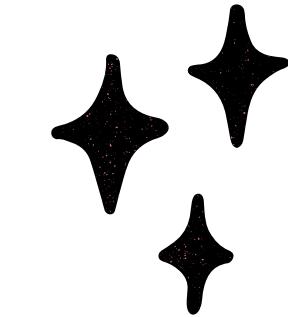


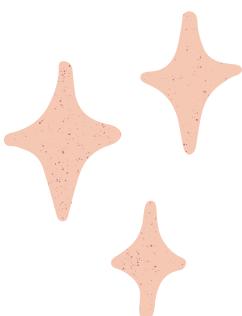
Analyzing Saudi Stock Exchange (Tadawul) with Hadoop

DaQuest Team



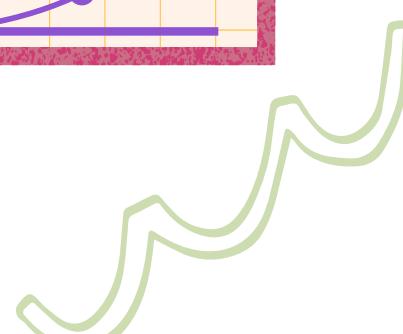
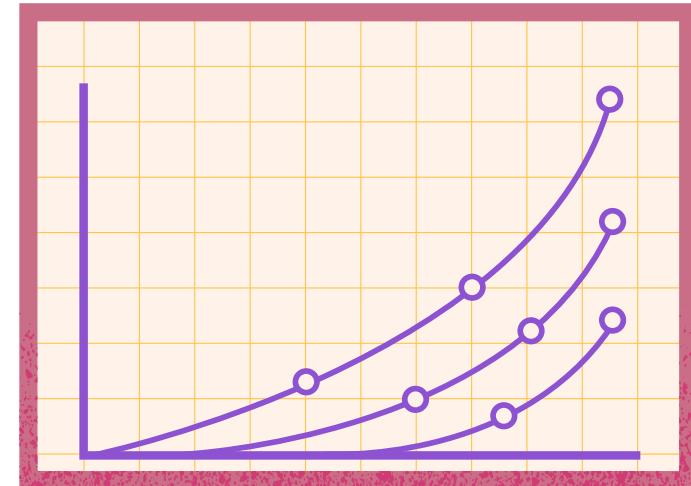
Today's Agenda

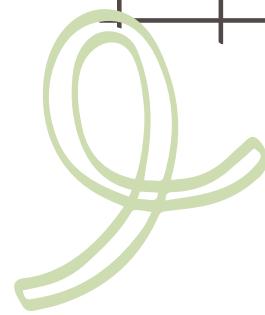
- 01** Introduction
- 02** Data Review
- 03** Data Preprocessing
- 04** Our Approach: Using Different MapReduce Functions
- 05** Results and What Can Be Improved



Introduction

- Saudi's 2030 Vision and Tadawul
- Business problem





Tadawal Dataset:

- This is the data of Saudi stock market companies from 2001-12-31 until 2020-04-16.
- Collected from Saudi Stock Exchange (Tadawul) website.
- This dataset has 14 columns and 593819 rows.



Tadawal Dataset:

Each row in the database represents the price of a specific stock at a specific date:

| | |
|---------------------------------|--|
| (symbol (Integer) | The symbol or the reference number of the company. |
| (name(String) | Name of the company. |
| (trading_name (String) | The trading name of the company |
| (sector (String) | The sector in which the company operates. |
| (date (Date) | The date of the stock price. |
| (open (Decimal) | The opening price. |
| (high (Decimal) | The highest price of the stock at that day. |
| (low (Decimal) | The lowest price of the stock at that day. |
| (close (Decimal) | The closing price. |
| (change (Decimal) | The change in price from the last day. |
| (perc_Change (Decimal) | The percentage of the change. |
| (volume_traded (Decimal) | The volume of the trades for the day. |
| (value_traded (Decimal) | The value of the trades for the day. |
| (no_trades (Decimal) | The number of trades for the day. |



Data preprocessing :

- Change the column Name
- Add year, month and day columns
- Drop unnecessary columns:

1-Symbol

2-Company Name

- Add the categorical target column

We added a new column from the per_change column in the dataset, which is the Change_category (the classification target).

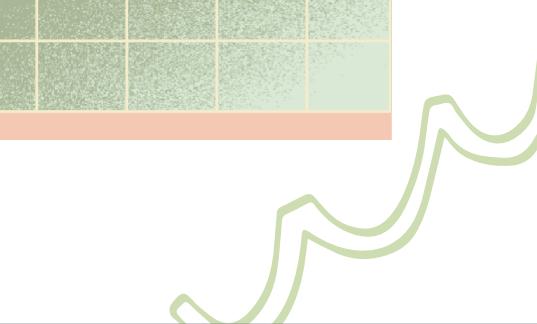
We create the classes from the per_change as below:

Good change: if the per_change > 0

Bad Change: if the per_change < 0

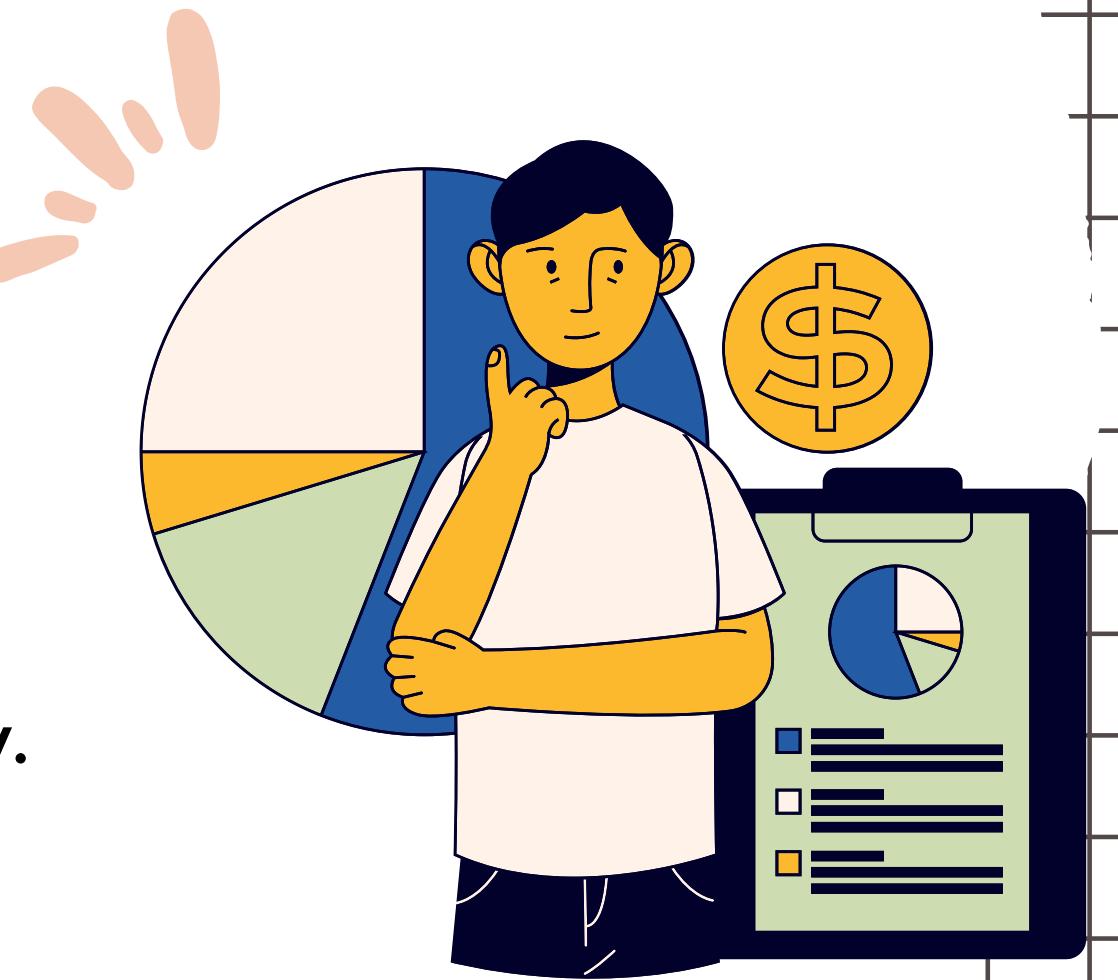
Stable: if the per_change = 0

- Handle the null values



Business Problem:

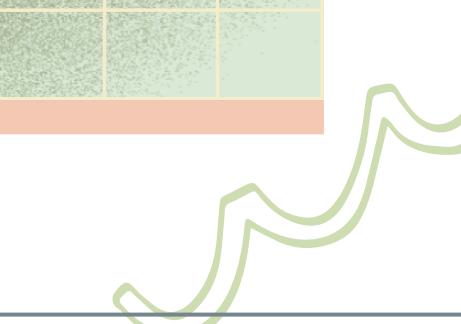
- Count change category for each sector.
- Count number of companies in each sector.
- Count the change categories for each year, month, and day.
- Count the number of shares in each year.



First Business Problem:

1-Count change category for each sector:

Our dataset contains a variety of sectors, it is critical to assess each one's performance using the change categories. We achieve this by grouping each sector's good, stable, and bad changes.



our approach



MapReduce Function:

```
from mrjob.job import MRJob # import the mrjob library
from mrjob.step import MRStep # import the mrStep library

class MRtarget(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_target,
                   reducer=self.reducer_count_target)
        ]

    # the map step: each line in the csv file is read as a key, value pair

    def mapper_get_target(self, _, target):
        yield target, 1

    # the reduce step: combine all tuples with the same key
    # then sum all the values of the tuple, which will give the total value

    def reducer_count_target(self, key, values):
        yield key, sum(values)

if __name__ == '__main__':
    MRtarget.run()
```

- Running Hadoop Cluster on a Sector_change data:

```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar data._Sector_change.csv
```



Output:

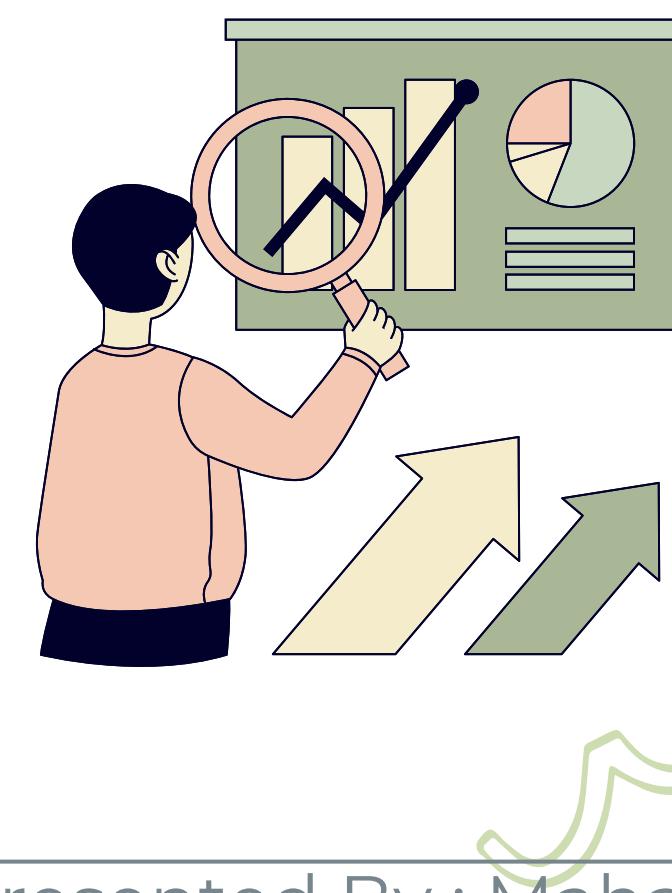
```
aria_dev.20221129.173214.064348/output...
"Communication Services,Bad Change" 10023
"Communication Services,Good Change" 9068
"Communication Services,Stable" 2257
"Consumer Discretionary,Bad Change" 24937
"Consumer Discretionary,Good Change" 23996
"Consumer Discretionary,Stable" 5923
"Consumer Staples,Bad Change" 28450
"Consumer Staples,Good Change" 27502
"Consumer Staples,Stable" 8734
"Energy,Bad Change" 7202
"Energy,Good Change" 6904
"Energy,Stable" 1625
"Financials,Bad Change" 67244
"Financials,Good Change" 63151
"Financials,Stable" 14455
"Health Care,Bad Change" 8259
"Health Care,Good Change" 8018
"Health Care,Stable" 1548
"IT,Bad Change" 394
"IT,Good Change" 360
"IT,Stable" 52
"Industrials,Bad Change" 28957
"Industrials,Good Change" 27568
"Industrials,Stable" 7412
"Materials,Bad Change" 63853
"Materials,Good Change" 61414
"Materials,Stable" 16277
"Real Estate,Bad Change" 19975
"Real Estate,Good Change" 18838
"Real Estate,Stable" 5650
"Utilities,Bad Change" 3804
"Utilities,Good Change" 3816
"Utilities,Stable" 1765
```



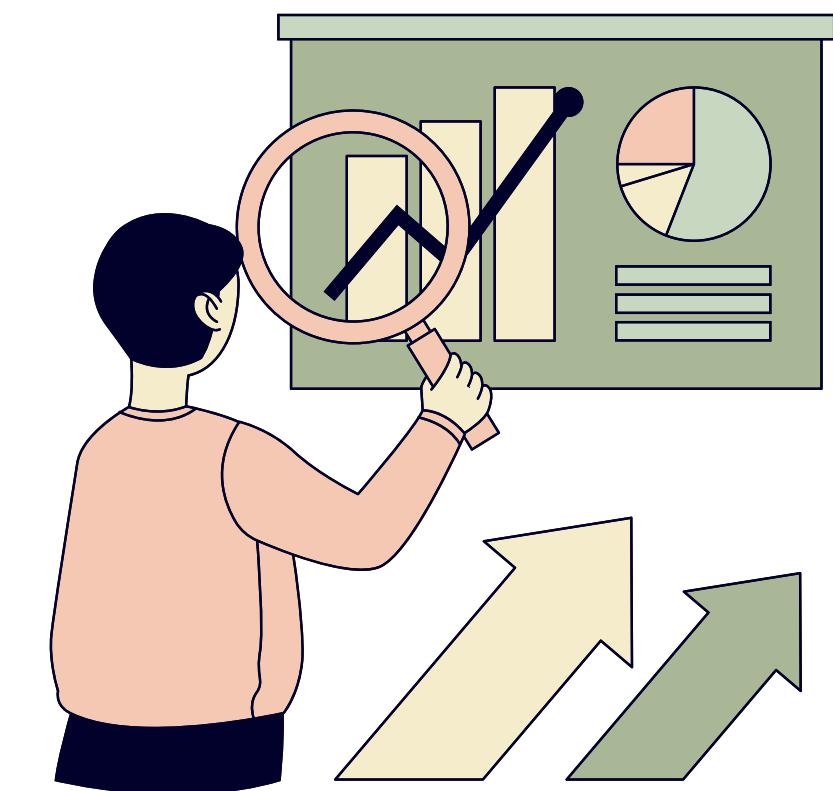
Second Business Problem:

2-Count number of companies in each sector:

We want to try another MapReduce problem that can help Tadawul to understand the Saudi Stock market better by knowing the number of trading companies in each sector to find the most popular sectors and the least popular ones.



Second Approach:

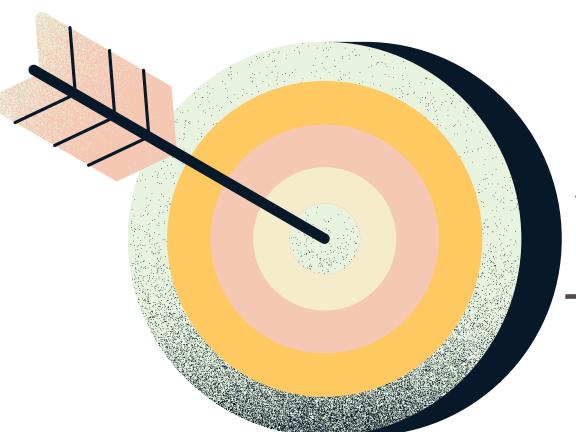
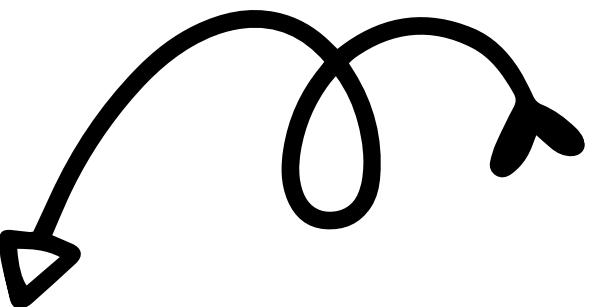


Running Hadoop Cluster on a Sector data:

```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar sector_data0.csv
```

Output:

```
Streaming final output from hdfs://
"Communication Services"      6
"Consumer Discretionary"      24
"Consumer Staples"           16
"Energy"                      5
"Financials"                  45
"Health Care"                 7
"IT"                           2
"Industrials"                 20
"Materials"                   42
"Real Estate"                 28
"Utilities"                   2
```



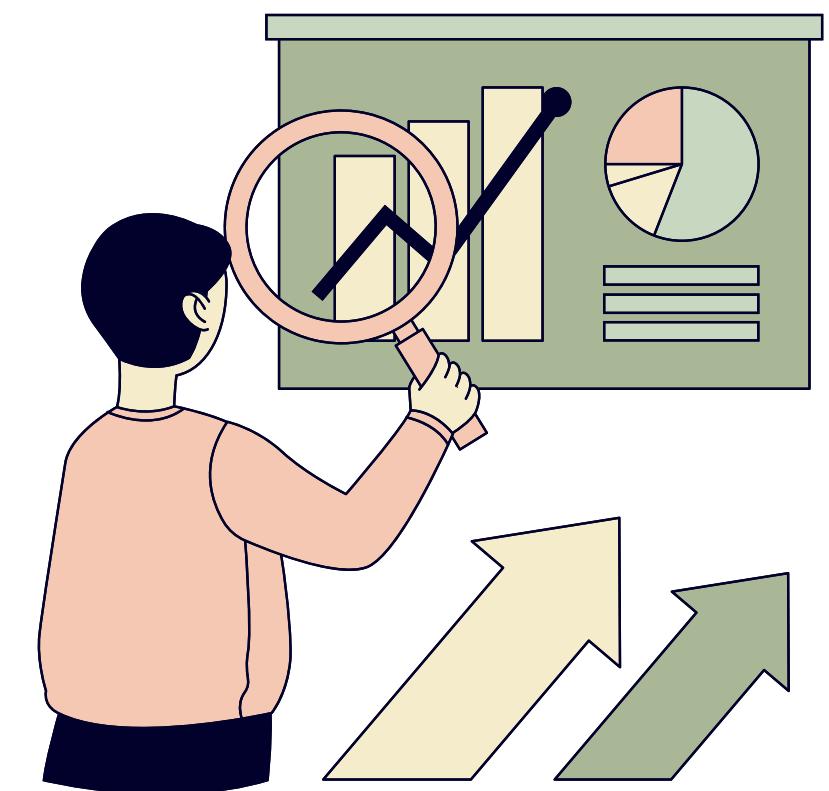
Business domain:

4-count of years to find out the year with the most shares:

In our data, each row expresses the information of one share. Among this information is the year in which this share was offered for trading so we can find the number of shares in each year, this could help Tadawul to understand the change in the stock market.



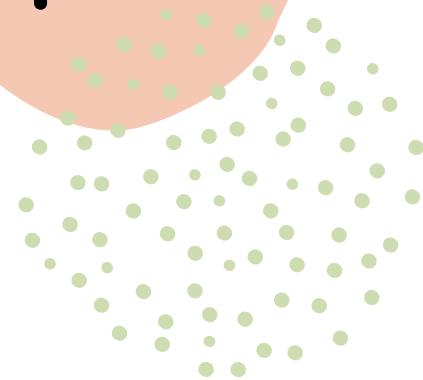
Third Approach:



Running Hadoop Cluster on a years data:

```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar years.csv
```

Output:



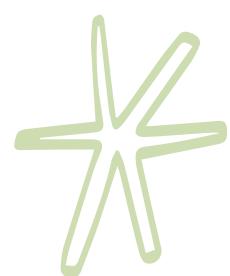
-2019 is the most year in which shares were offered for trading

| | |
|---------------|--------------|
| "2001" | 47 |
| "2002" | 14976 |
| "2003" | 17377 |
| "2004" | 19883 |
| "2005" | 21026 |
| "2006" | 20005 |
| "2007" | 22302 |
| "2008" | 28332 |
| "2009" | 30420 |
| "2010" | 33275 |
| "2011" | 34211 |
| "2012" | 36307 |
| "2013" | 37297 |
| "2014" | 38615 |
| "2015" | 40004 |
| "2016" | 40709 |
| "2017" | 42917 |
| "2018" | 45388 |
| "2019" | 47234 |
| "2020" | 9106 |

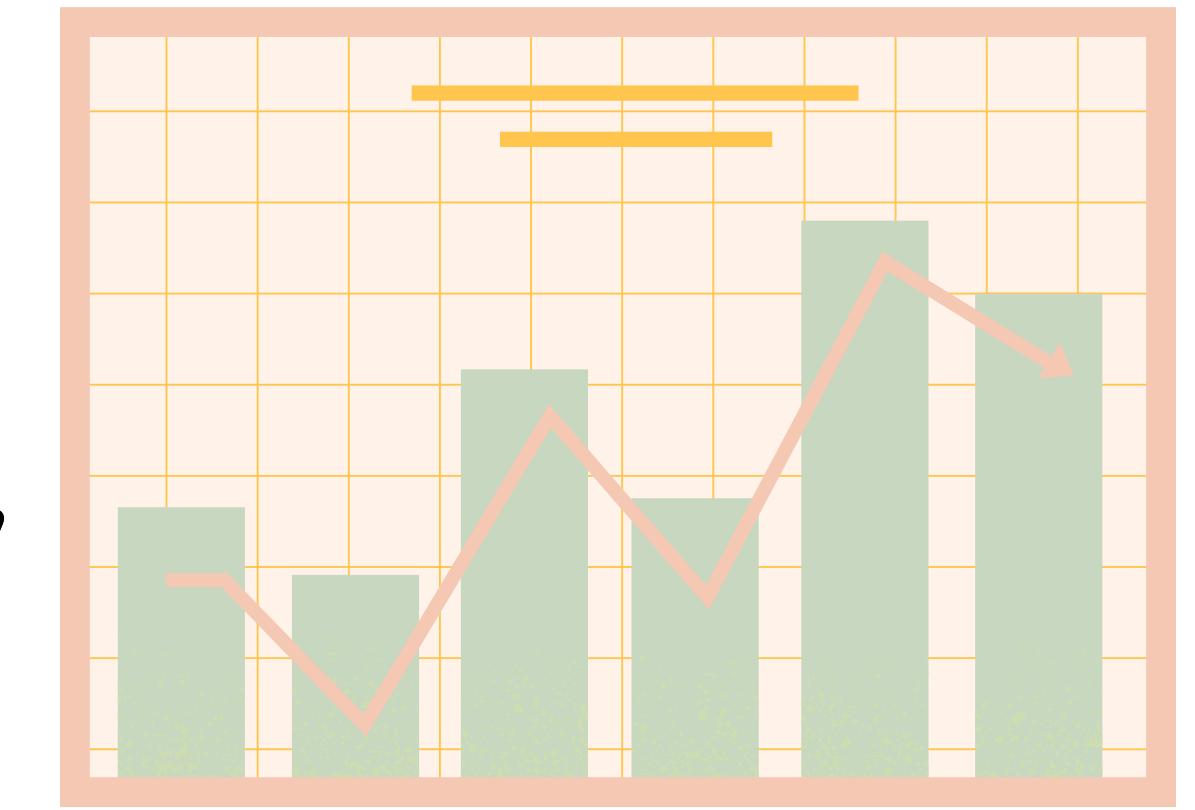
Fourth Business Problem:

4-Count the change categories for each year, month and day:

Our dataset is quite comprehensive and includes all the data necessary to derive information that is helpful to both businesses and investors. What is the ideal time to invest their money is the most crucial piece of knowledge they require. So in the following MapReduce functions, we count the occurrence of good change, stable change, and bad change with each year, each month, and each day.



Fourth Approach:



1-year

- Running Hadoop Cluster on a year_change data:

```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar years_change.csv
```



Output:

- **2019 is the year with the highest Good change (21139)**
- **2018 is the year with the highest Bad Change (23297)**
- **2002 is the most stable year (8720)**

| | | | |
|----------------------|-------------|---------------------------|--------------|
| "2001,Bad Change" | 21 | "2010,Bad Change" | 14927 |
| "2001,Good Change" | 1 | "2010,Good Change" | 13241 |
| "2001,Stable" | 25 | "2010,Stable" | 5107 |
| "2002,Bad Change" | 3111 | "2011,Bad Change" | 15084 |
| "2002,Good Change" | 3145 | "2011,Good Change" | 15046 |
| "2002,Stable" | 8720 | "2011,Stable" | 4081 |
| "2003,Bad Change" | 4334 | "2012,Bad Change" | 16664 |
| "2003,Good Change" | 4900 | "2012,Good Change" | 15771 |
| "2003,Stable" | 8143 | "2012,Stable" | 3872 |
| "2004,Bad Change" | 6539 | "2013,Bad Change" | 16037 |
| "2004,Good Change" | 7217 | "2013,Good Change" | 16179 |
| "2004,Stable" | 6127 | "2013,Stable" | 5081 |
| "2005,Bad Change" | 8654 | "2014,Bad Change" | 18736 |
| "2005,Good Change" | 9155 | "2014,Good Change" | 17975 |
| "2005,Stable" | 3217 | "2014,Stable" | 1904 |
| "2006,Bad Change" | 9758 | "2015,Bad Change" | 20623 |
| "2006,Good Change" | 9120 | "2015,Good Change" | 18744 |
| "2006,Stable" | 1127 | "2015,Stable" | 637 |
| "2007,Bad Change" | 9358 | "2016,Bad Change" | 19457 |
| "2007,Good Change" | 10025 | "2016,Good Change" | 20422 |
| "2007,Stable" | 2919 | "2016,Stable" | 830 |
| "2008,Bad Change" | 13685 | "2017,Bad Change" | 22128 |
| "2008,Good Change" | 11397 | "2017,Good Change" | 19447 |
| "2008,Stable" | 3250 | "2017,Stable" | 1342 |
| "2009,Bad Change" | 13461 | "2018,Bad Change" | 23297 |
| "2009,Good Change" | 13856 | "2018,Good Change" | 19766 |
| "2009,Stable" | 3103 | "2018,Stable" | 2325 |
| | | "2019,Bad Change" | 22676 |
| | | "2019,Good Change" | 21139 |
| | | "2019,Stable" | 3419 |
| | | "2020,Bad Change" | 4548 |
| | | "2020,Good Change" | 4089 |
| | | "2020,Stable" | 469 |

2-month

Running Hadoop Cluster on a month_change data:

```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar months_change.csv
```

Output:

-December is the month with the highest Good change (23217)

- May is the month with the highest Bad Change (24420)

- January is the most stable month (6418)

| | |
|-------------------------|-------|
| "April,Bad Change" | 20488 |
| "April,Good Change" | 22595 |
| "April,Stable" | 5305 |
| "August,Bad Change" | 21593 |
| "August,Good Change" | 18967 |
| "August,Stable" | 5876 |
| "December,Bad Change" | 21685 |
| "December,Good Change" | 23217 |
| "December,Stable" | 5476 |
| "February,Bad Change" | 21678 |
| "February,Good Change" | 20480 |
| "February,Stable" | 5502 |
| "January,Bad Change" | 23106 |
| "January,Good Change" | 22580 |
| "January,Stable" | 6418 |
| "July,Bad Change" | 22154 |
| "July,Good Change" | 20810 |
| "July,Stable" | 5861 |
| "June,Bad Change" | 20956 |
| "June,Good Change" | 19795 |
| "June,Stable" | 5175 |
| "March,Bad Change" | 22388 |
| "March,Good Change" | 22148 |
| "March,Stable" | 5993 |
| "May,Bad Change" | 24420 |
| "May,Good Change" | 20785 |
| "May,Stable" | 5242 |
| "November,Bad Change" | 21178 |
| "November,Good Change" | 20496 |
| "November,Stable" | 4788 |
| "October,Bad Change" | 23524 |
| "October,Good Change" | 19920 |
| "October,Stable" | 4981 |
| "September,Bad Change" | 19928 |
| "September,Good Change" | 18842 |
| "September,Stable" | 5081 |

3-day

- Running Hadoop Cluster on a days_change data:

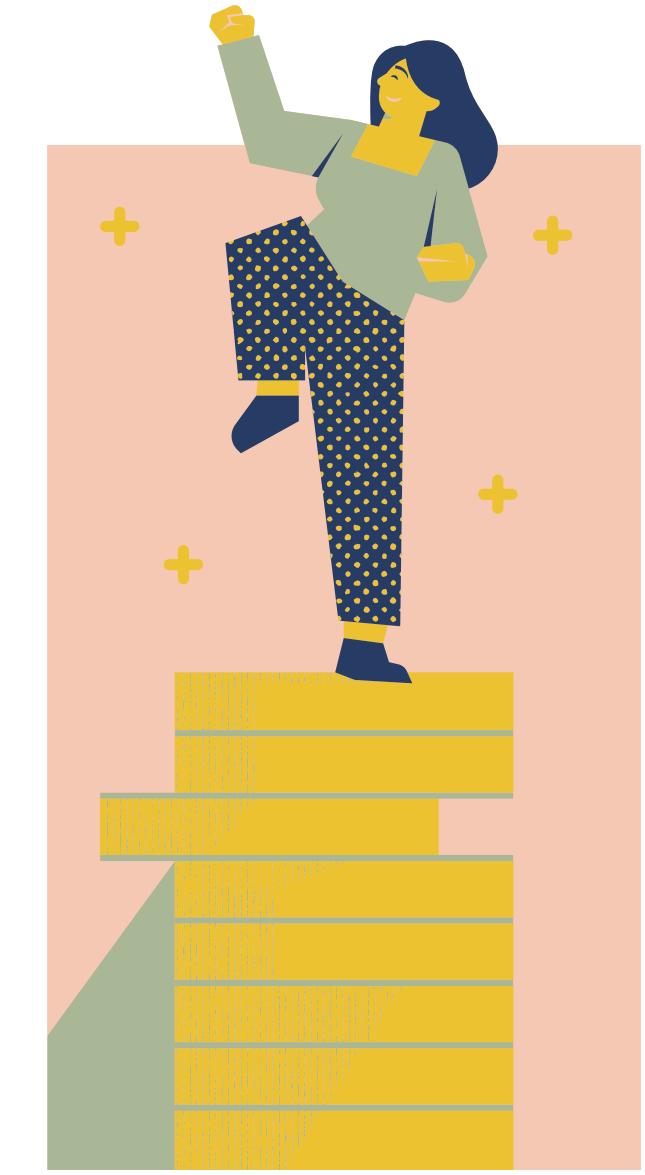
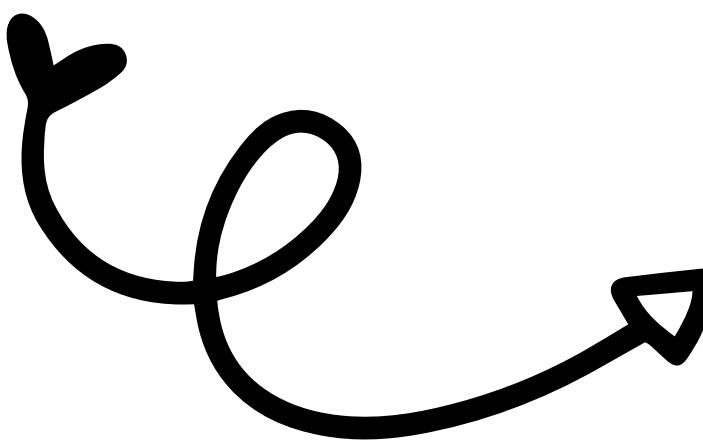
```
[maria_dev@sandbox-hdp ~]$ python project_MapReduce.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar days_change.csv
```

Output:

- Sunday is the day with the highest Good change (49147)
- Tuesday is the day with the highest Bad Change (53271)
- Wednesday is the most stable day (13029)

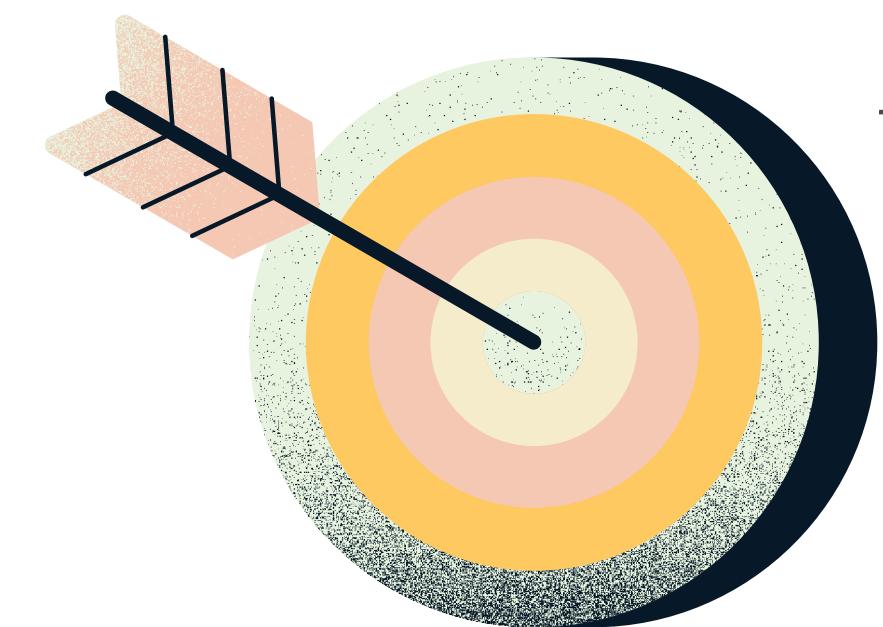
| | |
|------------------------------|--------------|
| "Monday, Bad Change" | 52606 |
| "Monday, Good Change" | 48386 |
| "Monday, Stable" | 12328 |
| "Saturday, Bad Change" | 21190 |
| "Saturday, Good Change" | 27040 |
| "Saturday, Stable" | 8145 |
| "Sunday, Bad Change" | 51603 |
| "Sunday, Good Change" | 49147 |
| "Sunday, Stable" | 12059 |
| "Thursday, Bad Change" | 31979 |
| "Thursday, Good Change" | 29523 |
| "Thursday, Stable" | 7905 |
| "Tuesday, Bad Change" | 53271 |
| "Tuesday, Good Change" | 48069 |
| "Tuesday, Stable" | 12232 |
| "Wednesday, Bad Change" | 52449 |
| "Wednesday, Good Change" | 48470 |
| "Wednesday, Stable" | 13029 |

Results



What Can Be Improved?

- combining problem 1 and 2 to which will count the number of companies in each change category in each sector as follow:
((Change category, count of change category), number of companies (reduced)).
- Compare the time for running this problem and the one that we did in the project and choose the faster one.



Any questions?



Thank you

