```
1 #Import packages
2 import numpy as np
3 import pandas as pd
4 from sklearn.tree import export_graphviz
5 import matplotlib.pyplot as plt
6 from sklearn.tree import plot_tree
7 import seaborn as sns
8 import plotly.express as px
9 from sklearn.model_selection import cross_val_score
10 from sklearn.model_selection import train_test_split
11 from sklearn.preprocessing import OneHotEncoder
12 from sklearn.ensemble import RandomForestRegressor
13 from sklearn.metrics import mean_absolute_error, r2_score
14 from matplotlib.patches import Rectangle
```

## ∨ Data Cleaning

```
1 #Import the dataset
2 data = pd.read_csv("vgsales.csv")
3 #Check the NA values
4 data.isnull().sum()
5 #Drop NA values and rename the dataset without NA values
6 data = data.dropna()
```
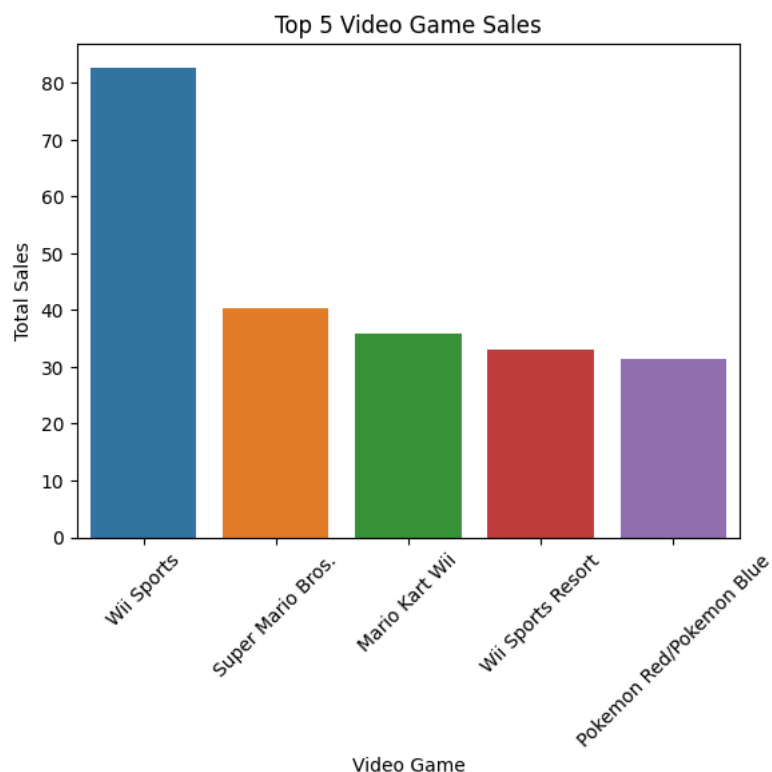
∨ **Question 1:** What are the top 5 video games with the most sales and in which region did they sell the most?

```
1 # Verifying Top 5 Video Games
2 top_five = data[:5]
3 print(top_five)
```

```
   Rank                     Name Platform   Year       Genre Publisher  \
0     1               Wii Sports      Wii 2006.0      Sports  Nintendo
1     2        Super Mario Bros.      NES 1985.0    Platform  Nintendo
2     3           Mario Kart Wii      Wii 2008.0      Racing  Nintendo
3     4         Wii Sports Resort      Wii 2009.0      Sports  Nintendo
4     5  Pokemon Red/Pokemon Blue       GB 1996.0 Role-Playing  Nintendo

   NA_Sales  EU_Sales  JP_Sales  Other_Sales  Global_Sales
0     41.49     29.02      3.77         8.46         82.74
1     29.08      3.58      6.81         0.77         40.24
2     15.85     12.88      3.79         3.31         35.82
3     15.75     11.01      3.28         2.96         33.00
4     11.27      8.89     10.22         1.00         31.37
```

```
1 # Top 5 Video Games
2 top_five = data[:5]
3 plot = sns.barplot(data=top_five, x='Name', y='Global_Sales', hue='Name')
4 plot.tick_params(axis="x", rotation=45)
5 plot.set(title='Top 5 Video Game Sales', xlabel='Video Game', ylabel='Total Sales')
6 plt.show()
```
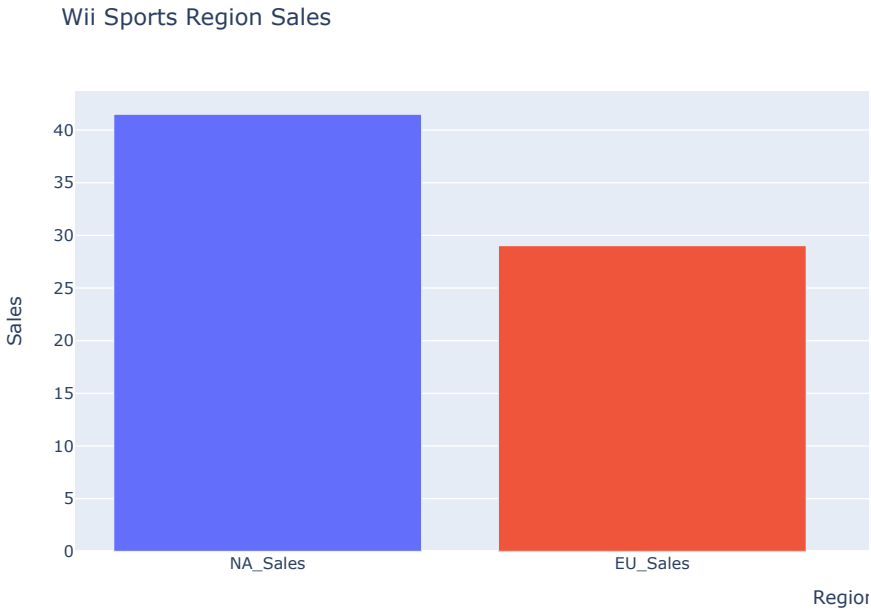


```
1 # Verifying Region Sale Values
2 regions = data[['Name','NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 game_region = regions[:5]
4 print(game_region)
```

```
                Name  NA_Sales  EU_Sales  JP_Sales  Other_Sales
0         Wii Sports     41.49     29.02      3.77         8.46
1  Super Mario Bros.     29.08      3.58      6.81         0.77
2     Mario Kart Wii     15.85     12.88      3.79         3.31
```
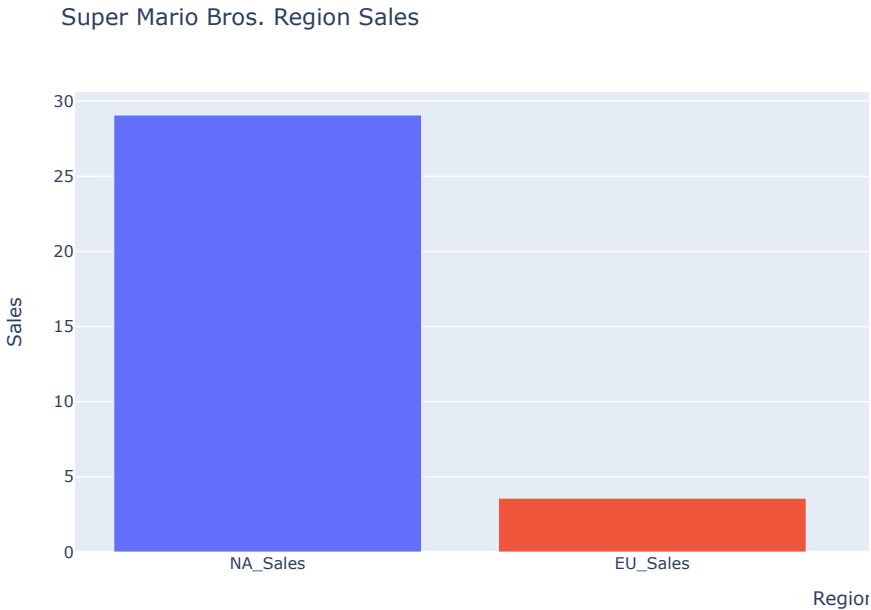
```
3         Wii Sports Resort     15.75    11.01    3.28    2.96
4   Pokemon Red/Pokemon Blue    11.27     8.89   10.22    1.00
```
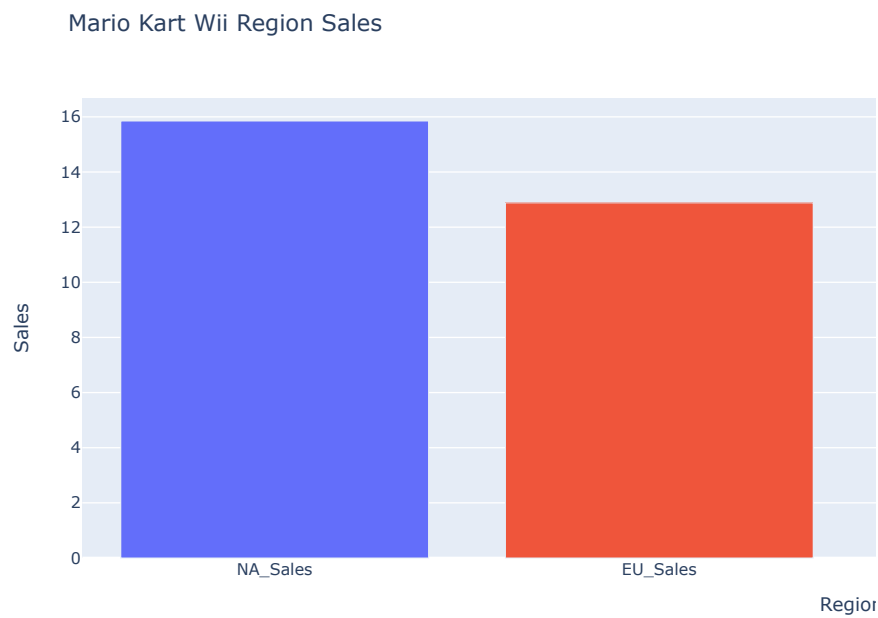
```python
1 # Wii Sports Region with most sales
2 regions = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 wii_region = regions[:1]
4
5 # making region sales into rows
6 melted_data_1 = wii_region.melt(
7         var_name='Region',
8         value_name='Sales',
9         ignore_index=False)
10
11 plot = px.bar(melted_data_1, x='Region', y='Sales', color='Region', title='Wii Sports Region Sales')
12 plot
```
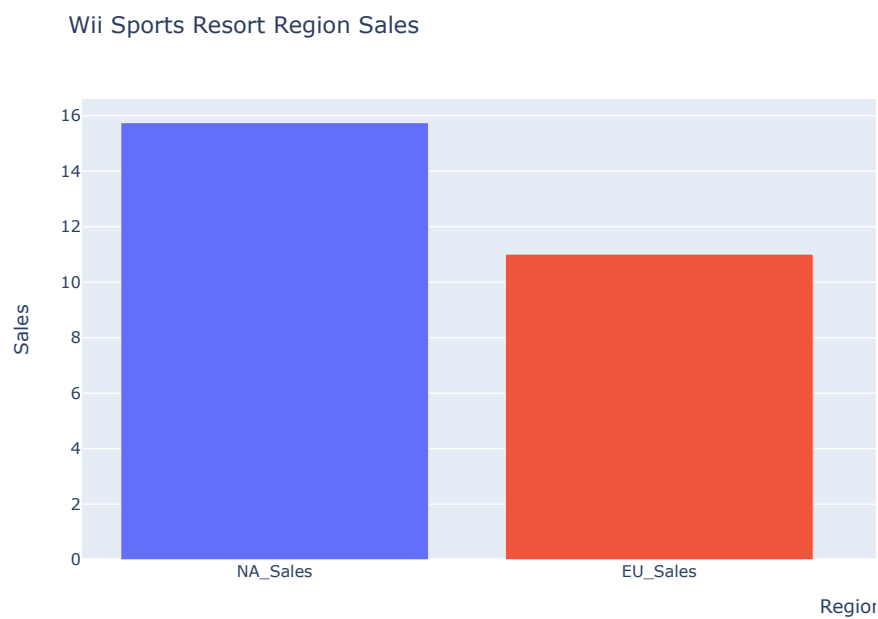


Wii Sports Region Sales

```python
1 # Super Mario Bros. Region with most sales
2 regions = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 mario_region = regions[1:2]
4
5 # making region sales into rows
6 melted_data_2 = mario_region.melt(
7         var_name='Region',
8         value_name='Sales',
9         ignore_index=False)
10
11 plot = px.bar(melted_data_2, x='Region', y='Sales', color='Region', title='Super Mario Bros. Region Sales')
12 plot
```



Super Mario Bros. Region Sales

```
1 # Mario Kart Wii Region with the most sales
2 regions = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 kart_region = regions[2:3]
4
5 # making region sales into rows
6 melted_data_3 = kart_region.melt(
7         var_name='Region',
8         value_name='Sales',
9         ignore_index=False)
10
11 plot = px.bar(melted_data_3, x='Region', y='Sales', color='Region', title='Mario Kart Wii Region Sales')
12 plot
```
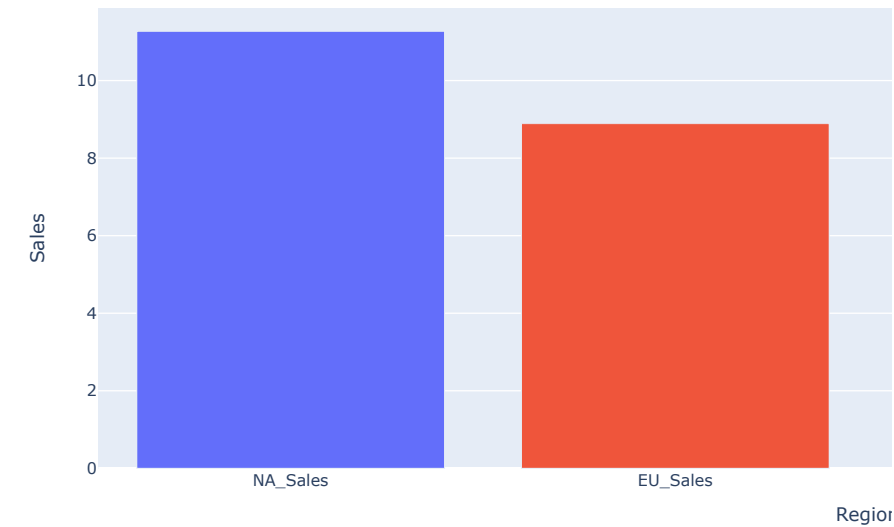
Mario Kart Wii Region Sales



```
1 # Wii Sports Resort Region with the most sales
2 regions = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 resort_region = regions[3:4]
4
5 # making region sales into rows
6 melted_data_4 = resort_region.melt(
7         var_name='Region',
8         value_name='Sales',
9         ignore_index=False)
10
11 plot = px.bar(melted_data_4, x='Region', y='Sales', color='Region', title='Wii Sports Resort Region Sales')
12 plot
```

Wii Sports Resort Region Sales



```
1 # Pokemon Red/Pokemon Blue Region with the most sales
2 regions = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 pokemon_region = regions[4:5]
4
5 # making region sales into rows
6 melted_data_5 = pokemon_region.melt(
7         var_name='Region',
8         value_name='Sales',
9         ignore_index=False)
10
11 plot = px.bar(melted_data_5, x='Region', y='Sales', color='Region', title='Pokemon Red/Pokemon Blue Region Sales')
12 plot
```
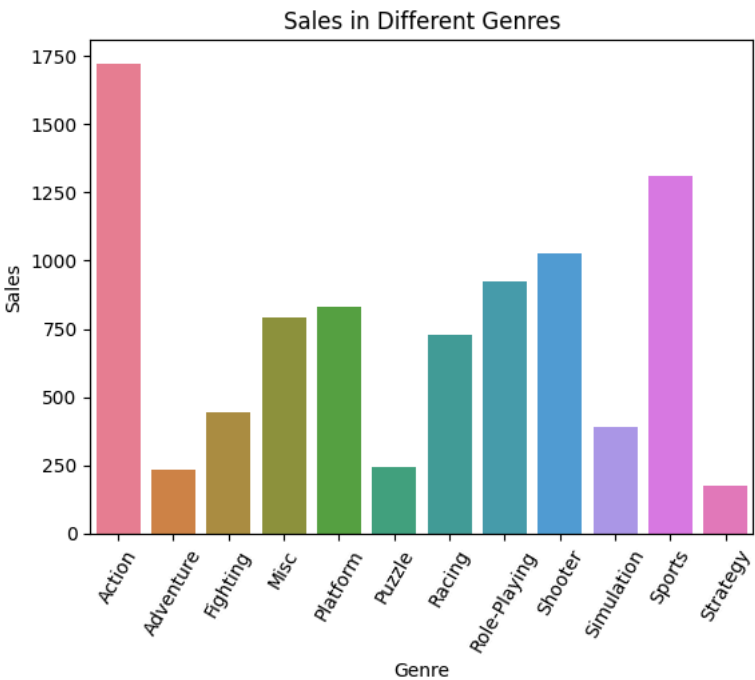
Pokemon Red/Pokemon Blue Region Sales



**Question 2:** Which genre of video game has the highest sales figures and what potential factors contribute to its appeal?

```
1 #sales of each genre
2 genre_data = data.groupby('Genre')
3 sales = genre_data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
4 total_sales = sales.sum()
5 total_sales['Sales'] = genre_data['Global_Sales'].sum()
6
7
8 bar = sns.barplot(total_sales, x='Genre', y = 'Sales', hue='Genre')
9 bar.tick_params(axis="x", rotation=60)
10 bar.set(title='Sales in Different Genres', xlabel='Genre', ylabel= 'Sales')

    [Text(0.5, 1.0, 'Sales in Different Genres'),
     Text(0.5, 0, 'Genre'),
     Text(0, 0.5, 'Sales')]
```



**Question 3:** Have the genres with the highest sales maintained consistent popularity over time?

```
1 #Top three most popular genre
2 sort = total_sales.sort_values(by='Sales', ascending=False)
3 top_three = sort.head(3)
4 top_three
```
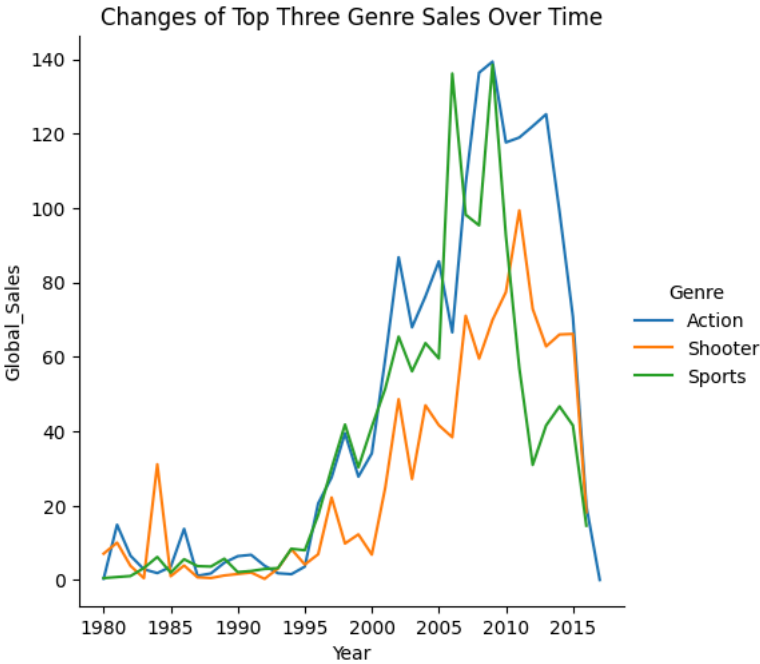
| Genre | NA_Sales | EU_Sales | JP_Sales | Other_Sales | Sales |
|---|---|---|---|---|---|
| Action | 861.77 | 516.48 | 158.65 | 184.92 | 1722.84 |
| Sports | 670.09 | 371.34 | 134.76 | 132.65 | 1309.24 |
| Shooter | 575.16 | 310.45 | 38.18 | 101.90 | 1026.20 |

Next steps:   Generate code with `top_three`      ◯ View recommended plots

```
1 #changes over years
2 top_three_genre = data.loc[data['Genre'].isin(['Action', 'Sports', 'Shooter'])]
3 group = top_three_genre.groupby(['Year', 'Genre']).sum()
4 line = sns.relplot(group, x="Year", y="Global_Sales", hue="Genre", kind="line")
5 line.set(title='Changes of Top Three Genre Sales Over Time')
6 line
```

    <ipython-input-13-9e11047ea8d6>:3: FutureWarning:

    The default value of numeric_only in DataFrameGroupBy.sum is deprecated. In a fu

    <seaborn.axisgrid.FacetGrid at 0x7f3f2f494f70>



Changes of Top Three Genre Sales Over Time

**Sales Region Analysis**: Analyze the factors (platform, genre, publisher, and year of release) in influencing global sales

### Data Preparation

```
1 # Convert 'Year' to integer
2 data['Year'] = data['Year'].astype(int)
3 # Initialize the OneHotEncoder
4 encoder = OneHotEncoder(sparse=False, handle_unknown='ignore')
5 # Transform the 'Platform', 'Genre', 'Publisher' columns, and convert it to a DataFrame
6 trasformed_data = encoder.fit_transform(data[['Platform', 'Genre', 'Publisher']])
7 feature_names = encoder.get_feature_names_out(['Platform', 'Genre', 'Publisher'])
8 encoded_df = pd.DataFrame(trasformed_data, columns=feature_names)
9 # Concatenate the encoded_df with 'Year' and 'Global_Sales' from the original data
10 new_data = pd.concat([data[['Year', 'Global_Sales']].reset_index(drop=True), encoded_df], axis=1)
11 # Define features X and target y using the new_data DataFrame
12 X = new_data.drop(['Global_Sales'], axis=1)
13 y = new_data['Global_Sales']
14 # Split the data into training and testing sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

    /usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_encoders.py:868: FutureWarning:

    `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you lea

### Train the model, Predict, and Evaluate

```
1 # Initialize and train the RandomForestRegressor
2 model = RandomForestRegressor(n_estimators = 20,random_state=42)
3 model.fit(X_train, y_train)
4 # Predict on the test set
5 y_predict = model.predict(X_test)
6 # Evaluate the model
7 MAE = mean_absolute_error(y_test, y_predict)
8 R2 = r2_score(y_test, y_predict)
9 scores = cross_val_score(model, X_train, y_train, cv=5, scoring='r2')
10 print(f"Mean Absolute Error: {MAE}")
11 print(f"R^2 Score: {R2}")
12 print(f"Average cross-validated R^2: {np.mean(scores)}")
```
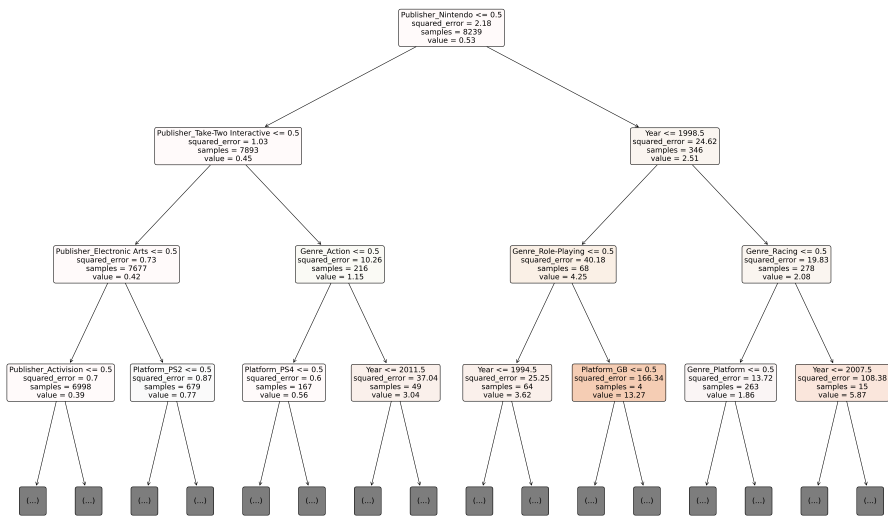
    Mean Absolute Error: 0.5232569474055581
    R^2 Score: 0.025417878866786925
    Average cross-validated R^2: -0.04068050592557313

### Feature Importance Analysis

```
1 # Feature Importance
2 feature_importances = model.feature_importances_
3 features = X.columns
4 # Combine feature names with their importance scores
5 feature_importance_dictionary = dict(zip(features, feature_importances))
6 # Sort the feature importances in descending order
7 sorted_feature_importances = sorted(feature_importance_dictionary.items(), key=lambda item: item[1], reverse=True)
8 #Total Importance
9 total_importance = sum(feature_importances)
10 # Display the top 10 most important features
11 print("Most Important Features in Percentage:")
12 for feature, importance in sorted_feature_importances[:10]:
13     print(f"{feature}: {importance / total_importance * 100:.2f}%")
```

```
Most Important Features in Percentage:
Year: 42.05%
Publisher_Nintendo: 10.30%
Genre_Role-Playing: 4.89%
Genre_Shooter: 4.04%
Genre_Platform: 2.72%
Genre_Racing: 2.67%
Genre_Action: 2.65%
Genre_Misc: 2.64%
Platform_GB: 2.59%
Platform_Wii: 2.26%
```

```
1 # Plot the top levels of the tree with a limited max_depth
2 single_tree = model.estimators_[0]
3
4 plt.figure(figsize=(30, 20), dpi=480)  # Set a large figure size and high DPI for better resolution
5 plot_tree(single_tree,
6           feature_names=X.columns,
7           filled=True,
8           rounded=True,
9           precision=2,
10          fontsize=14,
11          max_depth=3)
12 plt.show()
```



## Test

```
1 #test 1
2 def test_na_value(df):
3   """
4   This function called test_na_value take dataframe as an input and
5   return whether there is NA value in the dataset
6
7   >>> test_na_value(data)
8   False
```

```
 8
 9     """
10     if df.isna().any().any():
11        return True
12     else:
13        return False
14
15 #test 2
16 assert sorted((rectangle.get_height() for rectangle in bar.findobj(Rectangle)[:2]),
17            reverse=True) == [1722.84, 234.59], "Data does not match expected"
18
19 #test 3
20 assert isinstance(line, sns.axisgrid.FacetGrid), "Failed: Plot not created successfully"
21 ax = line.ax
22 assert ax.get_title() == "Changes of Top Three Genre Sales Over Time", "title does not match expected"
23 assert all(line.get_xydata().size == 0 for line in ax.get_lines()[37:]), "unexpected extra data"
24
25 #test 4
26 expected = len(data['Platform'].unique()) + len(data['Genre'].unique()) + len(data['Publisher'].unique()) + 1
27 actual = encoded_df.shape[1] + 1
28 assert actual == expected, f"The number of columns after encoding does not match the expected value)."
29
30 #test 5
31 def test_importance_scores(model, features):
32     """
33     Test if the feature importance scores from a trained model sum up to 1 (100%).
34
35     >>> test_importance_scores(model, X.columns)
36     (True, 'The importance scores is equal to one 1 (100%).')
37     """
38     feature_importances = model.feature_importances_
39
40     if np.isclose(np.sum(feature_importances), 1.0, atol=1e-6):
41         return True, "The importance scores is equal to one 1 (100%)."
42     else:
43         return False, f"The importance scores sum up to {np.sum(feature_importances):.2f} instead of 1."
```

```
1 import doctest
2 doctest.testmod()
```

```
TestResults(failed=0, attempted=2)
```