Guide to run the project:

- 1- Clone the repository.
- 2- Unzip the folder found in the git repository and place it under this path: 'C:/tmp/'
- 3- Open a spark session
- 4- Run the project

Code Walkthrough:

At first we use SparkSession to create our spark session or get a reference to an already created one.

Then we define the schema of the csv file.

Then we setup our spark session to watch for file changes in a specified path with a specific csv schema that we already created.

Then we create a streaming DataFrame. After that we setup our sql queries and then start the queries and instruct spark to show changes in the console.

Example of results:

- I- Video Games CSV:
- 1- First Query (global 2006 nintendo):

```
#display global_sales in 2006 where publisher is nintendo global_2006_nintendo = spark.sql("SELECT Global_Sales FROM VideoGames where Year=2006 and Publisher like 'Nintendo'")
```

Sélection Anaconda Prompt (py36) - jupyter notebook

2- Second Query (Total Global Sales):

```
#display total Global Sales for each year
Total_Global_Sales = spark.sql("SELECT Year, SUM(Global_Sales) FROM VideoGames GROUP BY Year")
                                                                                                                                              Sélection Anaconda Prompt (py36) - jupyter notebook
Batch: 0
  Year | sum(Global Sales)|
|1987.0|21.739999999999995
|2010.0| 600.2899999999948
1993.0
| 1992.0 | 260.38000000000013 | 1992.0 | 76.1599999999998 | 1996.0 | 49.38999999999999 | 1995.0 | 88.10999999999991 | 2009.0 | 667.2999999999947 | 2007.0 | 609.9199999999935
1996.0 199.14999999999995
2020.0
                          37.07
1986.0
 1998.0 256.46999999999963
| 1985.0 | 53.9400000000000005
2017.0
                          0 05
1982.0 28.85999999999999
 only showing top 20 rows
                                                                                  (125 + 8) / 200][I 23:25:33.680 NotebookApp]
```

3- Third Query (Avg_Global_sales):

#Display the average of global sales for each pair of plateform and video game name

Avg_Global_sales=spark.sql("SELECT Platform, Name ,MEAN(Global_Sales) FROM VideoGames GROUP BY Platform, Name")

Sélection Anaconda Prompt (py36) - jupyter notebook

```
Batch: 0
            Name avg(Global_Sales)
Platform
     PSP|Tiger Woods PGA T...|
                                          0.2
     Wii | I Am In The Movie
                                         0.08
    X360 Transformers: Fal...
                                         0.44
                                          0.41
     Wii | Monotaro Dentetsu... |
     PSP|Yu-Gi-Oh! GX: Tag...|
                                          0.14
     N64
              Donkey Kong 64
                                          5.27
     PS2|Jikkyou Powerful ...|
                                          0.35
     PSP | Pop'n Music Portable |
                                          0.13
     Wii Resident Evil: Th...
                                          1.08
      XB Pro Cast Sports F...
                                          0.03
     PS3 | Hail to the Chimp|
                                          0.05
                                          5.21
     PS2|Dragon Quest VIII...|
    X0ne
              Sniper Elite 3
                                          0.33
     PS3 Dynasty Warriors ...
                                          0.09
      DS Yu-Gi-Oh! GX: Spi...
                                          0.22
     PS3 Game of Thrones (...
                                          0.06
     GBA Eyeshield 21: Dev...
                                          0.03
      DS
             Mahjong Taikai
                                          0.04
     GBA MX 2002 Featuring...
                                          0.16
     3DS | Adventure Time: T...|
                                          0.02
only showing top 20 rows
```

I- Cars CSV:

1- First Query (top_fuel_eff_cars):

```
#Print Top Fuel Efficient cars
top_fuel_eff_cars = spark.sql("SELECT Model, AVG(Fuel_efficiency) FROM CarsTable group by Model ORDER BY AVG(Fuel_efficiency)
DESC")
```

Anaconda Prompt (py36) - jupyter notebook

Model	avg(Fuel_efficiency)
Metro	45.0
SC	33.0
SL	33.0
Corolla	33.0
Prizm	33.0
Civic	32.0
SW	31.0
Celica	31.0
Accent	31.0
Sentra	30.0
Escort	30.0
Cougar	30.0
Mirage	30.0
Neon	29.0
/ystique	28.0
Integra	28.0
Cirrus	27.0
Camry	27.0
Sunfire	27.0
Accord	27.0

2- Second Query (top_fuel_eff_manu):

#Print Top top manufacturer that have the best Fuel_efficiency average
top_fuel_eff_manu = spark.sql("SELECT Manufacturer, AVG(Fuel_efficiency) FROM CarsTable group by Manufacturer ORDER BY AVG(Fuel_efficiency) DESC")

Anaconda Prompt (py36) - jupyter notebook

```
Batch: 0
|Manufacturer|avg(Fuel_efficiency)|
      Saturn
               32.333333333333336
   Chevrolet
                          28.625
     Hyundai
               27.6666666666668
    Plymouth|
               26.666666666668
  Volkswagen
                            25.625
      Toyota
                             25.2
25.0
25.0
25.0
25.0
24.8
    Pontiac
Infiniti
       Honda
       Acura
    Chrysler
                              24.5
         BMW
       Buick
                             24.25
                             24.0
      Nissan
     Mercury
               23.66666666666668
        Audi
  Mercedes-B
                              23.0
  Mitsubishi
               22.857142857142858
               22.6666666666668
        Ford
nly showing top 20 rows
```

3- Third Query (Total Sales):

#Print Total Sales in Thousands grouped by Manufacturer and order by TotalSales

Total_Sales = spark.sql("SELECT Manufacturer, SUM(Sales_in_thousands) as TotalSales FROM CarsTable GROUP BY Manufacturer order
by TotalSales DESC")

Anaconda Prompt (py36) - jupyter notebook

```
Manufacturer
                     TotalSales
        Ford 1846.9650000000001
       Dodge
                        720.798
      Toyota
                        675.086
                        592.674
       Honda
   Chevrolet
                         446.37
     Pontiac|
                        330.962
        Jeep
                        293.153
                        280.472
      Nissan
       Buick
                         242.019
                         237.999
     Mercury
  Mitsubishi 180.895000000000004
  Volkswagen
                        159.749
     Hyundai
                        137.326
    Chrysler
                        117.545
      Saturn
                        110.389
    Cadillac
                          81.45
  Mercedes-B 66.07900000000001
       Acura 64.890999999999999
     Lincoln
                         62.709
    Plymouth 62.1290000000000005
only showing top 20 rows
```

4- Forth Query (Audi_cars):

5-

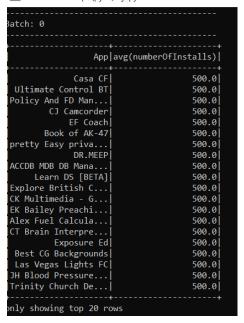
#Audi cars ordered by Year
Audi_cars=spark.sql("SELECT Manufacturer, Model ,SUBSTRING(Latest_Launch, length(Latest_Launch)-3 ,4) AS Year FROM CarsTable w
here Manufacturer like 'Audi' ")

II- Google Play CSV

1- First Query (top_installs):

#Top installs
top_installs = spark.sql("SELECT App,AVG(numberOfInstalls) FROM Googleplay group by App order by AVG(numberOfInstalls) DESC ")

Anaconda Prompt (py36) - jupyter notebook



2- Second Query (top_rate_medapp):

#top Rating apps where genre = medical
top_rate_medapp = spark.sql("SELECT App , AVG(Rating) from Googleplay where Genres like 'Medical' group by App order by AVG(Rating) DESC")

Anaconda Prompt (py36) - jupyter notebook

```
App|avg(Rating)|
           Sway Medical
     You're an Anime
Arrowhead AH App
                                          5.0
5.0
          FoothillsVet
Basics of Orthopa...
KBA-EZ Health Guide
                                          5.0
5.0
5.0
5.0
5.0
5.0
5.0
                Zen Leaf
    Clinic Doctor EHr
FHR 5-Tier 2.0
Super Hearing Sec...
CARDIAC CT TECHNIQUE
BP Journal - Bloo...
                                          5.0
5.0
5.0
5.0
          PrimeDelivery
        Labs on Demand
   CT Cervical Spine
NCLEX Multi-topic...
                                          5.0
Chenoweth AH
Cy-Fair VFD EMS P...
                                          5.0
Dermatology Atlas...
                                          5.0
     Galaxies of Hope
                                          5.01
nly showing top 20 rows
```

3- Third Query (top priced app):

```
#top priced apps
top_priced_app = spark.sql("SELECT App ,AVG(price) AS price from Googleplay group by App order by price DESC")
```

Anaconda Prompt (py36) - jupyter notebook

```
Batch: 0
                                        price
I'm Rich - Trump ...
                                        400.0
  I am Rich Plus 399.989990234375 I AM RICH PRO PLUS 399.989990234375
   I Am Rich Premium 399.989990234375
most expensive ap...|399.989990234375|
         I Am Rich Pro 399.989990234375
  I am Rich|399.989990234375|
I am Rich!|399.989990234375|
I am Rich!|399.989990234375|
            ? I'm rich 399.989990234375
I am rich (Most e...|399.989990234375|
I am rich|399.989990234375|
I Am Rich|389.989990234375|
 I am extremely Rich 379.989990234375
         I am rich VIP 299.989990234375
Chrome Canary (Un...
                                           0.0
free video calls ...
                                           0.0
Google Chrome: Fa...
                                           0.0
Mercari: The Sell...
                                           0.01
THE KING OF FIGHT...
                                           0.0
only showing top 20 rows
```

4- Forth Query (top_rated_cat):

#Top rated categories on Google play
top_rated_cat=spark.sql("SELECT Genres,AVG(Rating) as average_rating from Googleplay Group BY Genres order by average_rating DE
SC")

Anaconda Prompt (py36) - jupyter notebook

```
atch: 0
                          Genres average_rating
Comics;Creativity 4.800000190734863|
Board;Pretend Play 4.800000190734863|
Health & Fitness;... 4.699999809265137|
Adventure;Brain G... 4.599999904632568|
Strategy;Action &... 4.599999904632568|
Puzzle;Education 4.599999904632568|
Entertainment (Co. 4.6932320015441805)
 Entertainment;Cre...|4.5333333015441895
Music;Music & Video|4.5333333015441895
        Tools;Education
                                                                    4.5
  Arcade; Pretend Play
                                                                     4.5
  Racing; Pretend Play
                                                                     4.5
   Strategy; Education
    Casual;Brain Games 4.4692307985745945
                         Events 4.4355555640326605
Education;Brain G... | 4.425000071525574
Adventure;Action ... | 4.4230768863971415
Simulation;Action...| 4.418181766163219
Word| 4.410714294229235
Puzzle;Creativity| 4.400000095367432
       Card; Brain Games | 4.400000095367432|
only showing top 20 rows
```

5- Fifth Query (top inst arc):

```
#top installs where genre = Arcade
top_inst_arc=spark.sql("SELECT App, SUM(numberOfInstalls) from Googleplay where Genres like 'Arcade' Group BY App order by SUM
(numberOfInstalls) DESC limit 5")
```

```
Batch: 0

App|sum(numberOfInstalls)|

Mad Dash Fo' Cash| 100.0|

BL!TZ - Endless| 100.0|

B-52 Spirits of G...| 100.0|

Flippy Axe : Flip...| 100.0|

Galaxian(FC)| 100.0|
```

After your **feedback** we tried the following in Spark to add a column for the current timestamp so we can apply some window operations later:

```
#Adding a column for the current timestamp so we can apply some window opperations from pyspark.sql import functions as F dfCSV-dfCSV.withColumn('Time', F.current_timestamp())
```

As you can see in the following result 'Time' column is returned:

```
atch: 0
                                                                                                  Publisher|NA_Sales|EU_Sales|JP_Sales|Other_Sales|Global_Sales|
                             Name|Platform| Year|
                                                                        Genrel
                                             Wii|2006.0
                                                                                                                                                                                        82.74 2021-01-28 14:38:..
                    Wii Sports
                                                                       Sports
                                                                                                    Nintendo
                                                                                                                                                                                       40.24 | 2021-01-28 | 14:38:...
35.82 | 2021-01-28 | 14:38:...
33.0 | 2021-01-28 | 14:38:...
                                             Wii|2008.0
                                                                       Racing
                                                                                                    Nintendo
                                                                                                                       15.85
                                                                                                                                    12.88
                                              GB|1996.0|Role-Playing|
GB|1989.0| Puzzle|
                                                                                                                                     8.89
2.26
                                                                                                                                                                                       31.37 2021-01-28 14:38:..
30.26 2021-01-28 14:38:..
                                                                                                    Nintendo
                                                                                                                       11.27
                                                                                                    Nintendo
                                                                                                                                                   4.22
                                                                                                                        23.2
            Super Mario B...
Wii Play
                                             DS 2006.0
Wii 2006.0
                                                                     Platform
                                                                                                    Nintendo
Nintendo
                                                                                                                       11.38
14.03
                                                                                                                                     9.23
9.2
                                                                                                                                                                     2.9
2.85
                                                                                                                                                                                       30.01 2021-01-28 14:38:.
29.02 2021-01-28 14:38:.
                                                                         Misc
                    Mario B...
Duck Hunt
                                             Wii|2009.0
NES|1984.0
                                                                                                    Nintendo
Nintendo
                                                                                                                       14.59
26.93
                                                                                                                                     7.06
0.63
                                                                                                                                                   4.7
0.28
                                                                                                                                                                     2.26
0.47
                                                                                                                                                                                        28.62 2021-01-28 14:38:.
28.31 2021-01-28 14:38:.
                                                                      Shooter
                                                                                                                                                                                       24.76 | 2021-01-28 | 14:38:..
23.42 | 2021-01-28 | 14:38:..
23.1 | 2021-01-28 | 14:38:..
22.72 | 2021-01-28 | 14:38:..
                Nintendogs
Mario Kart DS
                                               DS 2005.0
DS 2005.0
                                                                                                                                     11.0
7.57
                                                                                                                                                   1.93
4.13
                                                                                                                                                                     2.75
                                                                       Racing
                                                                                                    Nintendo
                                                                                                                        9.81
                                                                                                                                     6.18
8.03
8.59
4.94
      Pokemon Gold/Poke...
Wii Fit
                                             GB|1999.0|Role-Playing
Wii|2007.0| Sports
                                                                                                    Nintendo
                                                                                                                                                   2.53
0.24
0.97
         Kinect Adventures!
                                             X360 2010.0
                                                                          Misc Microsoft Game St...
                                                                                                                                                                      1.67
                                                                                                                                                                                        21.82 2021-01-28 14:38:.
  18 Grand Theft Auto:...
19 Super Mario World
                                                                        Action Take-Two Interactive
                                                                                                                                                   0.41
3.54
                                                                                                                                                                    10.57
0.55
                                                                                                                                                                                       20.81 2021-01-28 14:38:..
20.61 2021-01-28 14:38:..
                                             PS2 2004 0
                                                                                                                        9.43
                                                                                                                                       0.4
  20|Brain Age: Train ...
                                                                                                    Nintendo
                                                                                                                                                                                        20.22 2021-01-28 14:38:.
nly showing top 20 rows
```

We then tried to apply a window operation:

And got this error:

We believe that this error happens due to spark version mismatch. We also tried the following instead with no luck:

```
from pyspark.sql import Window
window=Window.partitionBy(dfCSV["Publisher"]).orderBy(dfCSV["Year"].desc())
t=spark.sql("Select Name, Year, FIRST VALUE(Name) over (window) from VideoGames")
```

This following is our spark version: