

INTERNSHIP PROJECT REPORT
ON
“VIDEO GAMES ANALYSIS”

SUBMITTED BY:
PRAYAS SAMAL

EXECUTIVE SUMMARY

The "Video Game Sales Analysis" project report examines the variables influencing video games' financial performance. To find patterns, trends, and statistically significant relationships between different game qualities and their commercial performance, the analysis employs a quantitative research methodology to look at a sizable dataset of video game sales. The paper explores the market dynamics for video games and offers insights that industry stakeholders may find useful.

The study makes use of Kaggle data, which includes a range of video games from various platforms and time periods. Total sales, platform, genre, publisher, developer, year of release, and reviewer ratings are among the important characteristics that are examined. Descriptive statistics, grouping and aggregation, exploratory data analysis (EDA), statistical inference (including t-tests, ANOVA, and chi-squared tests), and correlation analysis are all part of the analytical approach. Python is used for the data analysis, along with libraries like scipy.stats, matplotlib, seaborn, and pandas.

Several significant insights are revealed by the analysis's findings:

- 64,016 video game titles with 14 features—including qualitative and numerical data—make up the dataset.
- The distribution of total sales is skewed, and a small number of best-selling titles dominate the market.
- Japan and other regions have smaller market shares than North America and the PAL region, which contribute the most to overall sales.
- Platforms (such as the PS2 and X360) and major publishers (like Activision and Electronic Arts) have shown impressive overall sales performance.
- Among the most popular game genres are action and adventure games.
- There is a weak positive association between critic scores and overall sales, and critic reviews are often positive.

Missing data, possible data discrepancies, the inability to prove causation, a limited collection of features, the scope of the study, and the potential for bias in statistical tests are just a few of the limitations noted in the report.

TABLE OF CONTENTS

SERIAL NO	PARTICULARS	PAGE NO
1.	EXECUTIVE SUMMARY	2
2.	OBJECTIVES	4
3.	DATA ANALYSIS	5-17
4.	CONCLUSION	18

OBJECTIVES

"Video Game Sales Analysis" is defined for the purposes of this project report as the methodical process of looking at and analysing quantitative data pertaining to the commercial performance of video games. This procedure entails:

- **Data Acquisition and Preparation:** Compiling unprocessed data on video game sales numbers, possibly including variables like genre, publisher, developer, release date, critic ratings, platform (console, PC, mobile), revenue generated (both globally and in particular regions like North America, Japan, and Europe), and units sold. Data organization and cleansing for analysis are also part of this step.
- **Descriptive statistics and exploratory data analysis (EDA):** determining summary metrics for sales data and other pertinent numerical variables, such as mean, median, and standard deviation. Using data visualizations (such as bar charts, histograms, and scatter plots) to comprehend the distribution of important variables and spot patterns, trends, and possible links.
- **Comparative analysis** is the study of sales success in many categories. For example, it compares the overall sales of various game consoles, genres, publishers, or time periods (e.g., year of release). This could entail figuring out and displaying variations in average or total sales.
- **Examining possible relationships or correlations** between various factors and video game sales is known as relationship analysis. For instance, employing statistical tests such as the chi-squared test to investigate the relationship between console and genre or correlation coefficients and scatter plots to investigate the relationship between critic scores and overall sales
- **Statistical Inference:** By using statistical tests (such as ANOVA and t-tests) to determine the significance of observed associations or differences in the data, conclusions can be made about the larger population of video games outside of the particular dataset.

The ultimate goal of this "Video Game Sales Analysis," as defined within this report, is to extract meaningful insights from the sales data to understand the factors influencing commercial success in the video game industry, identify trends, and potentially inform future strategies or predictions. The specific techniques and variables analysed will be detailed within the methodology section of this report.

DATA ANALYSIS

```
#Importing libraries
import numpy as np # linear algebra
import pandas as pd # data processing
import warnings
warnings.filterwarnings("ignore")
```

```
#Importing data
df = pd.read_csv('vgchartz-2024.csv')
```

```
#Looking into the data
df.head(5)
```

	img	title	console	genre	publisher	developer	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales	release_
0	/games/boxart/full_6510540AmericaFrontccc.jpg	Grand Theft Auto V	PS3	Action	Rockstar Games	Rockstar North	9.4	20.32	6.37	0.99	9.85	3.12	2013-0
1	/games/boxart/full_5563178AmericaFrontccc.jpg	Grand Theft Auto V	PS4	Action	Rockstar Games	Rockstar North	9.7	19.39	6.06	0.60	9.71	3.02	2014-1
2	/games/boxart/827563ccc.jpg	Grand Theft Auto: Vice City	PS2	Action	Rockstar Games	Rockstar North	9.6	16.15	8.41	0.47	5.49	1.78	2002-1
3	/games/boxart/full_9218923AmericaFrontccc.jpg	Grand Theft Auto V	X360	Action	Rockstar Games	Rockstar North	NaN	15.86	9.06	0.06	5.33	1.42	2013-0
4	/games/boxart/full_4990510AmericaFrontccc.jpg	Call of Duty: Black Ops 3	PS4	Shooter	Activision	Treyarch	8.1	15.09	6.18	0.41	6.05	2.44	2015-1

```
df.tail()
```

	img	title	console	genre	publisher	developer	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales
64011	/games/boxart/full_2779838AmericaFrontccc.jpg	XBlaze Lost: Memories	PC	Visual Novel	Aksys Games	Arc System Works	NaN	NaN	NaN	NaN	NaN	NaN
64012	/games/boxart/full_8031506AmericaFrontccc.jpg	Yoru, Tomosu	PS4	Visual Novel	Nippon Ichi Software	Nippon Ichi Software	NaN	NaN	NaN	NaN	NaN	NaN
64013	/games/boxart/full_6553045AmericaFrontccc.jpg	Yoru, Tomosu	NS	Visual Novel	Nippon Ichi Software	Nippon Ichi Software	NaN	NaN	NaN	NaN	NaN	NaN
64014	/games/boxart/full_6012940JapanFrontccc.png	Yunohana SpRING! ~Mellow Times~	NS	Visual Novel	Idea Factory	Otomate	NaN	NaN	NaN	NaN	NaN	NaN
64015	/games/boxart/default.jpg	Yurukill: The Calumniation Games	PS4	Visual Novel	Unknown	G.rev Ltd.	NaN	NaN	NaN	NaN	NaN	NaN

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 64016 entries, 0 to 64015
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   img                    64016 non-null  object 
1   title                  64016 non-null  object 
2   console                64016 non-null  object 
3   genre                  64016 non-null  object 
4   publisher               64016 non-null  object 
5   developer              63999 non-null  object 
6   critic_score           6678 non-null   float64
7   total_sales            18922 non-null  float64
8   na_sales               12637 non-null  float64
9   jp_sales               6726 non-null   float64
10  pal_sales              12824 non-null  float64
11  other_sales            15128 non-null  float64
12  release_date           56965 non-null  object 
13  last_update            17879 non-null  object 
dtypes: float64(6), object(8)
memory usage: 6.8+ MB
```

```
df.columns
```

```
Index(['img', 'title', 'console', 'genre', 'publisher', 'developer',
       'critic_score', 'total_sales', 'na_sales', 'jp_sales', 'pal_sales',
       'other_sales', 'release_date', 'last_update'],
      dtype='object')
```

```
df.dtypes
```

```
img                object
title              object
console            object
genre              object
publisher          object
developer          object
critic_score       float64
total_sales        float64
na_sales           float64
jp_sales           float64
pal_sales          float64
other_sales        float64
release_date       object
last_update        object
dtype: object
```

```
df.shape
```

```
(64016, 14)
```

```
df.describe()
```

	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales
count	6678.000000	18922.000000	12637.000000	6726.000000	12824.000000	15128.000000
mean	7.220440	0.349113	0.264740	0.102281	0.149472	0.043041
std	1.457066	0.807462	0.494787	0.168811	0.392653	0.126643
min	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	6.400000	0.030000	0.050000	0.020000	0.010000	0.000000
50%	7.500000	0.120000	0.120000	0.040000	0.040000	0.010000
75%	8.300000	0.340000	0.280000	0.120000	0.140000	0.030000
max	10.000000	20.320000	9.760000	2.130000	9.850000	3.120000

```
df.nunique()
```

```
img          56177
title        39798
console       81
genre        20
publisher    3383
developer    8862
critic_score  89
total_sales  482
na_sales     320
jp_sales     121
pal_sales    256
other_sales  133
release_date 7922
last_update  1545
dtype: int64
```

```
df.isnull().any()
```

```
img          False
title        False
console      False
genre        False
publisher    False
developer    True
critic_score True
total_sales  True
na_sales     True
jp_sales     True
pal_sales    True
other_sales  True
release_date True
last_update  True
dtype: bool
```

```
#Cleaning the dataset
```

```
df_cleaned = df.dropna(subset=['release_date', 'developer'])
sales_columns = ['total_sales', 'na_sales', 'jp_sales', 'pal_sales', 'other_sales']
df_cleaned[sales_columns] = df_cleaned[sales_columns].fillna(0)
df_cleaned
```

	img	title	console	genre	publisher	developer	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales
0	/games/boxart/full_6510540AmericaFrontccc.jpg	Grand Theft Auto V	PS3	Action	Rockstar Games	Rockstar North	9.4	20.32	6.37	0.99	9.85	3.12
1	/games/boxart/full_5563178AmericaFrontccc.jpg	Grand Theft Auto V	PS4	Action	Rockstar Games	Rockstar North	9.7	19.39	6.06	0.60	9.71	3.02
2	/games/boxart/827563ccc.jpg	Grand Theft Auto: Vice City	PS2	Action	Rockstar Games	Rockstar North	9.6	16.15	8.41	0.47	5.49	1.78
3	/games/boxart/full_9218923AmericaFrontccc.jpg	Grand Theft Auto V	X360	Action	Rockstar Games	Rockstar North	NaN	15.86	9.06	0.06	5.33	1.42
4	/games/boxart/full_4990510AmericaFrontccc.jpg	Call of Duty: Black Ops 3	PS4	Shooter	Activision	Treyarch	8.1	15.09	6.18	0.41	6.05	2.44
...
64010	/games/boxart/full_2294305JapanFrontccc.jpg	World End Syndrome	PS4	Visual Novel	Arc System Works	Arc System Works	NaN	0.00	0.00	0.00	0.00	0.00

```
df_cleaned['last_update'] = df_cleaned['last_update'].fillna(df_cleaned['release_date'])
df = df_cleaned
df
```

	img	title	console	genre	publisher	developer	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales
0	/games/boxart/full_6510540AmericaFrontccc.jpg	Grand Theft Auto V	PS3	Action	Rockstar Games	Rockstar North	9.4	20.32	6.37	0.99	9.85	3.12
1	/games/boxart/full_5563178AmericaFrontccc.jpg	Grand Theft Auto V	PS4	Action	Rockstar Games	Rockstar North	9.7	19.39	6.06	0.60	9.71	3.02
2	/games/boxart/827563ccc.jpg	Grand Theft Auto: Vice City	PS2	Action	Rockstar Games	Rockstar North	9.6	16.15	8.41	0.47	5.49	1.78
3	/games/boxart/full_9218923AmericaFrontccc.jpg	Grand Theft Auto V	X360	Action	Rockstar Games	Rockstar North	NaN	15.86	9.06	0.06	5.33	1.42
4	/games/boxart/full_4990510AmericaFrontccc.jpg	Call of Duty: Black Ops 3	PS4	Shooter	Activision	Treyarch	8.1	15.09	6.18	0.41	6.05	2.44
...
64010	/games/boxart/full_2294305JapanFrontccc.jpg	World End Syndrome	PS4	Visual Novel	Arc System Works	Arc System Works	NaN	0.00	0.00	0.00	0.00	0.00

```
#Libraries for plotting diagrams and graphs
```

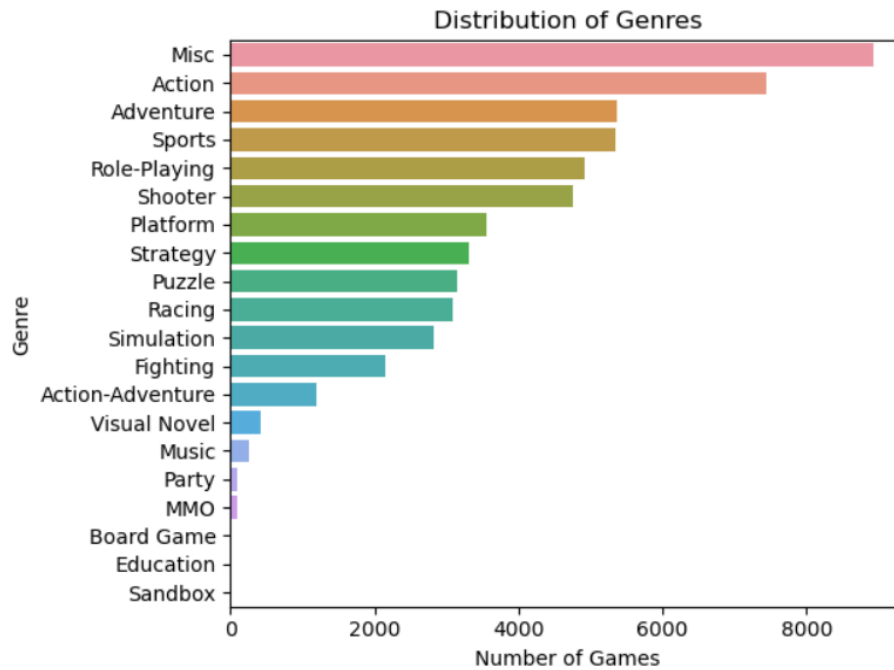
```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
#Top 10 games by sales
```

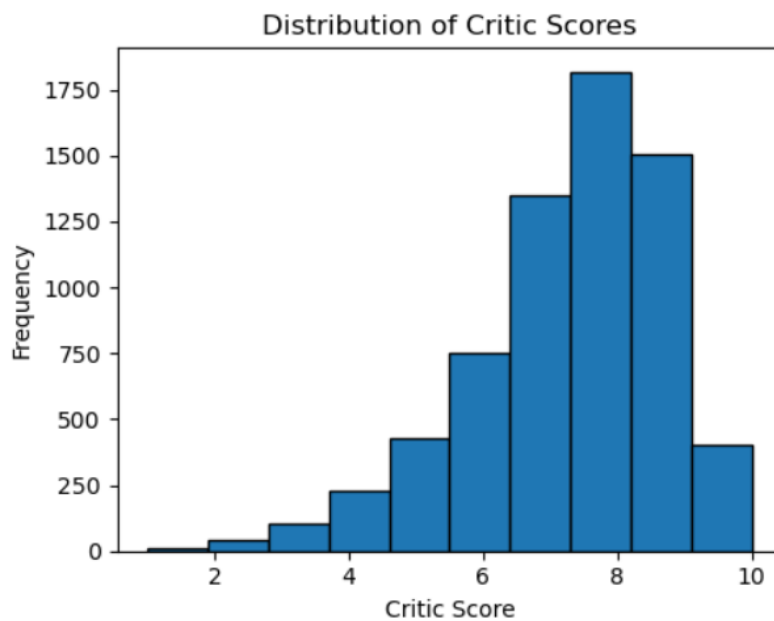
```
top_games = df[['title', 'total_sales', 'console', 'publisher']].sort_values(by='total_sales', ascending=False).head(10)
print(top_games)
```

	title	total_sales	console	publisher
0	Grand Theft Auto V	20.32	PS3	Rockstar Games
1	Grand Theft Auto V	19.39	PS4	Rockstar Games
2	Grand Theft Auto: Vice City	16.15	PS2	Rockstar Games
3	Grand Theft Auto V	15.86	X360	Rockstar Games
4	Call of Duty: Black Ops 3	15.09	PS4	Activision
5	Call of Duty: Modern Warfare 3	14.82	X360	Activision
6	Call of Duty: Black Ops	14.74	X360	Activision
7	Red Dead Redemption 2	13.94	PS4	Rockstar Games
8	Call of Duty: Black Ops II	13.86	X360	Activision
9	Call of Duty: Black Ops II	13.80	PS3	Activision

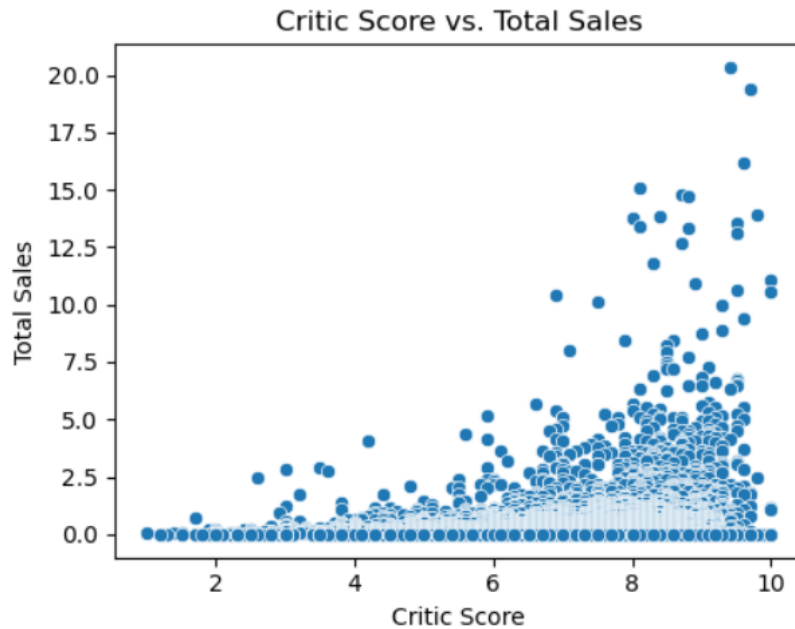

```
#Distribution of Genres
sns.countplot(data=df, y='genre', order=df['genre'].value_counts().index)
plt.title('Distribution of Genres')
plt.xlabel('Number of Games')
plt.ylabel('Genre')
plt.tight_layout()
plt.show()
```



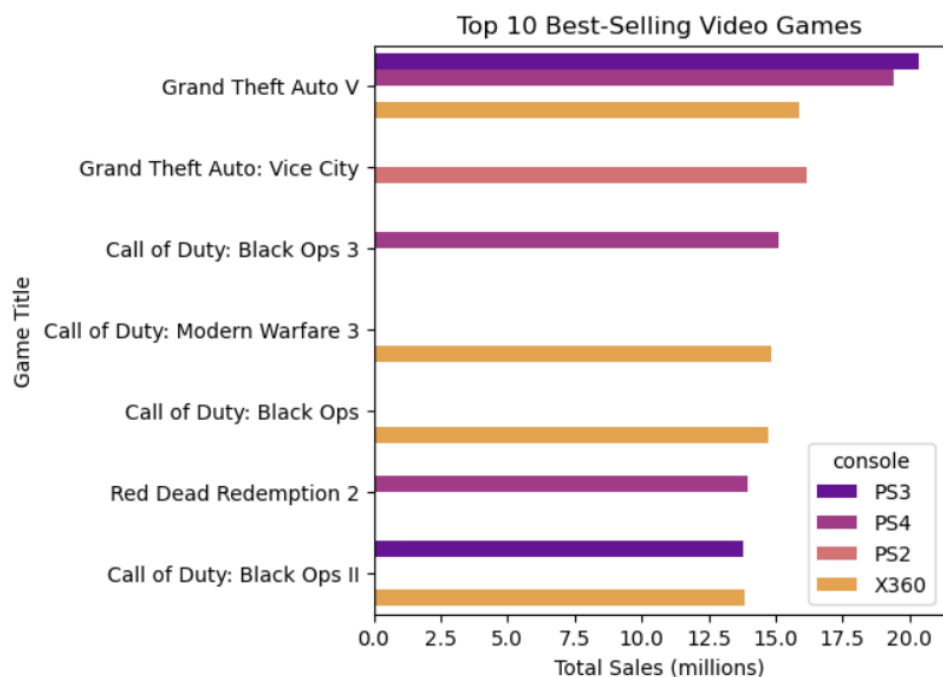
```
# Histogram of critic scores
df['critic_score'].plot(kind='hist', bins=10, edgecolor='black', figsize=(5, 4))
plt.title('Distribution of Critic Scores')
plt.xlabel('Critic Score')
plt.ylabel('Frequency')
plt.tight_layout()
plt.show()
```



```
# Scatter plot of critic score vs. total sales
plt.figure(figsize=(5, 4))
sns.scatterplot(data=df, x='critic_score', y='total_sales')
plt.title('Critic Score vs. Total Sales')
plt.xlabel('Critic Score')
plt.ylabel('Total Sales')
plt.tight_layout()
plt.show()
```



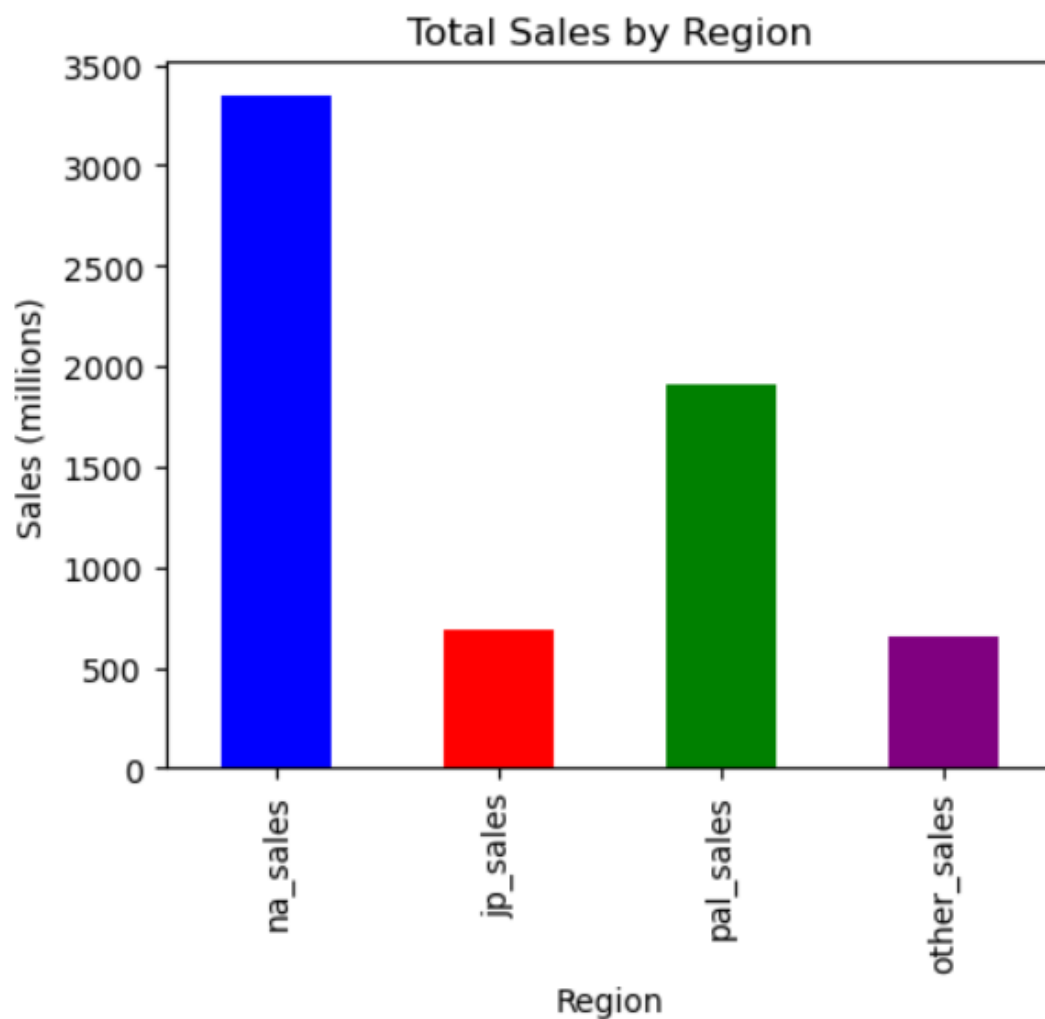
```
#Top 10 best selling games
plt.figure(figsize=(5, 5))
sns.barplot(data=top_games, x='total_sales', y='title', hue='console', palette='plasma')
plt.title("Top 10 Best-Selling Video Games")
plt.xlabel("Total Sales (millions)")
plt.ylabel("Game Title")
plt.show()
```



```
region_sales = df[['na_sales', 'jp_sales', 'pal_sales', 'other_sales']].sum()
print(region_sales)
```

```
na_sales      3344.54
jp_sales       684.99
pal_sales     1914.64
other_sales    650.49
dtype: float64
```

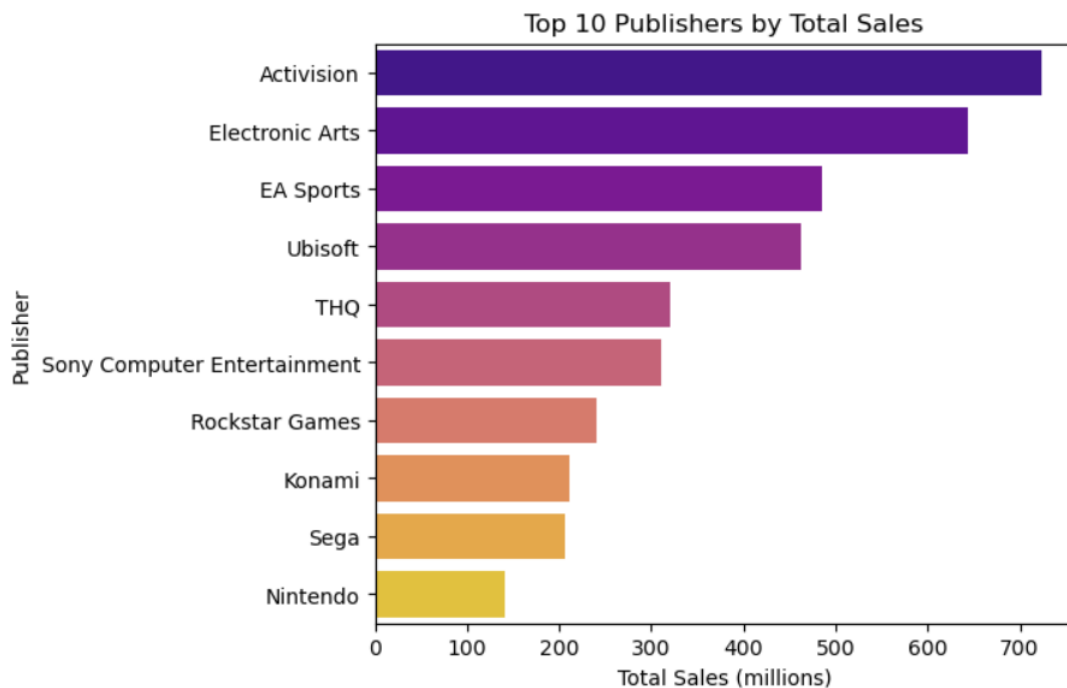
```
#Total sales by Region
plt.figure(figsize=(5, 4))
region_sales.plot(kind='bar', color=['blue', 'red', 'green', 'purple'])
plt.title("Total Sales by Region")
plt.ylabel("Sales (millions)")
plt.xlabel("Region")
plt.show()
```



```
top_publishers = df.groupby('publisher')['total_sales'].sum().sort_values(ascending=False).head(10)
print(top_publishers)
```

```
publisher
Activision          722.77
Electronic Arts     644.13
EA Sports           485.66
Ubisoft             462.43
THQ                 320.89
Sony Computer Entertainment 311.08
Rockstar Games      239.67
Konami              210.70
Sega                206.38
Nintendo            140.80
Name: total_sales, dtype: float64
```

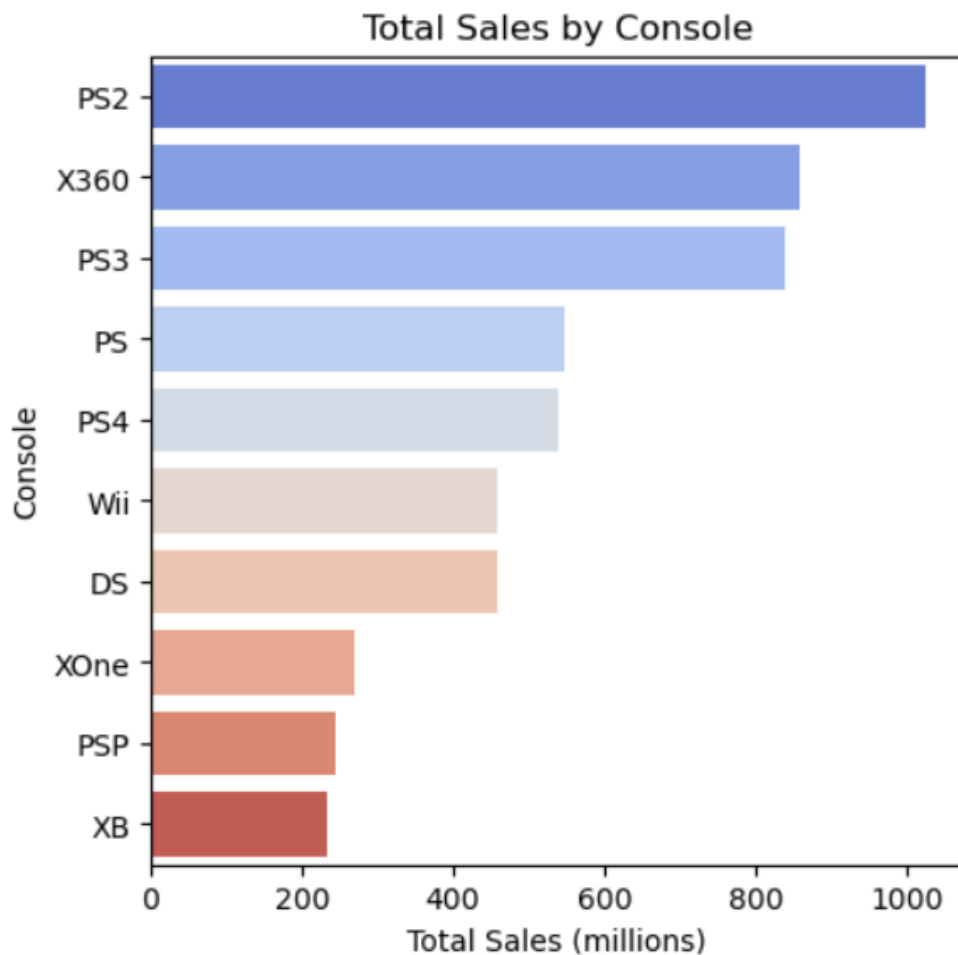
```
#Top 10 publisher by total sales
plt.figure(figsize=(6, 5))
sns.barplot(x=top_publishers.values, y=top_publishers.index, palette="plasma")
plt.title("Top 10 Publishers by Total Sales")
plt.xlabel("Total Sales (millions)")
plt.ylabel("Publisher")
plt.show()
```



```
console_sales = df.groupby('console')['total_sales'].sum().sort_values(ascending=False).head(10)
print(console_sales)
```

```
console
PS2      1025.38
X360      859.41
PS3       839.01
PS        546.21
PS4       539.92
Wii       458.92
DS         457.54
XOne       268.96
PSP        244.74
XB         232.05
Name: total_sales, dtype: float64
```

```
#Total sales by console
plt.figure(figsize=(5, 5))
sns.barplot(x=console_sales.values, y=console_sales.index, palette="coolwarm")
plt.title("Total Sales by Console")
plt.xlabel("Total Sales (millions)")
plt.ylabel("Console")
plt.show()
```

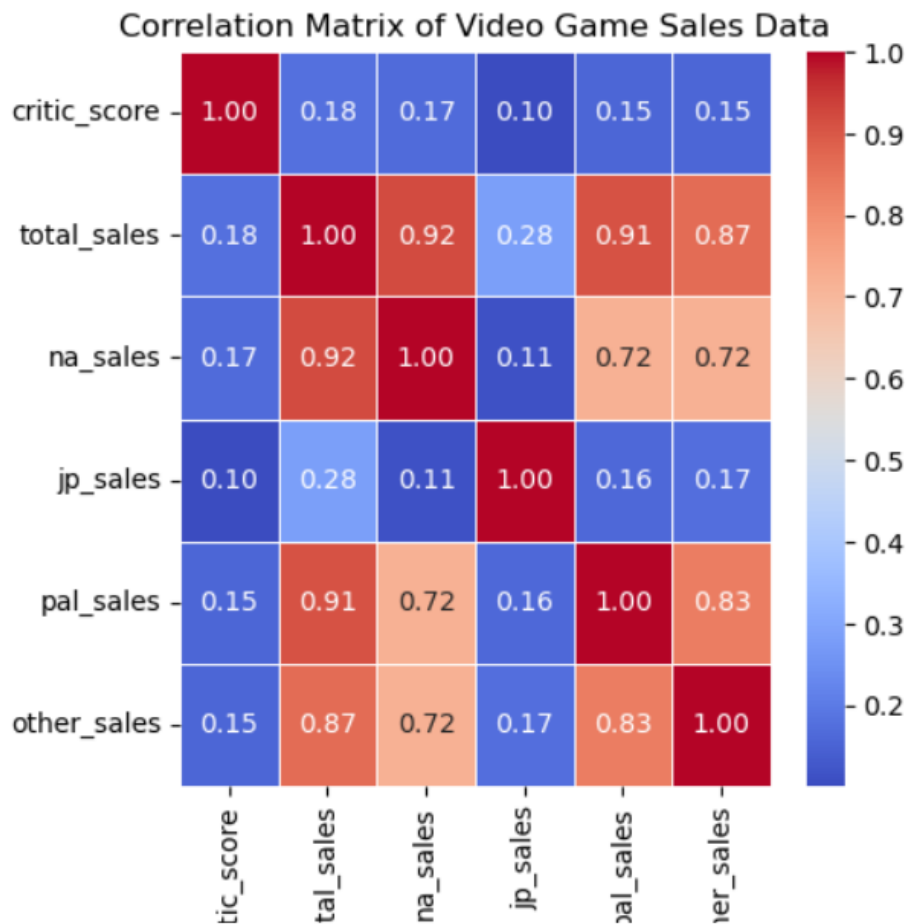


```
columns_to_drop = ['img', 'title', 'release_date', 'last_update', 'console', 'genre', 'publisher', 'developer']
df_dropped = df.drop(columns=columns_to_drop)
print("\nDataFrame after dropping specified columns:")
df_dropped.head()
```

DataFrame after dropping specified columns:

	critic_score	total_sales	na_sales	jp_sales	pal_sales	other_sales
0	9.4	20.32	6.37	0.99	9.85	3.12
1	9.7	19.39	6.06	0.60	9.71	3.02
2	9.6	16.15	8.41	0.47	5.49	1.78
3	NaN	15.86	9.06	0.06	5.33	1.42
4	8.1	15.09	6.18	0.41	6.05	2.44

```
#Correlation Matrix of Video Game Sales Data
correlation_matrix = df_dropped.corr()
plt.figure(figsize=(5, 5))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
plt.title('Correlation Matrix of Video Game Sales Data')
plt.show()
```



```
#Importing Libraries for wordcloud
from wordcloud import WordCloud
# Function to generate word cloud
def generate_wordcloud(text_series, title):
    text = ' '.join(text_series.dropna().astype(str))
    wordcloud = WordCloud(width=800, height=400, background_color='white').generate(text)
    plt.figure(figsize=(8, 5))
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis('off')
    plt.title(title)
    plt.tight_layout()
    plt.show()

# Word cloud for genres
generate_wordcloud(df['genre'], 'Common Game Genres')

# Word cloud for publishers
generate_wordcloud(df['publisher'], 'Common Game Publishers')

# Word cloud for developers
generate_wordcloud(df['developer'], 'Common Game Developers')
```


[illegible][illegible][illegible]

```
best_seller = df.groupby('title')['total_sales'].sum().sort_values(ascending=False).head(10)
print(best_seller)
```

```
title
Grand Theft Auto V          64.29
Call of Duty: Black Ops     30.99
Call of Duty: Modern Warfare 3  30.71
Call of Duty: Black Ops II   29.59
Call of Duty: Ghosts        28.80
Call of Duty: Black Ops 3    26.72
Call of Duty: Modern Warfare 2  25.02
Minecraft                  24.01
Grand Theft Auto IV         22.53
Call of Duty: Advanced Warfare 21.78
Name: total_sales, dtype: float64
```

```
best_seller2 = df.groupby('title')['total_sales'].sum().sort_values(ascending=False)[11:21]
print(best_seller2)
```

```
title
Call of Duty: WWII          19.82
Red Dead Redemption 2       19.71
Call of Duty 4: Modern Warfare 18.33
FIFA 15                     18.03
Battlefield 3               17.32
FIFA 14                     17.31
FIFA 17                     17.02
FIFA 18                     16.92
Guitar Hero III: Legends of Rock 16.38
Grand Theft Auto: Vice City 16.19
Name: total_sales, dtype: float64
```

```
from scipy import stats
```

```
#ANOVA Test
```

```
f_statistic, p_value = stats.f_oneway(best_seller, best_seller2)
print(f"F-statistic: {f_statistic}")
print(f"P-value: {p_value}")
```

```
F-statistic: 10.53986771760108
P-value: 0.004479657072584212
```

```
#T Test
```

```
ps4_sales = df[df['console'] == 'PS4']['total_sales']
ps3_sales = df[df['console'] == 'PS3']['total_sales']
tvalue, pvalue = stats.ttest_ind(ps3_sales, ps4_sales)
print(f"T-value: {tvalue}, P-value: {pvalue}")
```

```
T-value: 6.154334474099901, P-value: 8.312841588399693e-10
```

```
#ANOVA Test
```

```
f_statistic, p_value = stats.f_oneway(ps3_sales, ps4_sales)
print(f"F-statistic: {f_statistic}")
print(f"P-value: {p_value}")
```

```
F-statistic: 37.875832819094505
P-value: 8.312841588393562e-10
```



```

#Chi-square Test
contingency_table = pd.crosstab(df['console'], df['genre'])
chi2_statistic, p_value, dof, expected_frequencies = stats.chi2_contingency(contingency_table)

print(f"Chi2 Statistic: {chi2_statistic:.2f}")
print(f"P-value: {p_value:.3f}")
print("Degrees of Freedom:", dof)
print("Expected Frequencies:\n", expected_frequencies)

alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis: There is a significant association between console and genre.")
else:
    print("Fail to reject the null hypothesis: There is no significant association between console and genre.")

```

```

Chi2 Statistic: 32517.51
P-value: 0.000
Degrees of Freedom: 1482
Expected Frequencies:
[[6.33478772e+01 1.01336166e+01 4.57119781e+01 ... 4.54820557e+01
 2.81271728e+01 3.61063314e+00]
 [3.95760965e+01 6.33089862e+00 2.85582049e+01 ... 2.84145626e+01
 1.75722337e+01 2.25571514e+00]
 [1.46026653e+02 2.33595533e+01 1.05373178e+02 ... 1.04843172e+02
 6.48374829e+01 8.32306774e+00]
 ...
 [1.55430874e+01 2.48639253e+00 1.12159286e+01 ... 1.11595147e+01
 6.90130632e+00 8.85907926e-01]
 [2.48166942e+00 3.96987042e-01 1.79077852e+00 ... 1.78177125e+00
 1.10188924e+00 1.41447484e-01]
 [9.40422095e+00 1.50437195e+00 6.78610809e+00 ... 6.75197528e+00
 4.17558029e+00 5.36011518e-01]]
Reject the null hypothesis: There is a significant association between console and genre.

```

CONCLUSION

Several important insights into the elements affecting the industry's commercial performance have been uncovered by this examination of the video game sales dataset.

- **Data Characteristics:** The dataset includes 64,016 video game titles with 14 features. It includes both numerical (such as sales numbers and critic ratings) and categorical (such as genre, publisher, and console) data. Careful data cleaning and preparation are required because a sizable amount of the data, especially sales numbers and critic scores, have missing values.
- **Sales Distribution:** A few number of best-selling titles dominate the market, and the distribution of overall sales is severely skewed. One game that performs exceptionally well on a variety of platforms is Grand Theft Auto V.
- **Regional Sales:** While Japan and other regions have lesser market shares, North America and the PAL region make significant contributions to overall sales.
- **Performance of publishers and systems:** A number of publishers, including Activision and Electronic Arts, as well as systems such as the PS2 and X360, have shown impressive overall sales results.
- **Analysis of Genre:** Among the most popular gaming genres are action and adventure games.
- **Critic ratings:** The dataset generally shows a trend toward good evaluations, with critic ratings concentrated in the higher range. Total sales and critic scores have a weakly positive association.

Statistical Relationships:

- The genre of the game and the gaming system are statistically significantly correlated.
- The average reviewer evaluations for action and sports games differ statistically significantly.
- The difference between PS3 and PS4 system sales is statistically significant.