



# Microsoft Data Analysis

IN THE SLIDES ABOVE I AM PRESENTING METHODS I USED TO ANALYSE  
DATA FOR MICROSOFT



# Microsoft

MY CLIENT HAD IN MIND TO START MOVIE PRODUCTION AND I DID THE  
BELOW DATA ANALYSIS .

# The files provided for data analysis :

- ▶ 1.RT.REVIEWS.TSV
- ▶ 2. RT.MOVIES.INFO.TSV
- ▶ 3. MOVIES.CSV
- ▶ 4. MOVIES\_BUDGET.CSV
- ▶ 5. MOVIE\_GROSS

# I choose the following files to start my work

- ▶ RT.REVIEWS.TSV
- ▶ MOVIES.CSV
- ▶ MOVIES-BUDGET.CSV

# WORK PROCEDURE

- ▶ 1) I started my work with converting the files into data frames.
- ▶ eg movie\_budget.csv to movie\_budgetcsv\_df
- ▶ A data frame is a structure that organizes data into a two dimensional table of rows and columns.
- ▶ From there I printed the output to check my organized structure

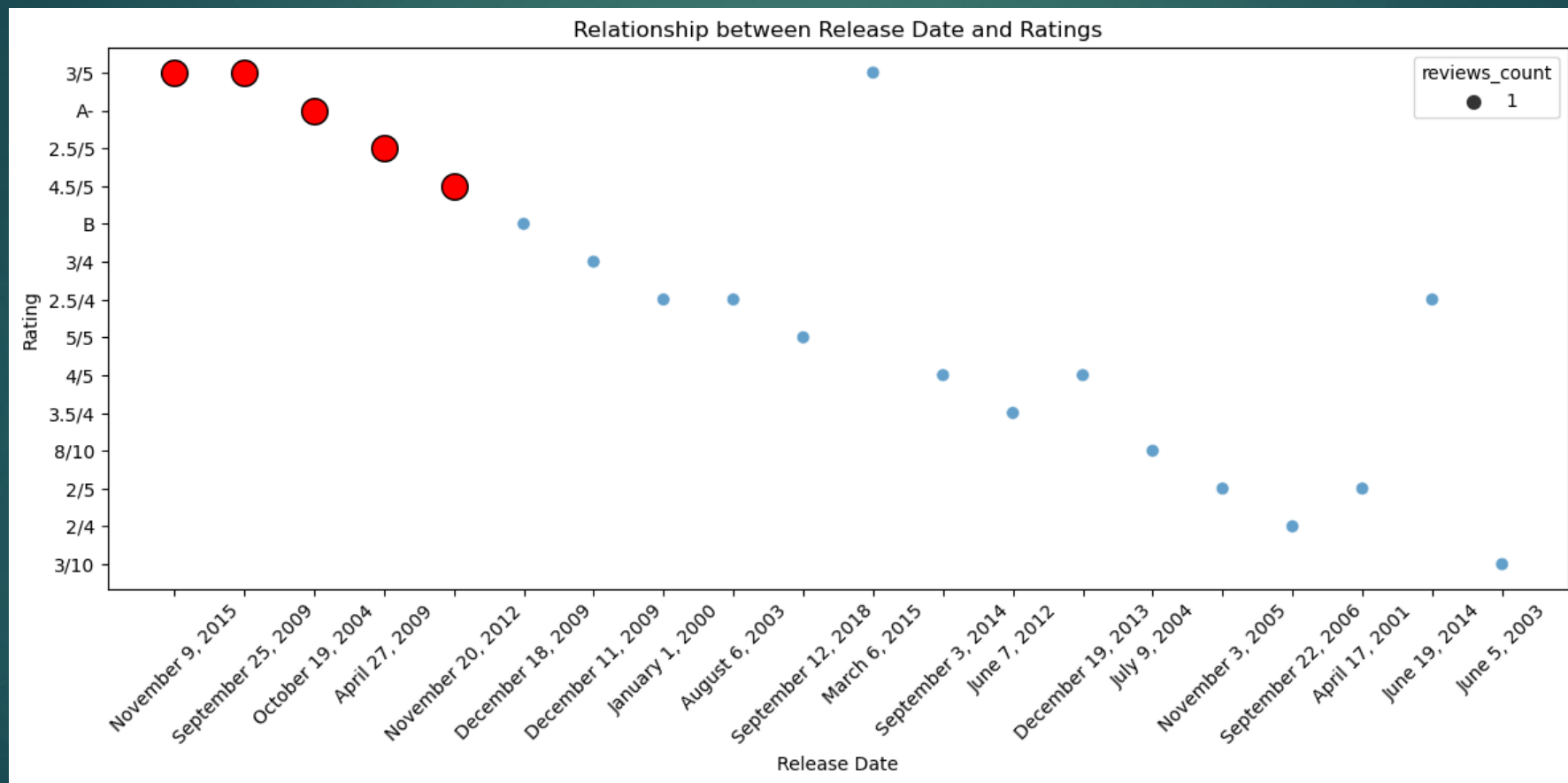
# PROCEDURE TWO:

- ▶ I conducted data cleaning to erase
  - ▶ -duplicates
  - ▶ - remove the incomplete data types
  - ▶ -remove the missing files.
  - ▶ -remove undesired formats
- ▶ - this enhances data accuracy while analyzing data.

# QUESTIONS

- ▶ I set my work into questions for easy understanding and flow of work.
- ▶ I used the `rt_reviews_df`
- ▶ It contained columns like;
- ▶ 'publisher', 'ratings,' 'release date' and 'reviews'
- ▶ :The question answered the below arguments:
- ▶ I used the question to calculate the reviews of each publisher ,
- ▶ from there i used data visualization by using a scatter graph to show the relationship between release date and rating . the last question was counting the number of ratings. –

# QUESTION ONE GRAPH

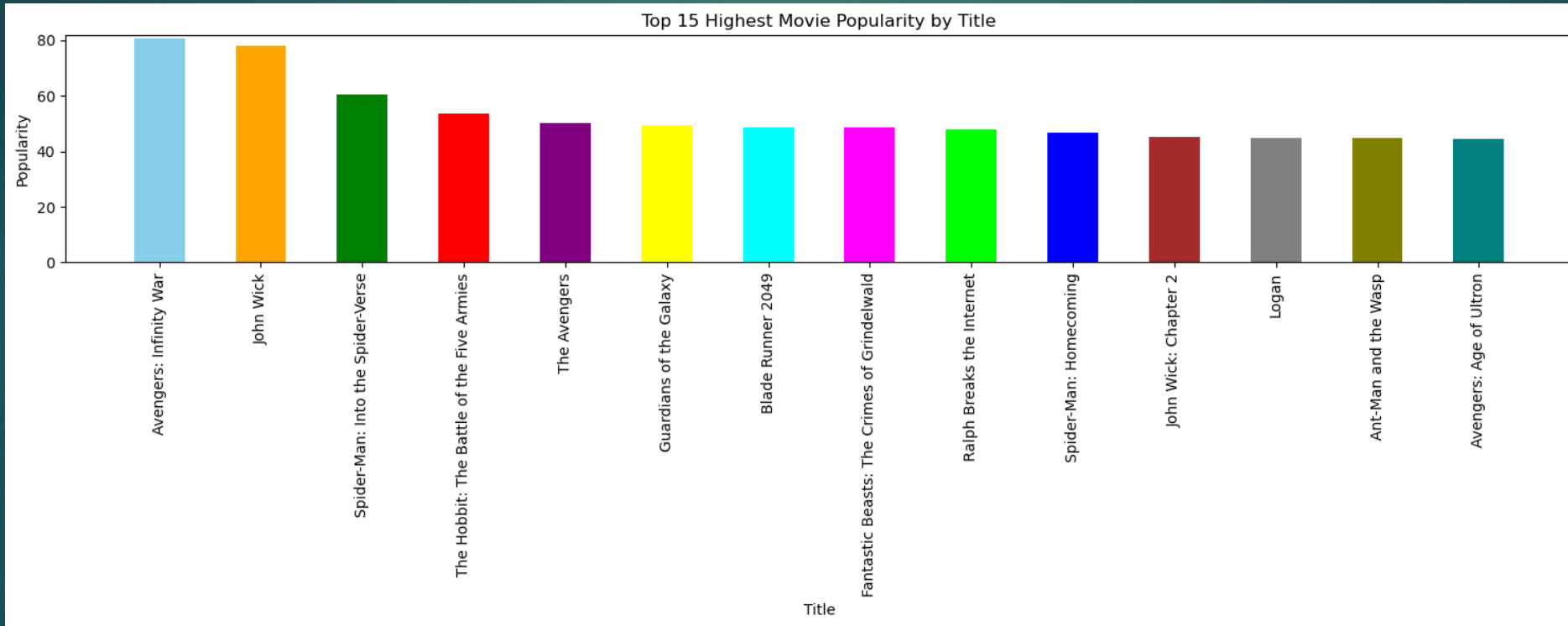




# QUESTION TWO

- ▶ In question two I followed the same directory as question one
- ▶ I used `movie_csv_df`
- ▶ The dataframe consist columns eg: 'popularity', 'genre\_ids', 'title'
- ▶ The question answered :the below arguments:
  - ▶ -check the movie title with the highest vote\_count
  - ▶ -check the movie title which was the most popular
  - ▶ -checked the ratings of the movie title

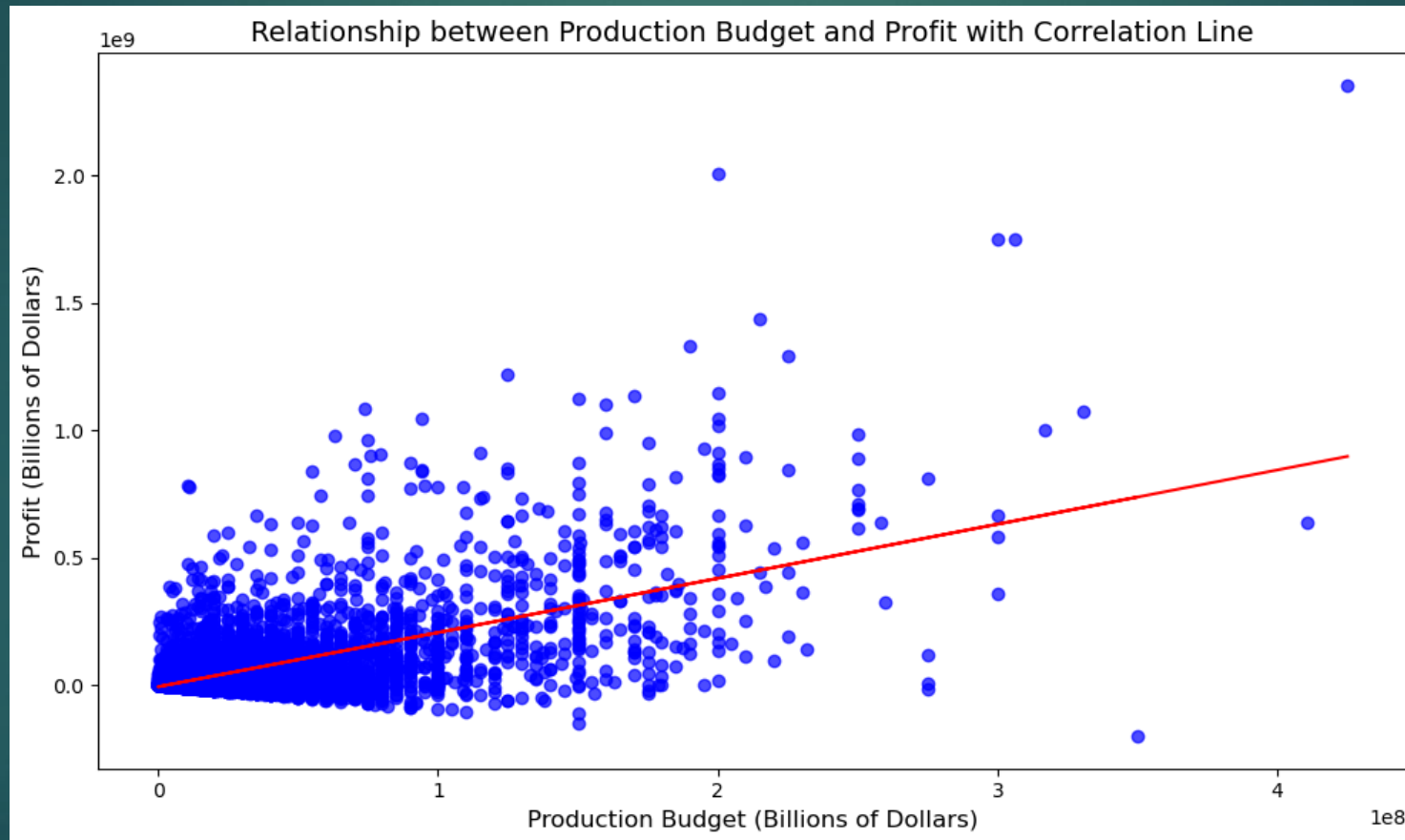
# QUEESTION TWO GRAPH



# QUESTION THREE

- in the above question I used the movie\_budgetcsv\_df
- The dataframe consists of the below columns:
  - 'worldwide\_gross', 'production\_budget', release date
- The answered the below arguments:
  - calculate the profit margin
  - show the relationship between profit and budget
  - relationship between release date and budget.

# QUESTION 3 GRAPH



# RECOMMENDATIONS

- ▶ Release movies during the last three months of the year.
- ▶ Hire publishers with the high rating for the movies,
- ▶ Choose a movie title that will be popular and get a high vote\_count
- ▶ Allocate a budget that is good enough to yeild high profits

# DATA SCIENTIST DETAILS:

ESTHER MUKAMI NJAGI

LINKEDLN PROFILE: Esther Mukami