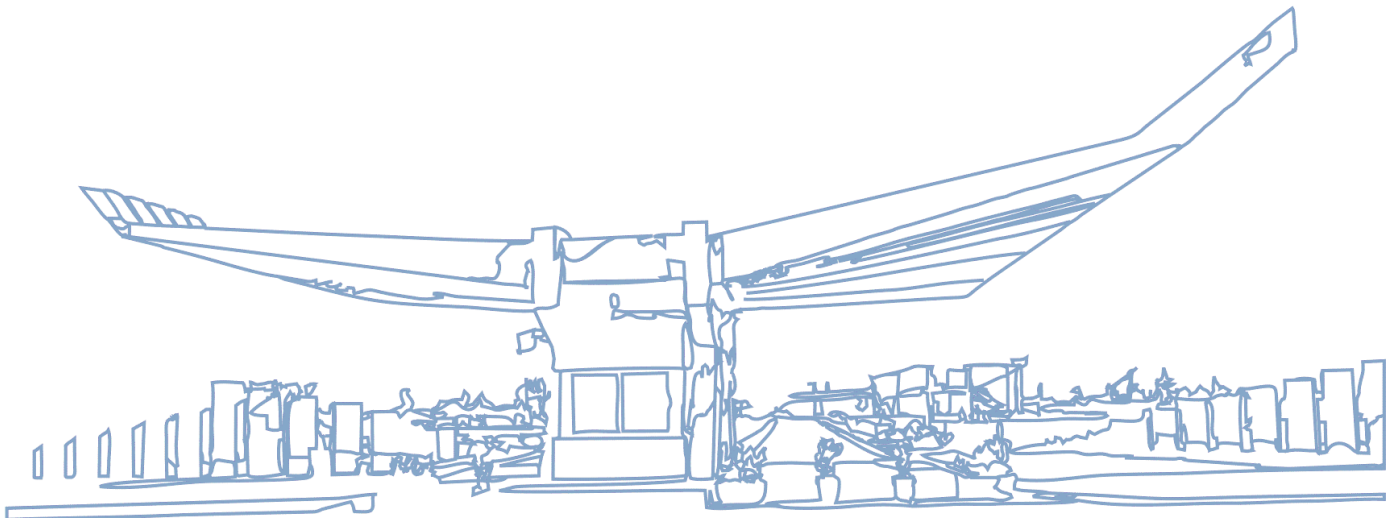# CEN 571 – Data Mining

## Data Mining with Spark

PREPARED:

**Baftjar TABAKU**

**21.6.2020**
Epoka University
Tirana, ALBANIA

ACCEPTED:

**Prof.Dr. Arben Asllani**

# 1. Dataset

The dataset was taken by Kaggle.com



With a size of 168 MB, the latest one, of 6 months, composed of 4 tables, by Stefano Leone. All in CSV format as shown below.

## 2. Cleaning the unwanted data, the data selection according to the project's goals.

Some redundant features will be removed, and data will be processed using Map Reduce, with corresponding code and jar files.

Removing some columns was used the Microsoft Excel, where the data was displayed better and modified.

Ex: According to my goals, I don't need a movie Description, cast list, reviews numbers from users and anything to do with the price, or the user's that rate professions, spouses' number and so on.

We also delete the cast data from dataset.



From all data of 4 tables, I reduced it to 2 and removed the unnecessary features for all of them.

# 3. Writing the Spark Application code in python

All the code is included as a single python file, included on the project submission folders.

```
"""
Main commands to proceed with, before the python

hadoop fs -ls
hadoop fs -mkdir /Spark_APP
hadoop fs -mkdir /Spark_APP/input
hadoop fs -copyFromLocal movies.csv /Spark_APP/input
hadoop fs -copyFromLocal ratings.csv /Spark_APP/input
hadoop fs -ls /Spark_APP/input

//run py Spark_APP
spark-submit IBDM_movies.py /Spark_APP/input/movies.csv
/Spark_APP/input/ratings.csv /Spark_APP
"""

"""
Author: Baftjar Tabaku
Epoka University
Data Mining
"""

import sys
from sched import scheduler

from pyspark.sql import SparkSession
from pyspark.sql.types import IntegerType
from pyspark.sql.types import FloatType

if __name__ == "__main__":
    if len(sys.argv) < 4:
        sys.stderr.write("Error: Usage: IBDM_movies.py <input-file 1> <input-
file 2> </Root of files>")
        sys.exit()

    spark = SparkSession.builder.getOrCreate()
    spark.sparkContext.setLogLevel("WARN")


    # Functions part to get the genders differences in vote for each movie
    def opinion_difference(num1, num2):
        return num1 - num2


    # Creating data frames
    # Movies table part, dealing with movies

    # IMDb_RDD_movies = spark.read.format("csv").option("header",
"true").load(
    #      "C:\\Users\\Baftjar
Tabaku\\PycharmProjects\\DataMining\\movies.csv")
    IMDb_RDD_movies = spark.read.format("csv").option("header",
"true").load(sys.argv[1])
```

```python
    # cast each variable to the current data type , except the strings
    IMDb_RDD_movies = IMDb_RDD_movies.withColumn("year",
IMDb_RDD_movies["year"].cast(IntegerType()))
    IMDb_RDD_movies = IMDb_RDD_movies.withColumn("duration",
IMDb_RDD_movies["duration"].cast(IntegerType()))
    IMDb_RDD_movies = IMDb_RDD_movies.withColumn("avg_vote",
IMDb_RDD_movies["avg_vote"].cast(FloatType()))
    IMDb_RDD_movies = IMDb_RDD_movies.withColumn("votes",
IMDb_RDD_movies["votes"].cast(IntegerType()))

    # Second RDD Ratings part
    IMDb_RDD_ratings = spark.read.format("csv").option("header",
"true").load(sys.argv[2])
    # Casting the variables
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("weighted_average_vote",

IMDb_RDD_ratings["weighted_average_vote"].cast(FloatType()))
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("total_votes",

IMDb_RDD_ratings["total_votes"].cast(IntegerType()))
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("mean_vote",
IMDb_RDD_ratings["mean_vote"].cast(FloatType()))
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("median_vote",

IMDb_RDD_ratings["median_vote"].cast(IntegerType()))
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("males_allages_avg_vote",

IMDb_RDD_ratings["males_allages_avg_vote"].cast(FloatType()))
    IMDb_RDD_ratings =
IMDb_RDD_ratings.withColumn("females_allages_avg_vote",

IMDb_RDD_ratings["females_allages_avg_vote"].cast(FloatType()))
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("top1000_voters_rating",

IMDb_RDD_ratings["top1000_voters_rating"].cast(FloatType()))

    # custom operations with functions
    IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("opinion_diff",

opinion_difference(IMDb_RDD_ratings["males_allages_avg_vote"],

IMDb_RDD_ratings[

"females_allages_avg_vote"]).cast(

                                        FloatType()))

    # IMDb_RDD_ratings = IMDb_RDD_ratings.withColumn("opinion_diff",
abs(["opinion_diff"]))

    # On the movies data we add a column, as the first task, the difference
in opinion

    IMDb_RDD_movies.printSchema()
    IMDb_RDD_movies.show()

    IMDb_RDD_ratings.printSchema()
    IMDb_RDD_ratings.show()

    # ---------------------------------------------------------------
    # Operations part, making them as tables for further operations
```

```python
# Creating two main tables
IMDb_RDD_movies.registerTempTable("IMDb_movies")
IMDb_RDD_ratings.registerTempTable("IMDb_ratings")

selected_all_movies = spark.sql("SELECT * FROM IMDb_movies")
selected_all_ratings = spark.sql("SELECT * FROM IMDb_ratings")
selected_all_movies.show(10)
selected_all_ratings.show(10)

# Sort all movies
selected_all_movies = spark.sql(
    "SELECT imdb_title_id, title,genre, year ,duration FROM IMDb_movies
SORT BY duration ASC")
# selected_all_movies.show(100)  # first 100 movies sorted to demonstrate
the query
selected_all_movies.show(10)

# counting movies of year 2010
count_2010_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2010 FROM
IMDb_movies where year=2010 GROUP BY year")
count_2010_movies.show()

# counting movies of year 2011
count_2011_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2011 FROM
IMDb_movies where year=2011 GROUP BY year")
count_2011_movies.show()

# counting movies of year 2012
count_2012_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2012 FROM
IMDb_movies where year=2012 GROUP BY year")
count_2012_movies.show()

# counting movies of year 2013
count_2013_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2013 FROM
IMDb_movies where year=2013 GROUP BY year")
count_2013_movies.show()

# counting movies of year 2014
count_2014_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2014 FROM
IMDb_movies where year=2014 GROUP BY year")
count_2014_movies.show()

# counting movies of year 2015
count_2015_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2015 FROM
IMDb_movies where year=2015 GROUP BY year")
count_2015_movies.show()

# counting movies of year 2016
count_2016_movies = spark.sql(
    "SELECT year , COUNT(imdb_title_id) as all_movies_2016 FROM
IMDb_movies where year=2016 GROUP BY year")
count_2016_movies.show()

# counting movies of year 2017
```

```python
    count_2017_movies = spark.sql(
        "SELECT year , COUNT(imdb_title_id) as all_movies_2017 FROM
IMDb_movies where year=2017 GROUP BY year")
    count_2017_movies.show()

    # counting movies of year 2018
    count_2018_movies = spark.sql(
        "SELECT year , COUNT(imdb_title_id) as all_movies_2018 FROM
IMDb_movies where year=2018 GROUP BY year")
    count_2018_movies.show()

    # Calculating the rating average of the movies for the years
2010,2011,2012,2013,2014,2015,2016, 2017,2018
    # 2010
    average_rating2010 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2010 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2010")
    average_rating2010.show()

    # 2011
    average_rating2011 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2011 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2011")
    average_rating2011.show()

    # 2012
    average_rating2012 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2012 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2012")
    average_rating2012.show()

    # 2013
    average_rating2013 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2013 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2013")
    average_rating2013.show()

    # 2014
    average_rating2014 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2014 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2014")
    average_rating2014.show()

    # 2015
    average_rating2015 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2015 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2015")
    average_rating2015.show()

    # 2016
    average_rating2016 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2016 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2016")
```

```
    average_rating2016.show()

    # 2017
    average_rating2017 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2017 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2017")
    average_rating2017.show()

    # 2018
    average_rating2018 = spark.sql(
        "SELECT AVG (weighted_average_vote) as total_avg_2018 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2018")
    average_rating2018.show()

    # Total Males and females ratting average for years 2010, 2011, 2012,
2013, 2014, 2015,2016,2017,2018
    # 2010
    # Females
    total_female_avg_rate_2010 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2010 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2010")
    total_female_avg_rate_2010.show()

    # Males
    total_male_avg_rate_2010 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2010 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2010")
    total_male_avg_rate_2010.show()

    # 2011
    # Females
    total_female_avg_rate_2011 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2011 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2011")
    total_female_avg_rate_2011.show()

    # Males
    total_male_avg_rate_2011 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2011 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2011")
    total_male_avg_rate_2011.show()

    # 2012
    # Females
    total_female_avg_rate_2012 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2012 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2012")
    total_female_avg_rate_2012.show()

    # Males
    total_male_avg_rate_2012 = spark.sql(
```

```
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2012 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2012")
    total_male_avg_rate_2012.show()

    # 2013
    # Females
    total_female_avg_rate_2013 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2013 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2013")
    total_female_avg_rate_2013.show()

    # Males
    total_male_avg_rate_2013 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2013 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2013")
    total_male_avg_rate_2013.show()

    # 2014
    # Females
    total_female_avg_rate_2014 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2014 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2014")
    total_female_avg_rate_2014.show()

    # Males
    total_male_avg_rate_2014 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2014 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2014")
    total_male_avg_rate_2014.show()

    # 2015
    # Females
    total_female_avg_rate_2015 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2015 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2015")
    total_female_avg_rate_2015.show()

    # Males
    total_male_avg_rate_2015 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2015 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2015")
    total_male_avg_rate_2015.show()

    # 2016
    # Females
    total_female_avg_rate_2016 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2016 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2016")
    total_female_avg_rate_2016.show()

    # Males
    total_male_avg_rate_2016 = spark.sql(
```

```python
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2016 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2016")
    total_male_avg_rate_2016.show()

    # 2017
    # Females
    total_female_avg_rate_2017 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2017 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2017")
    total_female_avg_rate_2017.show()

    # Males
    total_male_avg_rate_2017 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2017 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2017")
    total_male_avg_rate_2017.show()

    # 2018
    # Females
    total_female_avg_rate_2018 = spark.sql(
        "SELECT AVG(females_allages_avg_vote) as totalF_avg_2018 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2018")
    total_female_avg_rate_2018.show()

    # Males
    total_male_avg_rate_2018 = spark.sql(
        "SELECT AVG(males_allages_avg_vote) as totalM_avg_2018 FROM
IMDb_movies JOIN IMDb_ratings ON
IMDb_movies.imdb_title_id=IMDb_ratings.imdb_title_id WHERE year=2018")
    total_male_avg_rate_2018.show()

    # ============================Saving data===========================
    # TODO , finishing this part, saving the tables

    IMDb_RDD_movies.write.format('csv').option('header',
'true').save(sys.argv[3] + "/output_movies/")
    IMDb_RDD_ratings.write.format('csv').option('header',
'true').save(sys.argv[3] + "/output_ratings/")

    # udemy.write.format('csv').option('header',
'true').save("PROJ/OutputMain/")

    # simple selected and sorted
    selected_all_movies.coalesce(1).write.csv(sys.argv[3] + "/AllMovies")
    # Movies 2010-2018 counted
    count_2010_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2010Num/")
    count_2011_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2011Num/")
    count_2012_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2012Num/")
    count_2013_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2013Num/")
    count_2014_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2014Num/")
    count_2015_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2015Num/")
    count_2016_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2016Num/")
    count_2017_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2017Num/")
    count_2018_movies.coalesce(1).write.csv(sys.argv[3] + "/Movies2018Num/")

    # Rating average
```

```python
        average_rating2010.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2010/")
        average_rating2011.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2011/")
        average_rating2012.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2012/")
        average_rating2013.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2013/")
        average_rating2014.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2016/")
        average_rating2015.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2015/")
        average_rating2016.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2016/")
        average_rating2017.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2017/")
        average_rating2018.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2018/")

        # Males and females rating average for years 2010 - 2018
        # females
        total_female_avg_rate_2010.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2010/")
        total_female_avg_rate_2011.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2011/")
        total_female_avg_rate_2012.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2012/")
        total_female_avg_rate_2013.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2013/")
        total_female_avg_rate_2014.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2014/")
        total_female_avg_rate_2015.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2015/")
        total_female_avg_rate_2016.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2016/")
        total_female_avg_rate_2017.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2017/")
        total_female_avg_rate_2018.coalesce(1).write.csv(sys.argv[3] +
"/FemaleRatingAVG2018/")

        # Males
        total_male_avg_rate_2010.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2010/")
        total_male_avg_rate_2011.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2011/")
        total_male_avg_rate_2012.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2012/")
        total_male_avg_rate_2013.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2013/")
        total_male_avg_rate_2014.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2014/")
        total_male_avg_rate_2015.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2015/")
        total_male_avg_rate_2016.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2016/")
        total_male_avg_rate_2017.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2017/")
        total_male_avg_rate_2018.coalesce(1).write.csv(sys.argv[3] +
"/MaleRatingAVG2018/")

        # Movies 2010-2018 counted - END

        spark.stop()
```

# 4. Uploading the data into the cluster and processing calculations.

The data tables.csv were successfully added to the AWS cluster for further analyzing, using spark.

The output files and folders will be stored in the output folder, according to the ordering of the submission, below are the corresponding output screenshots.



Files successfully passed into the cluster

After execution, data files will be transferred again to local directory for further analyze,



# 5. Spark used commands and their output

I choose the python do these, I have included the commands in the code part, they are inside the python file, better than running them one by one.

Due to the fact that there are like 60.000 + records to millions, I decided to use the ".show(10)" records or each, then the results, with full selections are saved automatically in CSV files according to their records.



Tables and structures, the movies table,

The ratings table, there will be added also an 'opinion_diff' column which will be the total male and female opinion difference about a certain movie.



Simple select on movies,

An sorted selection, according to the duration of first 10 rows

Movies over the years,

```
+----+---------------+
|year|all_movies_2010|
+----+---------------+
|2010|           2253|
+----+---------------+


+----+---------------+
|year|all_movies_2011|
+----+---------------+
|2011|           2389|
+----+---------------+


+----+---------------+
|year|all_movies_2012|
+----+---------------+
|2012|           2517|
+----+---------------+


+----+---------------+
|year|all_movies_2013|
+----+---------------+
|2013|           2749|
+----+---------------+


+----+---------------+
|year|all_movies_2014|
+----+---------------+
|2014|           2851|
+----+---------------+
```

```
+----+---------------+
|year|all_movies_2017|
+----+---------------+
|2017|           3106|
+----+---------------+


+----+---------------+
|year|all_movies_2018|
+----+---------------+
|2018|           2880|
+----+---------------+
```

## Average rating for movies over the years 2010-2018

```
PuTTY (inactive)
+----------------+
|   total_avg_2010|
+----------------+
|5.650022191256351|
+----------------+


+----------------+
|   total_avg_2011|
+----------------+
|5.63784009545056|
+----------------+


+------------------+
|    total_avg_2012|
+------------------+
|5.6447755249900196|
+------------------+


+----------------+
|    total_avg_2013|
+----------------+
|5.641360492132152|
+----------------+


+----------------+
|    total_avg_2014|
+----------------+
|5.687022096773735|
+----------------+


+----------------+
|  total_avg_2015|
+----------------+
|5.63785738851381|
+----------------+


+----------------+
|    total_avg_2016|
+----------------+
|5.698186615061296|
+----------------+
```

```
+----------------+
|    total_avg_2017|
+----------------+
|5.759787507580698|
+----------------+


+----------------+
|    total_avg_2018|
+----------------+
|5.771666662684745|
+----------------+
```

Average ratings for booth males and females, for years 2010 – 2018,

```
+------------------+          totalF_avg_2011|          +------------------+
|    totalF_avg_2010|        ------------------+        |    totalF_avg_2012|
+------------------+          5.772331521270564|        +------------------+
|5.8027987567544885|        ------------------+         |5.794874853873755|
+------------------+                                    +------------------+


+----------------+            totalM_avg_2011|          +------------------+
|totalM_avg_2010|           ------------------+         |    totalM_avg_2012|
+----------------+           5.563164499625345|         +------------------+
|5.5832223673688|          ------------------+          |5.555462851427735|
+----------------+                                      +------------------+
```

```
+------------------+          +------------------+          +------------------+
|    totalF_avg_2013|         |    totalF_avg_2014|         |    totalF_avg_2015|
+------------------+          +------------------+          +------------------+
|5.785334791115104|          |5.849684217268961|          |5.812512931023249|
+------------------+          +------------------+          +------------------+


+------------------+          +------------------+          +------------------+
|    totalM_avg_2013|         |    totalM_avg_2014|         |    totalM_avg_2015|
+------------------+          +------------------+          +------------------+
|5.548235721005314|          |5.600876884851821|          |5.550017222715911|
+------------------+          +------------------+          +------------------+
```

```
+------------------+          +------------------+          +------------------+
|    totalF_avg_2016|         |    totalF_avg_2017|         |    totalF_avg_2018|
+------------------+          +------------------+          +------------------+
|5.915031384389945|          |5.991266520300719|          |6.0067826115981395|
+------------------+          +------------------+          +------------------+


+------------------+          +------------------+          +------------------+
|    totalM_avg_2016|         |    totalM_avg_2017|         |    totalM_avg_2018|
+------------------+          +------------------+          +------------------+
|5.588789978162549|          |5.614745652414642|          |5.5821180565903585|
+------------------+          +------------------+          +------------------+
```

Some data folders,

```
Traceback (most recent call last):
  File "/home/hadoop/IBDM_movies.py", line 312, in <module>
    average_rating2016.coalesce(1).write.csv(sys.argv[3] + "/RatingAVG2016/")
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/readwriter.py", line 932, in csv
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 1257, in __call__
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/utils.py", line 69, in deco
pyspark.sql.utils.AnalysisException: 'path hdfs://ip-172-31-66-54.ec2.internal:8020/SPRK_root/RatingAVG2016 already exists.;'
[hadoop@ip-172-31-66-54 ~]$ hadoop fs -ls /SPRK_root
Found 19 items
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/AllMovies
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2010Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2011Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2012Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2013Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2014Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2015Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2016Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2017Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/Movies2018Num
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2010
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2011
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2012
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2013
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2015
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/RatingAVG2016
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:38 /SPRK_root/input
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/output_movies
drwxr-xr-x   - hadoop hadoop          0 2020-06-20 13:39 /SPRK_root/output_ratings
[hadoop@ip-172-31-66-54 ~]$ hadoop fs -copyToLocal /SPRK_root
[hadoop@ip-172-31-66-54 ~]$
```