of movies produced in each year
3 // Used Group By year to count movies for each year ignored the null year by adding condition year is not null
4 // Used order by to sort the output in ascending order
5 val sqlWay = spark.sql("""
6 SELECT year,count(title) as count
7 FROM movies_table where year is not null
8 GROUP BY year ORDER BY year
9 """)
10 sqlWay.show()// showing the output year and count(number of movies)
```

▶ (2) Spark Jobs

▶  sqlWay: org.apache.spark.sql.DataFrame = [year: integer, count: long]

```
+-----+
|year|count|
+-----+
|1961| 2|
|1967| 2|
|1972| 12|
|1973| 5|
|1975| 5|
|1977| 40|
|1978| 30|
|1979| 37|
|1980| 47|
|1981| 53|
|1982| 103|
|1983| 119|
|1984| 149|
|1985| 133|
|1986| 174|
|1987| 126|
|1988| 111|
|1989| 152|
```

```

1 //task 15
2 // DataFrame based spark code to find the title and year for every movie that Tom Hanks acted in
3 val dataframeWay3=df
4 .where(col("actor").equalTo("Hanks, Tom")) // filtered actor column with actor equal "Hanks, Tom"
5 .select("year", "title").orderBy("year")// Select the year and title to show in output sorted by year
6
7 dataframeWay3.show(false)// showing the out put year and title show(false) to show the full title

```

► (1) Spark Jobs

► dataframeWay3: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [year: integer, title: string]

year	title
1984	Bachelor Party
1993	Sleepless in Seattle
1993	Philadelphia
1994	Forrest Gump
1995	Apollo 13
1995	Toy Story
1998	Saving Private Ryan
1998	You've Got Mail
1999	The Green Mile
1999	Toy Story 2
2000	Cast Away
2002	Catch Me If You Can
2002	Road to Perdition
2004	The Polar Express
2004	The Terminal
2004	The Ladykillers
2005	Magnificent Desolation: Walking on the Moon 3D

```

1 //task 14
2 // DataFrame based spark code find the five top most actors who acted in the most number of movies
3 val dataframeWay2 =df
4 .groupBy("actor")// Used Group By actor to count movies by each actor
5 .count()
6 .withColumnRenamed("count", "number_of_movies")// Renamed the count with number_of_movies
7 .sort(desc("number_of_movies"))// Used sort function with asc for sorting them in descending order based on the number of movies
8 .limit(5)// Used Limit 5 for showing the top 5 actors who acted in the most number of movies
9 dataframeWay2.show()// showing the output actor and number_of_movies

```

► (2) Spark Jobs

► dataframeWay2: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [actor: string, number_of_movies: long]

actor	number_of_movies
Tatasciore, Fred	38
Welker, Frank	38
Jackson, Samuel L.	32
Harnell, Jess	31
Damon, Matt	27

dataframeWay2: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [actor: string, number_of_movies: bigint]

Command took 1.51 seconds -- by sc26s@login.missouristate.edu at 10/12/2023, 4:53:04 PM on HW 2

```

1 //Task 13
2 // SQL based spark code find the five top most actors who acted in the most number of movies
3 // Used Group By actor to count movies by each actor
4 // Renamed the count(title) with number_of_movies
5 // Used Order by count(title) and DESC for sorting them in descending order based on the number of movies
6 // Used Limit 5 for showing the top 5 actors who acted in the most number of movies
7
8 val sqlWay2 = spark.sql("""
9 SELECT actor,count("title") as number_of_movies
10 FROM movies_table
11 GROUP BY actor ORDER BY count("title") DESC LIMIT 5
12 """)
13 sqlWay2.show()// showing the output actor and number_of_movies

```

▶ (2) Spark Jobs

▶  sqlWay2: org.apache.spark.sql.DataFrame = [actor: string, number_of_movies: long]

```

+-----+
|          actor|number_of_movies|
+-----+
|  Tatasciore, Fred|          38|
|    Welker, Frank|          38|
|Jackson, Samuel L.|          32|
|   Harnell, Jess|          31|
|    Damon, Matt|          27|
+-----+

```

sqlWay2: org.apache.spark.sql.DataFrame = [actor: string, number_of_movies: bigint]

Command took 1.93 seconds -- by sc26s@Login.missouristate.edu at 10/12/2023, 4:53:04 PM on HW 2

```

1 //task 12
2 // dataframe based spark code compute the number of movies produced in each year
3 // Used Group By year to count movies for each year ignored the null year by adding condition year is not null
4 // Used sort function with asc to show the output in ascending order
5 val dataframeWay =df.where("year is not null")
6 .groupBy("year")
7 .count()
8 .withColumnRenamed("count", "count")
9 .sort(asc("year"))
10 dataframeWay.show()// showing the output year and count(number of movies)

```

▶ (2) Spark Jobs

▼  dataframeWay: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row]

year: integer
count: long

```

+-----+
|year|count|
+-----+
|1961|    2|
|1967|    2|
|1972|   12|
|1973|    5|
|1975|    5|
|1977|   49|
|1978|   39|
|1979|   37|
|1980|   47|

```