# Data Mining (DSC550-T301_2245_1)

Assignement Week 1; Author: Zemelak Goraga; Date: 03/16/2024

In [1]:
```python
# Step 1: Import necessary libraries
import pandas as pd
```

In [4]:
```python
# Step 2: Load the dataset into a Pandas DataFrame and save it as df
df = pd.read_csv("Video_Games_Sales_as_at_22_Dec_2016.csv")
```

In [5]:
```python
# Step 3: Display the first ten rows of the df dataset
print("First ten rows of the dataset:")
print(df.head(10))
```

```
First ten rows of the dataset:
                      Name Platform  Year_of_Release        Genre  \
0                Wii Sports      Wii           2006.0       Sports
1         Super Mario Bros.     NES           1985.0     Platform
2            Mario Kart Wii     Wii           2008.0       Racing
3         Wii Sports Resort     Wii           2009.0       Sports
4   Pokemon Red/Pokemon Blue      GB           1996.0  Role-Playing
5                    Tetris      GB           1989.0       Puzzle
6      New Super Mario Bros.      DS           2006.0     Platform
7                  Wii Play     Wii           2006.0         Misc
8   New Super Mario Bros. Wii     Wii          2009.0     Platform
9                 Duck Hunt     NES           1984.0      Shooter

   Publisher  NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  \
0   Nintendo     41.36     28.96      3.77         8.45         82.53
1   Nintendo     29.08      3.58      6.81         0.77         40.24
2   Nintendo     15.68     12.76      3.79         3.29         35.52
3   Nintendo     15.61     10.93      3.28         2.95         32.77
4   Nintendo     11.27      8.89     10.22         1.00         31.37
5   Nintendo     23.20      2.26      4.22         0.58         30.26
6   Nintendo     11.28      9.14      6.50         2.88         29.80
7   Nintendo     13.96      9.18      2.93         2.84         28.92
8   Nintendo     14.44      6.94      4.70         2.24         28.32
9   Nintendo     26.93      0.63      0.28         0.47         28.31

   Critic_Score  Critic_Count User_Score  User_Count Developer Rating
0          76.0          51.0          8       322.0  Nintendo      E
1           NaN           NaN        NaN         NaN       NaN    NaN
2          82.0          73.0        8.3       709.0  Nintendo      E
3          80.0          73.0          8       192.0  Nintendo      E
4           NaN           NaN        NaN         NaN       NaN    NaN
5           NaN           NaN        NaN         NaN       NaN    NaN
6          89.0          65.0        8.5       431.0  Nintendo      E
7          58.0          41.0        6.6       129.0  Nintendo      E
8          87.0          80.0        8.4       594.0  Nintendo      E
9           NaN           NaN        NaN         NaN       NaN    NaN
```

In [6]:
```python
# Step 4: Display the columns of the df dataset
print("\nColumns of the dataset:")
print(df.columns)
```

```
Columns of the dataset:
Index(['Name', 'Platform', 'Year_of_Release', 'Genre', 'Publisher', 'NA_Sales',
       'EU_Sales', 'JP_Sales', 'Other_Sales', 'Global_Sales', 'Critic_Score',
       'Critic_Count', 'User_Score', 'User_Count', 'Developer', 'Rating'],
      dtype='object')
```

```
In [7]:  # Step 5: Find the dimensions (number of rows and columns) in the df data frame
         num_rows, num_cols = df.shape
         print(f"\nDimensions of the dataset: {num_rows} rows x {num_cols} columns")
```

```
Dimensions of the dataset: 16719 rows x 16 columns
```

```
In [9]:  # Step 6: Find the top five games by critic score
         top_games_by_critic_score = df.nlargest(5, 'Critic_Score')[['Name', 'Critic_Score']
         print("\nTop five games by critic score:")
         print(top_games_by_critic_score)
```

```
Top five games by critic score:
                        Name  Critic_Score
51           Grand Theft Auto IV          98.0
57           Grand Theft Auto IV          98.0
227    Tony Hawk's Pro Skater 2          98.0
5350                 SoulCalibur          98.0
16             Grand Theft Auto V          97.0
```

```
In [11]:  # Step 7: Find the number of video games in the df data frame in each genre
          genre_counts = df['Genre'].value_counts()
          print("\nNumber of video games in each genre:")
          print(genre_counts)
```

```
Number of video games in each genre:
Action          3370
Sports          2348
Misc            1750
Role-Playing    1500
Shooter         1323
Adventure       1303
Racing          1249
Platform         888
Simulation       874
Fighting         849
Strategy         683
Puzzle           580
Name: Genre, dtype: int64
```

```
In [13]:  # Step 8: Find the first five games in the df data frame on the SNES platform
          snes_games = df[df['Platform'] == 'SNES'].head(5)
          print("\nFirst five games on the SNES platform:")
          print(snes_games)
```

```
First five games on the SNES platform:
                                Name Platform  Year_of_Release     Genre  \
18                   Super Mario World     SNES           1990.0  Platform
56                Super Mario All-Stars     SNES           1993.0  Platform
71                Donkey Kong Country     SNES           1994.0  Platform
76                   Super Mario Kart     SNES           1992.0    Racing
137  Street Fighter II: The World Warrior     SNES           1992.0  Fighting

     Publisher  NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  \
18    Nintendo     12.78      3.75      3.54         0.55         20.61
56    Nintendo      5.99      2.15      2.12         0.29         10.55
71    Nintendo      4.36      1.71      3.00         0.23          9.30
76    Nintendo      3.54      1.24      3.81         0.18          8.76
137     Capcom      2.47      0.83      2.87         0.12          6.30

     Critic_Score  Critic_Count User_Score  User_Count Developer Rating
18            NaN           NaN        NaN         NaN       NaN    NaN
56            NaN           NaN        NaN         NaN       NaN    NaN
71            NaN           NaN        NaN         NaN       NaN    NaN
76            NaN           NaN        NaN         NaN       NaN    NaN
137           NaN           NaN        NaN         NaN       NaN    NaN
```

In [15]:
```python
# Step 9: Find the five publishers with the highest total global sales
publisher_sales = df.groupby('Publisher')['Global_Sales'].sum().nlargest(5)
print("\nFive publishers with the highest total global sales:")
print(publisher_sales)
```

```
Five publishers with the highest total global sales:
Publisher
Nintendo                       1788.81
Electronic Arts                1116.96
Activision                      731.16
Sony Computer Entertainment     606.48
Ubisoft                         471.61
Name: Global_Sales, dtype: float64
```

In [16]:
```python
# Step 10: Create a new column for the percentage of global sales from North Americ
df['NA_Sales_Percentage'] = (df['NA_Sales'] / df['Global_Sales']) * 100
```

In [17]:
```python
# Step 11: Display the first five rows of the new DataFrame
print("\nFirst five rows with the new column:")
print(df.head())
```

```
First five rows with the new column:
                   Name Platform  Year_of_Release        Genre Publisher  \
0            Wii Sports      Wii           2006.0       Sports  Nintendo
1       Super Mario Bros.      NES           1985.0     Platform  Nintendo
2         Mario Kart Wii      Wii           2008.0       Racing  Nintendo
3      Wii Sports Resort      Wii           2009.0       Sports  Nintendo
4  Pokemon Red/Pokemon Blue       GB           1996.0  Role-Playing  Nintendo

   NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  Critic_Score  \
0     41.36     28.96      3.77         8.45         82.53          76.0
1     29.08      3.58      6.81         0.77         40.24           NaN
2     15.68     12.76      3.79         3.29         35.52          82.0
3     15.61     10.93      3.28         2.95         32.77          80.0
4     11.27      8.89     10.22         1.00         31.37           NaN

   Critic_Count User_Score  User_Count Developer Rating  NA_Sales_Percentage
0          51.0          8       322.0  Nintendo      E            50.115110
1           NaN        NaN         NaN       NaN    NaN            72.266402
2          73.0        8.3       709.0  Nintendo      E            44.144144
3          73.0          8       192.0  Nintendo      E            47.635032
4           NaN        NaN         NaN       NaN    NaN            35.926044
```

```
In [18]:  # Step 12: Find the number of NaN entries in each column
          nan_counts = df.isna().sum()
          print("\nNumber of NaN entries in each column:")
          print(nan_counts)

Number of NaN entries in each column:
Name                      2
Platform                  0
Year_of_Release         269
Genre                     2
Publisher                54
NA_Sales                  0
EU_Sales                  0
JP_Sales                  0
Other_Sales               0
Global_Sales              0
Critic_Score           8582
Critic_Count           8582
User_Score             6704
User_Count             9129
Developer              6623
Rating                 6769
NA_Sales_Percentage       0
dtype: int64
```

```
In [19]:  # Step 13: Replace non-numerical user score entries with NaN
          df['User_Score'] = pd.to_numeric(df['User_Score'], errors='coerce')
```

```
In [21]:  # Step 14: Calculate the median user score
          median_user_score = df['User_Score'].median()
          median_user_score
```

Out[21]:  7.5

```
In [22]:  # Step 15: Replace NaN entries in the user score column with the median value
          df['User_Score'].fillna(median_user_score, inplace=True)
```

```
In [24]:  # Step 16: Display the updated DataFrame
          print("\nUpdated DataFrame with NaN replaced:")
          print(df)
```

```
Updated DataFrame with NaN replaced:
                                 Name Platform  Year_of_Release        Genre  \
0                           Wii Sports     Wii           2006.0       Sports
1                     Super Mario Bros.     NES           1985.0     Platform
2                        Mario Kart Wii     Wii           2008.0       Racing
3                      Wii Sports Resort     Wii           2009.0       Sports
4                 Pokemon Red/Pokemon Blue      GB          1996.0  Role-Playing
...                                  ...     ...              ...          ...
16714   Samurai Warriors: Sanada Maru     PS3           2016.0       Action
16715              LMA Manager 2007    X360           2006.0       Sports
16716        Haitaka no Psychedelica     PSV           2016.0    Adventure
16717               Spirits & Spells     GBA           2003.0     Platform
16718             Winning Post 8 2016     PSV           2016.0   Simulation

            Publisher  NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  \
0            Nintendo     41.36     28.96      3.77         8.45         82.53
1            Nintendo     29.08      3.58      6.81         0.77         40.24
2            Nintendo     15.68     12.76      3.79         3.29         35.52
3            Nintendo     15.61     10.93      3.28         2.95         32.77
4            Nintendo     11.27      8.89     10.22         1.00         31.37
...               ...       ...       ...       ...          ...           ...
16714    Tecmo Koei      0.00      0.00      0.01         0.00          0.01
16715   Codemasters      0.00      0.01      0.00         0.00          0.01
16716   Idea Factory      0.00      0.00      0.01         0.00          0.01
16717       Wanadoo      0.01      0.00      0.00         0.00          0.01
16718    Tecmo Koei      0.00      0.00      0.01         0.00          0.01

       Critic_Score  Critic_Count  User_Score  User_Count Developer Rating  \
0              76.0          51.0         8.0       322.0  Nintendo      E
1               NaN           NaN         7.5         NaN       NaN    NaN
2              82.0          73.0         8.3       709.0  Nintendo      E
3              80.0          73.0         8.0       192.0  Nintendo      E
4               NaN           NaN         7.5         NaN       NaN    NaN
...             ...           ...         ...         ...       ...    ...
16714           NaN           NaN         7.5         NaN       NaN    NaN
16715           NaN           NaN         7.5         NaN       NaN    NaN
16716           NaN           NaN         7.5         NaN       NaN    NaN
16717           NaN           NaN         7.5         NaN       NaN    NaN
16718           NaN           NaN         7.5         NaN       NaN    NaN

       NA_Sales_Percentage
0                50.115110
1                72.266402
2                44.144144
3                47.635032
4                35.926044
...                    ...
16714             0.000000
16715             0.000000
16716             0.000000
16717           100.000000
16718             0.000000

[16719 rows x 17 columns]
```

In [ ]:

Title: Comprehensive Analysis of Video Game Sales with Ratings Dataset

Summary: This report provides a comprehensive analysis of the Video Game Sales with Ratings dataset, focusing on various aspects such as top games by critic score, genre distribution, publisher sales, and user score data wrangling. Through thorough examination

and analysis, valuable insights into the video game industry are derived, aiding stakeholders in making informed decisions.

Introduction: The video game industry has witnessed exponential growth over the years, with the rise of various platforms and genres catering to diverse audiences. Understanding the dynamics of this industry is crucial for stakeholders, including developers, publishers, and investors. The Video Game Sales with Ratings dataset offers a wealth of information that can be leveraged to gain insights into consumer preferences, market trends, and more.

Statement of the Problem: The dataset presents several challenges and opportunities for analysis. Key issues include missing data entries, inconsistent formats, and the need to derive meaningful insights from the available information. The goal is to extract actionable insights that can inform business strategies and decision-making processes.

Methodology:

Data Acquisition: The dataset was obtained from Kaggle using the Kaggle API. Data Preprocessing: Data cleaning and wrangling techniques were applied to handle missing values, format inconsistencies, and prepare the data for analysis. Exploratory Data Analysis (EDA): Various statistical and visual methods were employed to explore the dataset and uncover patterns, trends, and relationships. Data Analysis: Quantitative analysis techniques were used to derive insights into key metrics such as top games, genre distribution, publisher sales, and user scores.

Dimensions of the dataset: 16719 rows x 16 columns, representing the number of observations (video games) and studied variables about the games, respectively.

Results:

Top five games by critic score:

Grand Theft Auto IV - 98.0; Grand Theft Auto IV - 98.0; Tony Hawk's Pro Skater 2 - 98.0; SoulCalibur - 98.0; Grand Theft Auto V - 97.0

Genre distribution:

Action: 3370; Sports: 2348; Misc: 1750; Role-Playing: 1500; Shooter: 1323; Adventure: 1303; Racing: 1249; Platform: 888; Simulation: 874; Fighting: 849; Strategy: 683; Puzzle: 580;

Publisher sales:

Nintendo: 1788.81; Electronic Arts: 1116.96; Activision: 731.16; Sony Computer Entertainment: 606.48; Ubisoft: 471.61

User score data wrangling:

Median user score: 7.5 NaN entries replaced

Discussion of Results:

Top Games: The analysis reveals that Grand Theft Auto IV, Tony Hawk's Pro Skater 2, and SoulCalibur are among the top-rated games by critics, suggesting a strong market demand

for immersive gaming experiences. This indicates potential opportunities for game developers and publishers to focus on creating high-quality titles that resonate with players.

Genre Distribution: Action games dominate the market, followed by sports and miscellaneous genres. This highlights the diverse preferences of gamers and underscores the importance of catering to various interests to maximize market reach and revenue potential. Developers may benefit from targeting specific genres based on consumer demand and emerging trends.

Publisher Sales: Nintendo emerges as the top publisher with the highest total global sales, emphasizing the significance of brand reputation and quality content in driving sales. Other leading publishers such as Electronic Arts and Activision also play a significant role in shaping the gaming landscape, indicating fierce competition within the industry.

User Score Data: The median user score of 7.5 reflects a generally positive reception among gamers, although challenges such as missing data entries necessitate robust data cleaning and preprocessing techniques. Addressing these issues is essential to ensure the accuracy and reliability of user score data, enabling stakeholders to make informed decisions based on consumer feedback.

Conclusions: The analysis of the Video Game Sales with Ratings dataset provides valuable insights into the video game industry, offering stakeholders a deeper understanding of market dynamics and consumer behavior. By leveraging these insights, stakeholders can make informed decisions to enhance product development, marketing strategies, and overall business performance.

Recommendations:

Data Quality Assurance: Implement robust data cleaning and validation processes to ensure data accuracy and consistency. Market Segmentation: Utilize genre preferences and user demographics to tailor marketing campaigns and product offerings. Investment Strategies: Consider partnering with top publishers or investing in genres with high market demand and potential for growth. Way Forward: Further research could focus on longitudinal analysis to track trends over time, sentiment analysis of user reviews to gauge consumer sentiment, and predictive modeling to forecast future sales trends and game popularity.

In [ ]: