

```
import pandas as pd
import numpy as np

data = pd.read_csv("Video_Games_Sales_as_at_22_Dec_2016.csv")

Out[2]:
   NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  Critic_Score  Critic_Count  User_Score  User_Count  Developer  Rating
0           Wii Sports      Wii      2009.0      Sports      Nintendo      41.36      28.96      3.77      8.45      82.53      76.0      51.0      8.0      322.0  Nintendo  E
1           Super Mario Bros. NES      1985.0      Platform      Nintendo      29.08      15.78      6.81      0.77      40.24      NaN      NaN      NaN      NaN      NaN      NaN
2           Mario Kart Wii      Wii      2009.0      Racing      Nintendo      15.68      12.76      3.79      3.29      35.52      42.0      73.0      8.3      709.0  Nintendo  E
3           Wii Sports Resort      Wii      2009.0      Sports      Nintendo      15.61      10.93      3.28      2.95      36.77      80.0      73.0      8.0      192.0  Nintendo  E
4           Pokemon Red/Pokemon Blue GB      1996.0      Role-Playing      Nintendo      11.27      6.89      10.22      1.00      31.37      NaN      NaN      NaN      NaN      NaN      NaN
...
16714 Samurai Warriors: Sanada Maru PS3      2016.0      Action      Tecmo Koei      0.00      0.00      0.00      0.00      0.01      NaN      NaN      NaN      NaN      NaN      NaN
16715 LMA Manager 2007 X360      2006.0      Sports      Codemasters      0.00      0.01      0.00      0.00      0.01      NaN      NaN      NaN      NaN      NaN      NaN
16716 Hatake no Psychodokka PS3      2016.0      Adventure      Idea Factory      0.00      0.00      0.01      0.00      0.01      NaN      NaN      NaN      NaN      NaN      NaN
16717 Spitz & Spels GBA      2003.0      Platform      Wanadoo      0.01      0.00      0.00      0.00      0.01      NaN      NaN      NaN      NaN      NaN      NaN
16718 Winning Post 2016 PSV      2016.0      Simulation      Tecmo Koei      0.00      0.00      0.01      0.00      0.01      NaN      NaN      NaN      NaN      NaN      NaN
16719 rows x 16 columns

In [3]:
data.info()
Out[3]:
Out[4]:
data.describe()
Out[5]:
data.fillna(0, inplace=True)
In [6]:
data.info()
Out[6]:
data.describe()
In [7]:
data.fillna(0, inplace=True)
In [8]:
data.info()
Out[8]:
data.describe()
In [9]:
data['Platform'].unique()
Out[9]:
array(['Wii', 'NES', 'GB', 'DS', 'X360', 'PS3', 'PSP', 'WES', 'GBA', 'PS4', '3DS', 'N64', 'PS', 'XB', 'PC', '2608', 'PSP', 'XOne', 'WiiU', 'GC', 'GBN', 'DC', 'PSV', 'SAT', 'SCD', 'WS', 'NS', 'TG16', '3DS', 'GC', 'PCCX'], dtype=object)
In [9]:
data['Genre'].unique()
Out[10]:
array(['Sports', 'Platform', 'Racing', 'Role-Playing', 'Puzzle', 'Misc', 'Shooter', 'Simulation', 'Action', 'Fighting', 'Adventure', 'Strategy', ''], dtype=object)
In [13]:
import pandas as pd
import matplotlib.pyplot as plt

# Select the columns for the scatter plots
x_na_sales = data['NA_Sales']
y_na_sales = data['Global_Sales']

x_eu_sales = data['EU_Sales']
y_eu_sales = data['Global_Sales']

x_jp_sales = data['JP_Sales']
y_jp_sales = data['Global_Sales']

x_other_sales = data['Other_Sales']
y_other_sales = data['Global_Sales']

# Plot the scatter plots
plt.figure(figsize=(15, 8))

plt.subplot(221)
plt.scatter(x_na_sales, y_na_sales)
plt.xlabel('NA Sales')
plt.ylabel('Global Sales')
plt.title('Scatter Plot: NA Sales vs Global Sales')

plt.subplot(222)
plt.scatter(x_eu_sales, y_eu_sales)
plt.xlabel('EU Sales')
plt.ylabel('Global Sales')
plt.title('Scatter Plot: EU Sales vs Global Sales')

plt.subplot(223)
plt.scatter(x_jp_sales, y_jp_sales)
plt.xlabel('JP Sales')
plt.ylabel('Global Sales')
plt.title('Scatter Plot: JP Sales vs Global Sales')

plt.subplot(224)
plt.scatter(x_other_sales, y_other_sales)
plt.xlabel('Other Sales')
plt.ylabel('Global Sales')
plt.title('Scatter Plot: Other Sales vs Global Sales')

plt.tight_layout()
plt.show()

Scatter Plot: NA Sales vs Global Sales
Scatter Plot: EU Sales vs Global Sales
Scatter Plot: JP Sales vs Global Sales
Scatter Plot: Other Sales vs Global Sales

In [14]:
import pandas as pd
import matplotlib.pyplot as plt

# Load the dataset
data = pd.read_csv("Video_Games_Sales_as_at_22_Dec_2016.csv")

# Group the data by genre and calculate the sum of global sales
genre_sales = data.groupby('Genre')['Global_Sales'].sum().reset_index()

# Sort the genres based on global sales in descending order
genre_sales = genre_sales.sort_values('Global_Sales', ascending=False)

# Plot the bar chart
plt.bar(genre_sales['Genre'], genre_sales['Global_Sales'])
plt.xlabel('Genre')
plt.ylabel('Global Sales')
plt.title('Genre vs Global Sales')
plt.xticks(rotation=50) # Rotate x-axis labels for better visibility
plt.show()

Genre vs Global Sales

In [15]:
import pandas as pd
import matplotlib.pyplot as plt

# Group the data by genre and calculate the sum of global sales
genre_sales = data.groupby('Genre')['JP_Sales'].sum().reset_index()

# Sort the genres based on global sales in descending order
genre_sales = genre_sales.sort_values('JP_Sales', ascending=False)

# Plot the bar chart
plt.bar(genre_sales['Genre'], genre_sales['JP_Sales'])
plt.xlabel('Genre')
plt.ylabel('Japan Sales')
plt.title('Genre vs Japan Sales')
plt.xticks(rotation=50) # Rotate x-axis labels for better visibility
plt.show()

Genre vs Japan Sales

In [16]:
# Group the data by genre and calculate the sum of global sales
genre_sales = data.groupby('Genre')['NA_Sales'].sum().reset_index()

# Sort the genres based on global sales in descending order
genre_sales = genre_sales.sort_values('NA_Sales', ascending=False)

# Plot the bar chart
plt.bar(genre_sales['Genre'], genre_sales['NA_Sales'])
plt.xlabel('Genre')
plt.ylabel('North America Sales')
plt.title('Genre vs North America Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

Genre vs North America Sales

In [17]:
# Group the data by genre and calculate the sum of global sales
genre_sales = data.groupby('Genre')['EU_Sales'].sum().reset_index()

# Sort the genres based on global sales in descending order
genre_sales = genre_sales.sort_values('EU_Sales', ascending=False)

# Plot the bar chart
plt.bar(genre_sales['Genre'], genre_sales['EU_Sales'])
plt.xlabel('Genre')
plt.ylabel('Europe Sales')
plt.title('Genre vs Europe Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

Genre vs Europe Sales

In [18]:
# Group the data by genre and calculate the sum of global sales
genre_sales = data.groupby('Genre')['Other_Sales'].sum().reset_index()

# Sort the genres based on global sales in descending order
genre_sales = genre_sales.sort_values('Other_Sales', ascending=False)

# Plot the bar chart
plt.bar(genre_sales['Genre'], genre_sales['Other_Sales'])
plt.xlabel('Genre')
plt.ylabel('Other Sales')
plt.title('Genre vs Other Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

Genre vs Other Sales

In [19]:
# Group the data by year and calculate the sum of other sales
year_sales_others = data.groupby('Year_of_Release')['Other_Sales'].sum().reset_index()

# Sort the year based on other sales in descending order
year_sales_others = year_sales_others.sort_values('Other_Sales', ascending=False)

# Group the data by year and calculate the sum of EUROPE sales
year_sales_eu = data.groupby('Year_of_Release')['EU_Sales'].sum().reset_index()

# Sort the year based on other sales in descending order
year_sales_eu = year_sales_eu.sort_values('EU_Sales', ascending=False)

# Group the data by year and calculate the sum of JAPAN sales
year_sales_jp = data.groupby('Year_of_Release')['JP_Sales'].sum().reset_index()

# Sort the year based on other sales in descending order
year_sales_jp = year_sales_jp.sort_values('JP_Sales', ascending=False)

# Group the data by year and calculate the sum of NORTH AMERICA sales
year_sales_na = data.groupby('Year_of_Release')['NA_Sales'].sum().reset_index()

# Sort the year based on other sales in descending order
year_sales_na = year_sales_na.sort_values('NA_Sales', ascending=False)

# Plot the bar chart
plt.bar(year_sales_others['Year_of_Release'], year_sales_others['Other_Sales'])
plt.xlabel('Year')
plt.ylabel('Other Sales')
plt.title('Year vs Other Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

# Plot the bar chart
plt.bar(year_sales_eu['Year_of_Release'], year_sales_eu['EU_Sales'])
plt.xlabel('Year')
plt.ylabel('EU Sales')
plt.title('Year vs EU Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

# Plot the bar chart
plt.bar(year_sales_jp['Year_of_Release'], year_sales_jp['JP_Sales'])
plt.xlabel('Year')
plt.ylabel('JAPAN Sales')
plt.title('Year vs JAPAN Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

# Plot the bar chart
plt.bar(year_sales_na['Year_of_Release'], year_sales_na['NA_Sales'])
plt.xlabel('Year')
plt.ylabel('NA Sales')
plt.title('Year vs NA Sales')
plt.xticks(rotation=90) # Rotate x-axis labels for better visibility
plt.show()

Year vs Other Sales
Year vs EU Sales
Year vs JAPAN Sales
Year vs NA Sales

In [20]:
data = data.drop(['Genre', 'Platform', 'Year_of_Release', 'Name', 'Publisher', 'Developer', 'Rating'], axis=1)
In [21]:
data.fillna(0, inplace=True)
In [32]:
import seaborn as sns
from sklearn.model_selection import train_test_split

y = data['Global_Sales']

# Load X variables into a Pandas DataFrame with columns
X = data.drop(['Global_Sales'], axis=1)

# y is dependent variable and X is independent variable
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=100, test_size=0.3)

In [33]:
cor = X_train.corr()
plt.figure(figsize=(10, 10))
sns.heatmap(cor, cmap=plt.cm.CmRmap_r, annot=True)
plt.show()

NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales  Critic_Score  Critic_Count  User_Score  User_Count
NA_Sales  1.00  0.77  0.42  0.63  0.15  0.23  0.14  0.19
EU_Sales  0.77  1.00  0.44  0.69  0.16  0.24  0.15  0.25
JP_Sales  0.42  0.44  1.00  0.19  0.04  0.03  0.02  0.02
Other_Sales  0.63  0.69  0.19  1.00  0.16  0.24  0.15  0.22
Global_Sales  0.15  0.16  0.04  0.16  1.00  0.76  0.82  0.25
Critic_Score  0.23  0.24  0.03  0.24  0.76  1.00  0.69  0.19
Critic_Count  0.14  0.15  0.02  0.15  0.82  0.69  1.00  0.12
User_Score  0.19  0.25  0.02  0.22  0.25  0.19  0.12  1.00
User_Count  0.19  0.25  0.02  0.22  0.25  0.19  0.12  1.00

In [34]:
def correlation(dataset, threshold):
    cor_corr = list()
    cor_data = dataset.corr()
    for i in range(len(cor_data.columns)):
        for j in range(i):
            if abs(cor_data.iloc[i, j]) > threshold:
                colname = cor_data.columns[i]
                cor_corr.append(colname)
    return cor_corr

cor_features = correlation(X_train, 0.7)
cor_features

['Critic_Count', 'EU_Sales', 'User_Score']

In [36]:
X_train.drop(cor_features, axis=1)
X_test.drop(cor_features, axis=1)
X_train
X_test

Out[36]:
   NA_Sales  EU_Sales  JP_Sales  Other_Sales  Critic_Score  Critic_Count  User_Score  User_Count
1162      1.15      0.36      0.04      0.04      48.0      50.0      5.0      53.0
13175     0.05      0.00      0.00      0.00      60.0      10.0      5.8      9.0
1669      1.74      0.45      0.00      0.18      42.0      10.0      5.5      16.0
16495     0.00      0.01      0.00      0.00      0.00      0.0      0.0      0.0
...
16304     0.00      0.00      0.01      0.00      0.0      0.0      0.0      0.0
79      1.75      3.47      2.49      0.07      66.0      40.0      7.4      54.0
12119     0.06      0.00      0.00      0.00      65.0      16.0      6.3      4.0
14147     0.00      0.00      0.00      0.00      0.0      0.0      0.0      0.0
5640      0.27      0.03      0.00      0.03      79.0      80.0      7.3      117.0
11703 rows x 8 columns

In [37]:
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
y_train = label_encoder.fit_transform(y_train)

In [38]:
#CHI-square
#feature chi2 test
from sklearn.feature_selection import chi2
#calculating feature and p values
f_p_values=chi2(X_train,y_train)

Out[38]:
(array([2.87376486e+04, 2.87428781e+04, 1.02933661e+04, 0.22899808e+03,
        6.03993209e+02, 8.48248989e+04, 1.4273744e+03, 0.6061339e+05]),
array([0., 0., 0., 0., 0., 0., 0., 0.]))

In [39]:
import pandas as pd
p_values=pd.Series(f_p_values[1])
p_values.index=X_train.columns
p_values

Out[39]:
NA_Sales      0.0
EU_Sales      0.0
JP_Sales      0.0
Other_Sales    0.0
Critic_Score   0.0
Critic_Count   0.0
User_Score     0.0
User_Count     0.0
dtype: float64

In [40]:
p_values.sort_index(ascending=False)

In [40]:
User_Score      0.0
User_Count      0.0
Other_Sales     0.0
NA_Sales        0.0
JP_Sales        0.0
EU_Sales        0.0
Critic_Score    0.0
Critic_Count    0.0
dtype: float64

In [41]:
#Mutual information gain
#exporting mutual information gain
from sklearn.feature_selection import mutual_info_classif
# determine the mutual information
mutual_info = mutual_info_classif(X_train, y_train)
mutual_info

Out[41]:
array([1.07806838, 0.63053717, 0.54841337, 0.65351435, 0.0642521 ,
        0.06215876, 0.84514413, 0.8562982 ]))

In [42]:
#exporting in list form
mutual_info = list(mutual_info)
mutual_info.index = X_train.columns
mutual_info.sort_values(ascending=False)

Out[42]:
NA_Sales      1.078068
Other_Sales    0.653534
EU_Sales      0.630537
JP_Sales      0.548413
Critic_Score  0.064252
Critic_Count  0.062158
User_Count    0.062089
User_Score    0.045144
dtype: float64

In [43]:
mutual_info.sort_values(ascending=False).plot.bar(figsize=(20, 8))

In [43]:
<AxesSubplot>

In [44]:
columns = ['NA_Sales', 'Other_Sales', 'EU_Sales', 'JP_Sales', 'Critic_Score', 'Critic_Count']
X = data[columns]
y = data['Global_Sales']

In [45]:
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)
print(X.shape, X_train.shape, X_test.shape)
(16719, 6) (13375, 6) (3344, 6)

In [46]:
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.metrics import accuracy_score

In [47]:
#fit sklearn.linear_model import LinearRegression
# Create and fit the linear regression model
model = LinearRegression()
model.fit(X_train, Y_train)

# Calculate the R-squared score for training data
train_score = model.score(X_train, Y_train)
print("Training score:", train_score)

# Calculate the R-squared score for testing data
test_score = model.score(X_test, Y_test)
print("Testing score:", test_score)
Training score: 0.999993286959928
Testing score: 0.999984843679464

In [48]:
y_pred = model.predict(X_test)

In [48]:
#Mean Squared Error (MSE)
mse = mean_squared_error(Y_test, y_pred)
print("Mean Squared Error:", mse)
Mean Squared Error: 2.6647849826636276e-05

In [48]:
#Mean Absolute Error (MAE)
from sklearn.metrics import mean_absolute_error
# Assuming you have the true target values in y_true and the predicted values in y_pred
mae = mean_absolute_error(Y_test, y_pred)
print("Mean Absolute Error:", mae)
Mean Absolute Error: 0.0628386184298871547

In [51]:
from sklearn.metrics import r2_score

# Assuming you have the true target values in y_true and the predicted values in y_pred
r2 = r2_score(Y_test, y_pred)
print("R-squared:", r2)
R-squared: 0.999984843679464

In [ ]:
```