```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

    Mounted at /content/drive

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
1 # from google.colab import files
2
3
4 # uploaded = files.upload()
```

```
1 data = pd.read_csv("/content/drive/ MyDrive/BIProject/vgsales.csv")
2 data.head(5)
```

| | Name | Platform | Year_of_Release | Genre | Publisher | NA_Sales | EU_Sales | JP_S |
|---|---|---|---|---|---|---|---|---|
| 0 | Wii Sports | Wii | 2006.0 | Sports | Nintendo | 41.36 | 28.96 | |
| 1 | Super Mario Bros. | NES | 1985.0 | Platform | Nintendo | 29.08 | 3.58 | |
| 2 | Mario Kart Wii | Wii | 2008.0 | Racing | Nintendo | 15.68 | 12.76 | |
| 3 | Wii Sports Resort | Wii | 2009.0 | Sports | Nintendo | 15.61 | 10.93 | |
| 4 | Pokemon Red/Pokemon Blue | GB | 1996.0 | Role-Playing | Nintendo | 11.27 | 8.89 | 1 |

```
1 print(data.isnull().sum())
```

```
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
dtype: int64
```

```
1 data = data.dropna()
2 data.dtypes
3 data['User_Score'] = data['User_Score'].apply(pd.to_numeric)
4 data.dtypes
```

```
Name               object
Platform           object
Year_of_Release    float64
Genre              object
Publisher          object
NA_Sales           float64
EU_Sales           float64
JP_Sales           float64
Other_Sales        float64
Global_Sales       float64
Critic_Score       float64
Critic_Count       float64
User_Score         float64
User_Count         float64
Developer          object
Rating             object
dtype: object
```

```
1 data.describe()
```

|  | Year_of_Release | NA_Sales | EU_Sales | JP_Sales | Other_Sales | Global_Sales |
|---|---|---|---|---|---|---|
| count | 6825.000000 | 6825.000000 | 6825.000000 | 6825.000000 | 6825.000000 | 6825.000000 |
| mean | 2007.436777 | 0.394484 | 0.236089 | 0.064158 | 0.082677 | 0.777590 |
| std | 4.211248 | 0.967385 | 0.687330 | 0.287570 | 0.269871 | 1.963443 |
| min | 1985.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.010000 |
| 25% | 2004.000000 | 0.060000 | 0.020000 | 0.000000 | 0.010000 | 0.110000 |
| 50% | 2007.000000 | 0.150000 | 0.060000 | 0.000000 | 0.020000 | 0.290000 |
| 75% | 2011.000000 | 0.390000 | 0.210000 | 0.010000 | 0.070000 | 0.750000 |
| max | 2016.000000 | 41.360000 | 28.960000 | 6.500000 | 10.570000 | 82.530000 |

```
1 data['Genre'].value_counts()
```

```
Action           1630
Sports            943
Shooter           864
```

```
Role-Playing      712
Racing            581
Platform          403
Misc              384
Fighting          378
Simulation        297
Strategy          267
Adventure         248
Puzzle            118
Name: Genre, dtype: int64
```
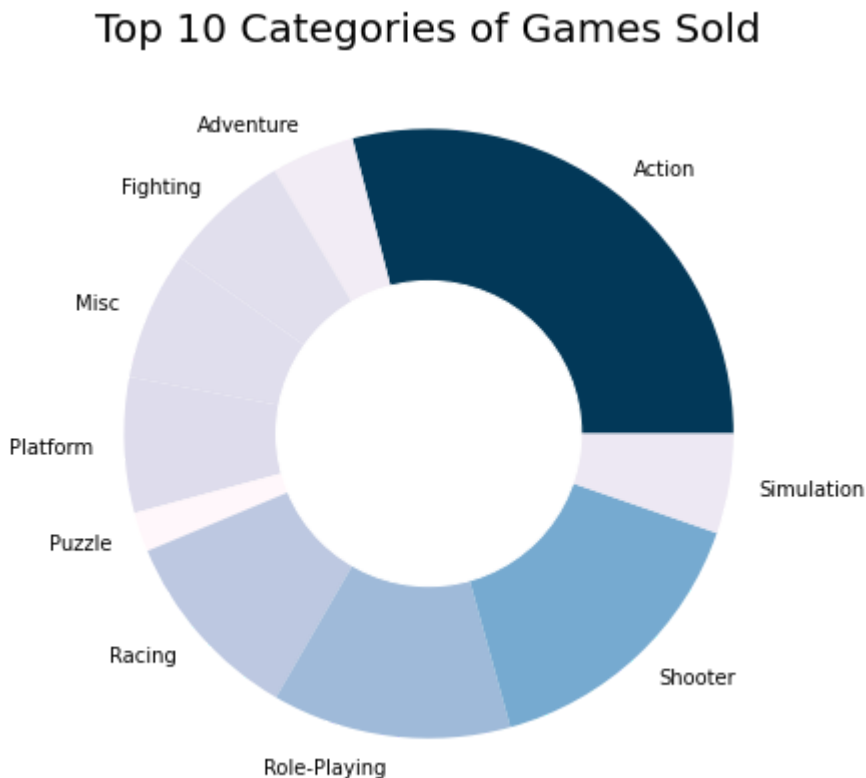
```
1  import matplotlib as mpl
2  game = data.groupby("Genre")["Global_Sales"].count().head(10)
3  custom_colors = mpl.colors.Normalize(vmin=min(game), vmax=max(game))
4  colours = [mpl.cm.PuBu(custom_colors(i)) for i in game]
5  plt.figure(figsize=(7,7))
6  plt.pie(game, labels=game.index, colors=colours)
7  central_circle = plt.Circle((0, 0), 0.5, color='white')
8  fig = plt.gcf()
9  fig.gca().add_artist(central_circle)
10 plt.rc('font', size=12)
11 plt.title("Top 10 Categories of Games Sold", fontsize=20)
12 plt.show()
```

## Top 10 Categories of Games Sold



```
1  data_platform = data.groupby(by=['Platform'])['Global_Sales'].sum()
2  data_platform = data_platform.reset_index()
3  data_platform = data_platform.sort_values(by=['Global_Sales'], ascending=False)
```
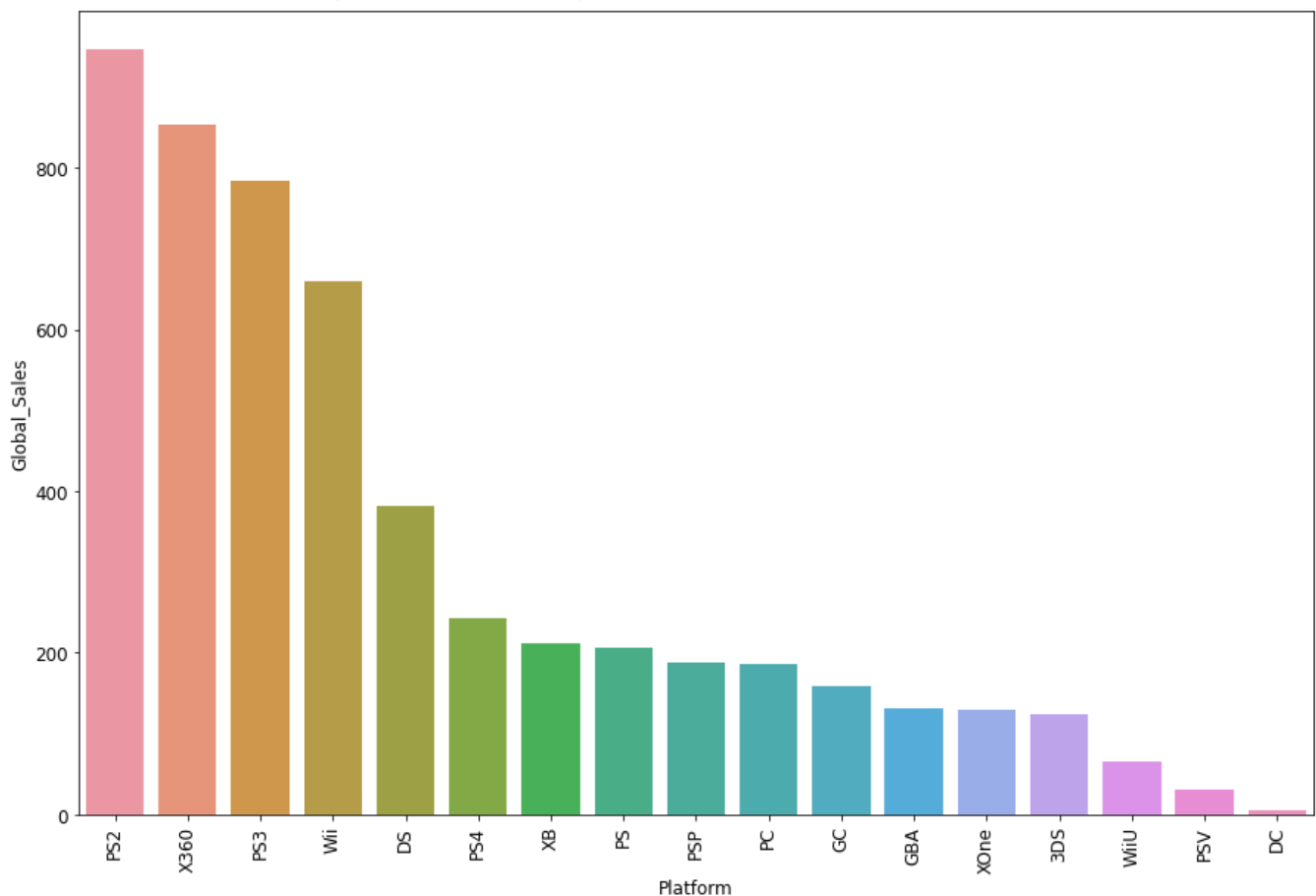
```
1 plt.figure(figsize=(15, 10))
2 sns.barplot(x="Platform", y="Global_Sales", data=data_platform)
3 plt.xticks(rotation=90)
```

```
(array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16]),
 <a list of 17 Text major ticklabel objects>)
```
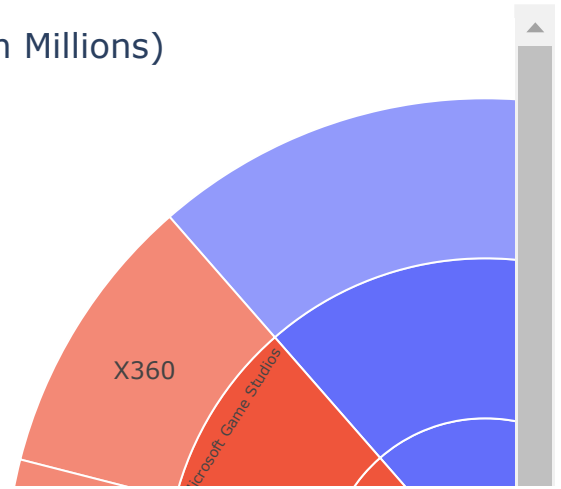


```
1 from plotly import express as px
2 top_sales = data.sort_values(by=['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales'], ascen
3
4 # ['NA_Sales', '', '', '']
5 dicts_name = {
6     'NA_Sales' : "North America Sales ( In Millions)",
7     'EU_Sales' : "Europe Sales ( In Millions)",
8     'JP_Sales' : "Japan Sales ( In Millions)",
9     'Other_Sales' : "Other Sales ( In Millions)",
10 }
11
```

```
12  for (key, title) in dicts_name.items():
13
14      fig = px.sunburst(top_sales, path=['Genre', 'Publisher', 'Platform'], values=key, titl
15
16      fig.update_layout(
17          grid= dict(columns=2, rows=2),
18          margin = dict(t=40, l=2, r=2, b=5)
19      )
20
21      fig.show()
```

## Top Selling by North America Sales ( In Millions)



```
1
2 comp_genre = data[['Genre', 'NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
3 # comp_genre
4 comp_map = comp_genre.groupby(by=['Genre']).sum()
5 # comp_map
```
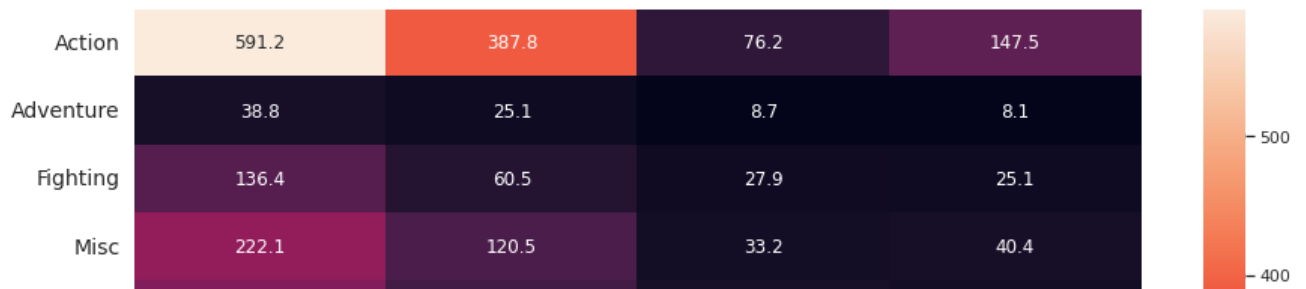


```
1 plt.figure(figsize=(15, 10))
2 sns.set(font_scale=1)
3 sns.heatmap(comp_map, annot=True, fmt = '.1f')
4
5 plt.xticks(fontsize=14)
6 plt.yticks(fontsize=14)
7 plt.show()
```

| | | | |
|---|---|---|---|
| Action | 591.2 | 387.8 | 76.2 | 147.5 |
| Adventure | 38.8 | 25.1 | 8.7 | 8.1 |
| Fighting | 136.4 | 60.5 | 27.9 | 25.1 |
| Misc | 222.1 | 120.5 | 33.2 | 40.4 |

```
1 rating = data[['Rating', 'NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
2 # comp_genre
3 comp_rat = rating.groupby(by=['Rating']).sum()
```

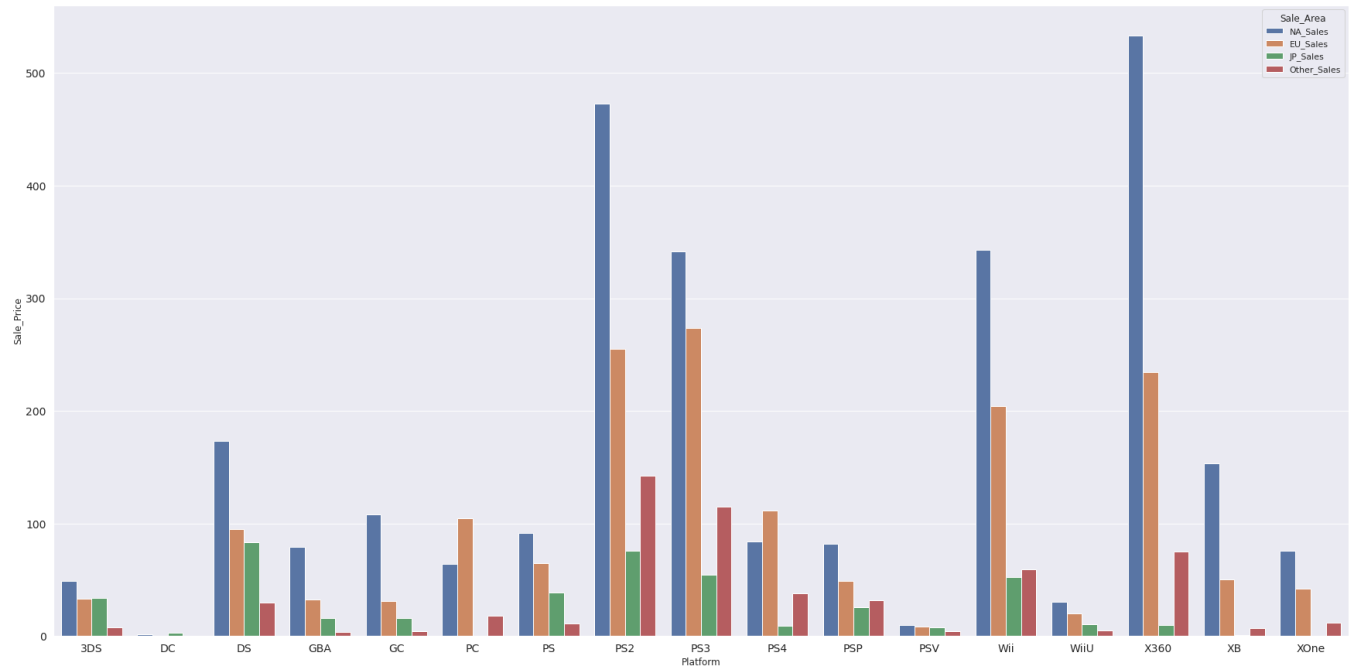| | | | |
|---|---|---|---|
| Racing | 225.6 | 164.7 | 27.8 | 58.3 |

```
1 plt.figure(figsize=(15, 10))
2 sns.set(font_scale=1)
3 sns.heatmap(comp_rat, annot=True, fmt = '.1f')
4
5 plt.xticks(fontsize=14)
6 plt.yticks(fontsize=14)
7 plt.show()
```

```
1  comp_platform = data[['Platform', 'NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
2  comp_platform.head()
3  comp_platform = comp_platform.groupby(by=['Platform']).sum().reset_index()
4  comp_table = pd.melt(comp_platform, id_vars=['Platform'], value_vars=['NA_Sales', 'EU_Sale
5  comp_table.head()
6  plt.figure(figsize=(30, 15))
7  sns.barplot(x='Platform', y='Sale_Price', hue='Sale_Area', data=comp_table)
8  plt.xticks(fontsize=14)
9  plt.yticks(fontsize=14)
10 plt.show()
```

```
1 top_publisher = data.groupby(by=['Publisher'])['Year_of_Release'].count().sort_values(asce
2 top_publisher = pd.DataFrame(top_publisher).reset_index()
3 plt.figure(figsize=(15, 10))
4 sns.countplot(x="Publisher", data=data, order = data.groupby(by=['Publisher'])['Year_of_Re
5 plt.xticks(rotation=90)
```
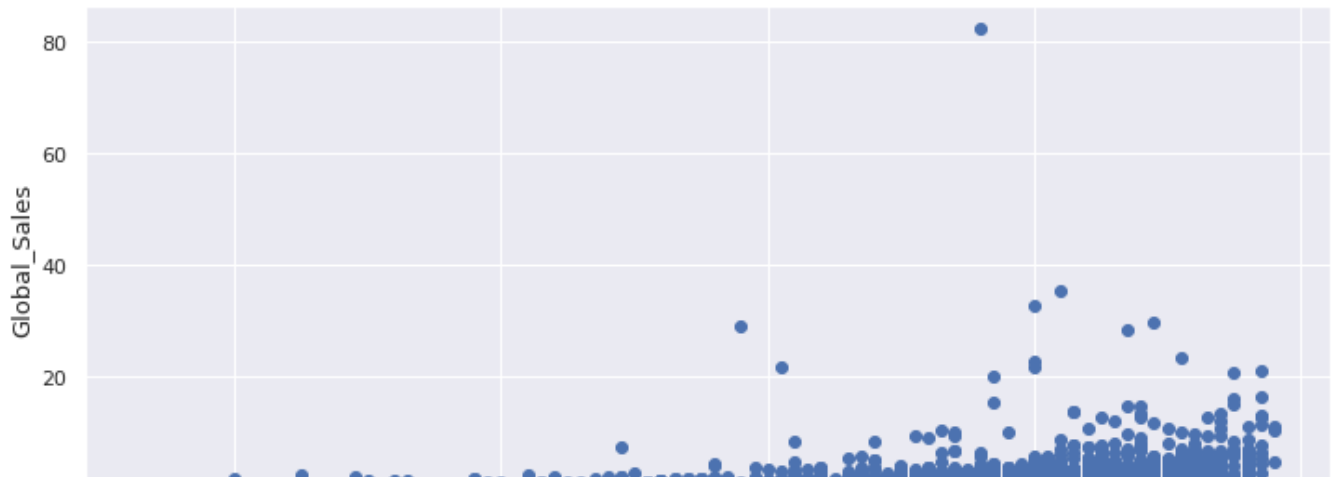
```
(array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,
        17, 18, 19]), <a list of 20 Text major ticklabel objects>)
```



```python
1 top_sale_reg = data[['NA_Sales', 'EU_Sales', 'JP_Sales', 'Other_Sales']]
2 # pd.DataFrame(top_sale_reg.sum(), columns=['a', 'b'])
3 top_sale_reg = top_sale_reg.sum().reset_index()
4 top_sale_reg = top_sale_reg.rename(columns={"index": "region", 0: "sale"})
5 top_sale_reg
6
7 labels = top_sale_reg['region']
8 sizes = top_sale_reg['sale']
9 plt.figure(figsize=(10, 8))
10 plt.pie(sizes, labels=labels, autopct='%1.1f%%', shadow=True, startangle=90)
```

```
([<matplotlib.patches.Wedge at 0x7ff1a5756a90>,
  <matplotlib.patches.Wedge at 0x7ff1a57562d0>,
  <matplotlib.patches.Wedge at 0x7ff1a574f610>,
  <matplotlib.patches.Wedge at 0x7ff1a9029590>],
```

```
1 #data = data.drop(data[data['Global_Sales']>60].index)
```

```
Text(0.8801245340527122, 0.6598339219519532, 'JP Sales'),
```

```
1 corrmat = data.corr()
2 top_corr_features = corrmat.index
3 plt.figure(figsize=(10,10))
4 #Plotting heat map
5 g=sns.heatmap(data[top_corr_features].corr(),annot=True,linewidths=.5)
6 b, t = plt.ylim() # Finding the values for bottom and top
7 b += 0.5
8 t -= 0.5
9 plt.ylim(b, t)
10 plt.show()
```

```
1 feat=['Critic_Score', 'User_Score','User_Count','Critic_Count','Year_of_Release']
2 for i in feat:
3   fig, ax = plt.subplots(1,1, figsize=(12,5))
4   ax.scatter(x = data[i], y = data['Global_Sales'])
5   plt.ylabel('Global_Sales', fontsize=13)
6   plt.xlabel(i, fontsize=13)
7   plt.show()
8   sns.regplot(x=i, y="Global_Sales", data=data,truncate=True, x_bins=15, color="#75556c").
```

```
1 def rm_outliers(df, list_of_keys):
2     df_out = df
3     for key in list_of_keys:
4         # Calculate first and third quartile
5         first_quartile = df_out[key].describe()["25%"]
6         third_quartile = df_out[key].describe()["75%"]
7
8         # Interquartile range
9         iqr = third_quartile - first_quartile
10
11        # Remove outliers
12        removed = df_out[(df_out[key] <= (first_quartile - 1.5 * iqr)) |
13                   (df_out[key] >= (third_quartile + 1.5 * iqr))]
14        df_out = df_out[(df_out[key] > (first_quartile - 1.5 * iqr)) &
15                   (df_out[key] < (third_quartile + 1.5 * iqr))]
16    return df_out, removed
```

```
1 data2, outliers = rm_outliers(data, ["Global_Sales"])
```

```
1 data2
```

| | Name | Platform | Year_of_Release | Genre | Publisher | NA_Sales | EU_Sales |
|---|---|---|---|---|---|---|---|
| **1050** | Deal or No Deal | DS | 2007.0 | Misc | Mindscape | 1.15 | 0.40 |
| **1052** | Portal 2 | PS3 | 2011.0 | Shooter | Valve | 0.83 | 0.60 |

```
1 counts = data['Publisher'].value_counts()
2
3 data['Publisher'] = data['Publisher'].apply(lambda x: 'Small Publisher' if counts[x] < 50
4
5 counts = data2['Publisher'].value_counts()
6
7 data2['Publisher'] = data2['Publisher'].apply(lambda x: 'Small Publisher' if counts[x] < 5
8
9 counts = outliers['Publisher'].value_counts()
10
11 outliers['Publisher'] = outliers['Publisher'].apply(lambda x: 'Small Publisher' if counts[
```

```
    /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:7: SettingWithCopyWarning:


    A value is trying to be set on a copy of a slice from a DataFrame.
    Try using .loc[row_indexer,col_indexer] = value instead

    See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_

    /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:11: SettingWithCopyWarning


    A value is trying to be set on a copy of a slice from a DataFrame.
    Try using .loc[row_indexer,col_indexer] = value instead

    See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_
```

```
1 features=["Critic_Score","User_Count","Critic_Count"]
2 d1=data.copy()
3 d2=data2.copy()
4 d3=outliers.copy()
5 feature = ["Platform","Genre","Rating","Publisher"]
6 y1 = data["Global_Sales"]
7 y2 = data2["Global_Sales"]
8 y3 = outliers["Global_Sales"]
9 features
```

```
    ['Critic_Score', 'User_Count', 'Critic_Count']
```

```
1 from sklearn.preprocessing import LabelEncoder
2 from sklearn.model_selection import train_test_split
```
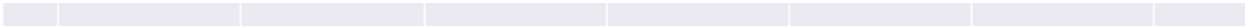
```
3 import numpy as np
4
5 le = LabelEncoder()
6 for col in feature:
7     d1[col]=le.fit_transform(d1[col])
8     d2[col]=le.fit_transform(d2[col])
9     d3[col]=le.fit_transform(d3[col])
10
```

```
1 features.extend(feature)
```

```
1 X = d1[features]
2 y = data['Global_Sales']
3 X2 = d2[features]
4 y2 = data2['Global_Sales']
5 X3 = d3[features]
6 y3 = outliers['Global_Sales']
```

```
1 X
```

| | Critic_Score | User_Count | Critic_Count | Platform | Genre | Rating | Publisher |
|---|---|---|---|---|---|---|---|
| 0 | 76.0 | 322.0 | 51.0 | 12 | 10 | 1 | 15 |
| 2 | 82.0 | 709.0 | 73.0 | 12 | 6 | 1 | 15 |
| 3 | 80.0 | 192.0 | 73.0 | 12 | 10 | 1 | 15 |
| 6 | 89.0 | 431.0 | 65.0 | 2 | 4 | 1 | 15 |
| 7 | 58.0 | 129.0 | 41.0 | 12 | 3 | 1 | 15 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 16667 | 46.0 | 21.0 | 4.0 | 3 | 0 | 1 | 19 |
| 16677 | 81.0 | 9.0 | 12.0 | 3 | 2 | 4 | 13 |
| 16696 | 80.0 | 412.0 | 20.0 | 5 | 0 | 4 | 10 |
| 16700 | 61.0 | 43.0 | 12.0 | 5 | 8 | 6 | 19 |
| 16706 | 60.0 | 13.0 | 12.0 | 5 | 11 | 2 | 19 |

6825 rows × 7 columns

```
1 from sklearn.model_selection import train_test_split
2 x_train,x_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=0)
3 x2_train,x2_test,y2_train,y2_test = train_test_split(X2,y2,test_size=0.2,random_state=0)
4 x3_train,x3_test,y3_train,y3_test = train_test_split(X3,y3,test_size=0.2,random_state=0)
```

```
1 ## Training the multiple linear regression on the training set
2 from sklearn.linear_model import LinearRegression
3 regressor_MultiLinear1 = LinearRegression()
4 regressor_MultiLinear2 = LinearRegression()
5 regressor_MultiLinear3 = LinearRegression()
6 regressor_MultiLinear1.fit(x_train,y_train)
7 regressor_MultiLinear2.fit(x2_train,y2_train)
8 regressor_MultiLinear3.fit(x3_train,y3_train)
```

```
LinearRegression()
```

```
1 y_pred1 = regressor_MultiLinear1.predict(x_test)
2 #print("for all dataset ",y_pred1)
3 y_pred2 = regressor_MultiLinear2.predict(x2_test)
4 #print("for dataset with removed outliers ",y_pred2)
5 y_pred3 = regressor_MultiLinear3.predict(x3_test)
6 #print("for outlier dataset ",y_pred3)
```

```
1 # Calculating r2 score
2 from sklearn.metrics import r2_score
3 r2_MultiLinear = r2_score(y_test,y_pred1)
4 print(r2_MultiLinear)
5 r2_MultiLinear = r2_score(y2_test,y_pred2)
6 print(r2_MultiLinear)
7 r2_MultiLinear = r2_score(y3_test,y_pred3)
8 print(r2_MultiLinear)
```

```
0.10492455277481427
0.20209660390201634
0.06377777866553735
```

```
 1 ## Finding out the optimal degree of polynomial regression
 2 from sklearn.preprocessing import PolynomialFeatures
 3 sns.set_style('darkgrid')
 4 def polyplot(xtrain,ytrain,xtest,ytest):
 5   scores_list = []
 6   pRange = range(2,6)
 7   for i in pRange :
 8       poly_reg = PolynomialFeatures(degree=i)
 9       x_poly = poly_reg.fit_transform(xtrain)
10       poly_regressor = LinearRegression()
11       poly_regressor.fit(x_poly,ytrain)
12       y_pred = poly_regressor.predict(poly_reg.fit_transform(xtest))
13       scores_list.append(r2_score(ytest,y_pred))
14   plt.plot(pRange,scores_list,linewidth=2)
15   plt.xlabel('Degree of polynomial')
16   plt.ylabel('r2 score with varying degrees')
17   plt.show()
18 polyplot(x_train,y_train ,x_test,y_test)
```
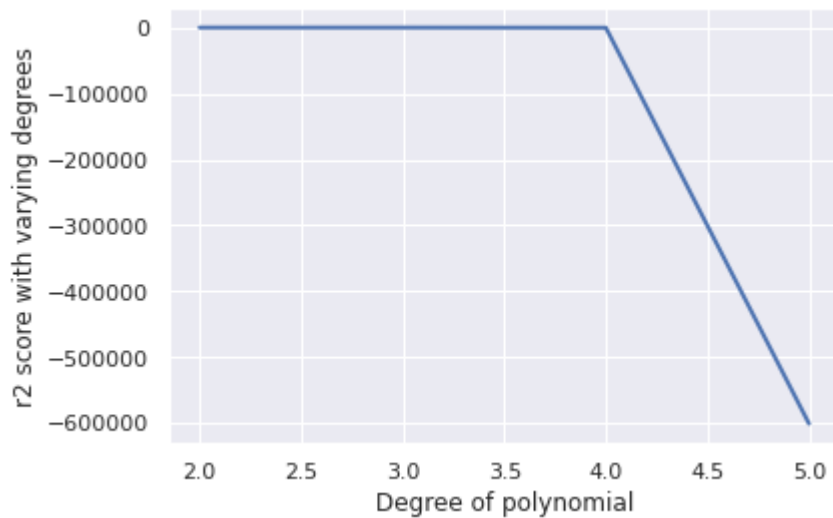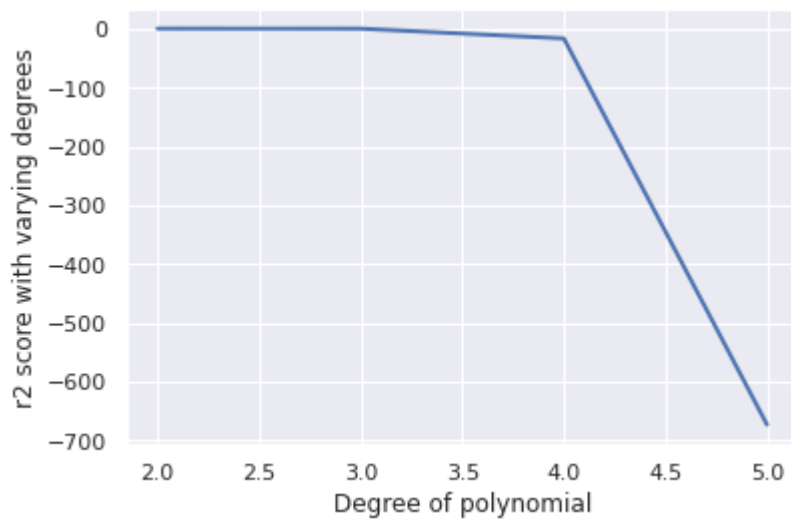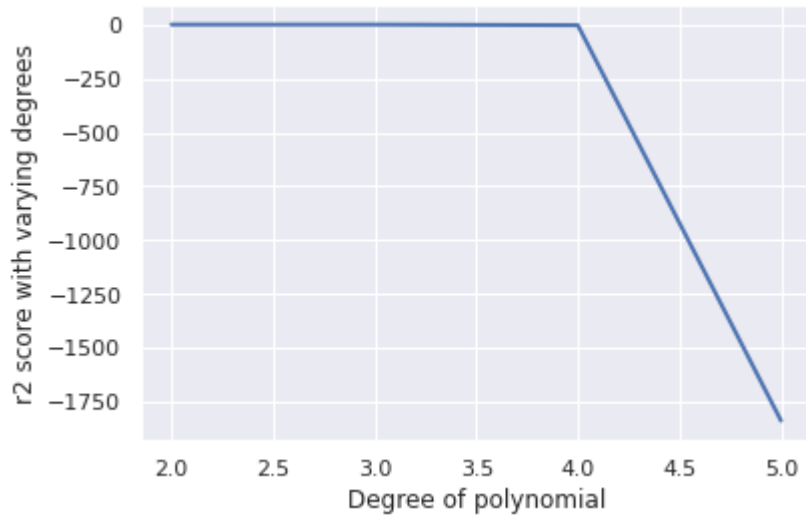
```
19 polyplot(x2_train,y2_train ,x2_test,y2_test)
20 polyplot(x3_train,y3_train ,x3_test,y3_test)
```
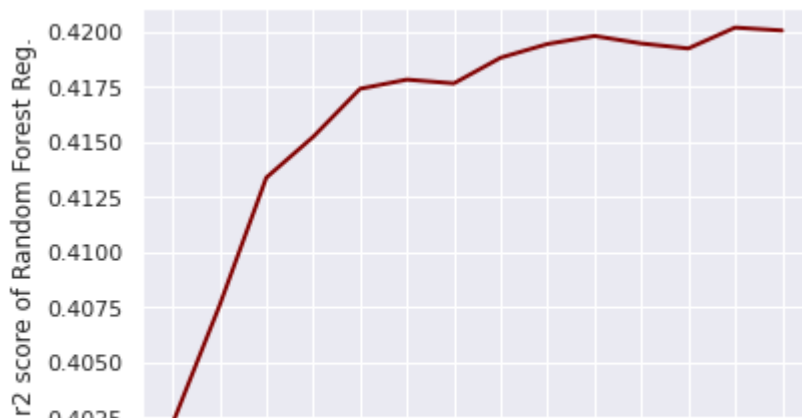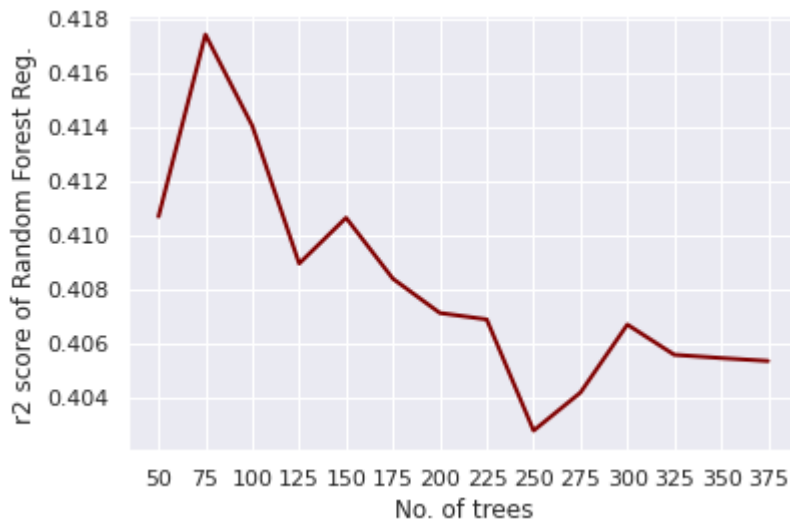






```
1 ## Training the polynomial regression on the training model
2 def poly(xtrain, ytrain , xtest ,ytest):
3    poly_reg = PolynomialFeatures(degree=2)
```

```
 4    x_poly = poly_reg.fit_transform(xtrain)
 5    poly_regressor = LinearRegression()
 6    poly_regressor.fit(x_poly,ytrain)
 7    y_pred = poly_regressor.predict(poly_reg.fit_transform(xtest))
 8    r2_poly = r2_score(ytest,y_pred)
 9    return r2_poly
10 print(poly(x_train,y_train ,x_test,y_test))
11 print(poly(x2_train,y2_train ,x2_test,y2_test))
12 print(poly(x3_train,y3_train ,x3_test,y3_test))
```

```
    0.21695986796292321
    0.24684587182313034
    0.15610619368864342
```

```
 1 # Finding out the optimal number of trees for Random Forest Regression
 2 from sklearn.ensemble import RandomForestRegressor
 3 forestRange=range(50,300,25)
 4 def forests(xtrain,ytrain,xtest,ytest):
 5    scores_list=[]
 6    for i in forestRange:
 7        regressor_Forest = RandomForestRegressor(n_estimators=i,random_state=0)
 8        regressor_Forest.fit(xtrain,ytrain)
 9        y_pred = regressor_Forest.predict(xtest)
10        scores_list.append(r2_score(ytest,y_pred))
11    plt.plot(forestRange,scores_list,linewidth=2,color='maroon')
12    plt.xticks(forestRange)
13    plt.xlabel('No. of trees')
14    plt.ylabel('r2 score of Random Forest Reg.')
15    plt.show()
16 forests(x_train,y_train ,x_test,y_test)
17 forests(x2_train,y2_train ,x2_test,y2_test)
18 forests(x3_train,y3_train ,x3_test,y3_test)
```

```
 1 # Training the Random Forest regression on the training model
 2 def forest(xtrain, ytrain , xtest ,ytest,est):
 3   regressor_Forest = RandomForestRegressor(n_estimators=est,random_state=0)
 4   regressor_Forest.fit(xtrain,ytrain)
 5   y_pred = regressor_Forest.predict(xtest)
 6   r2_forest = r2_score(ytest,y_pred)
 7   return r2_forest
 8 print(forest(x_train,y_train ,x_test,y_test,75))
 9 print(forest(x2_train,y2_train ,x2_test,y2_test,350))
10 print(forest(x3_train,y3_train ,x3_test,y3_test,275))
```

```
    0.4173995514782516
    0.42019352237037844
    0.1209474383767134
```

```
 1   ## Applying XGBoost Regression model on the training set
 2 from xgboost import XGBRegressor
 3 def xgb(xtrain, ytrain , xtest ,ytest):
 4     regressor_xgb = XGBRegressor(objective ='reg:squarederror')
 5     regressor_xgb.fit(xtrain,ytrain)
 6     ## Predicting test results
 7     y_pred = regressor_xgb.predict(xtest)
 8     ## Calculating r2 score
 9     r2_xgb = r2_score(ytest,y_pred)
```

```
10     return (r2_xgb)
11 print(xgb(x_train,y_train ,x_test,y_test))
12 print(xgb(x2_train,y2_train ,x2_test,y2_test))
13 print(xgb(x3_train,y3_train ,x3_test,y3_test))

   0.3766346103180298
   0.42063805335977456
   0.19319476303447647
```
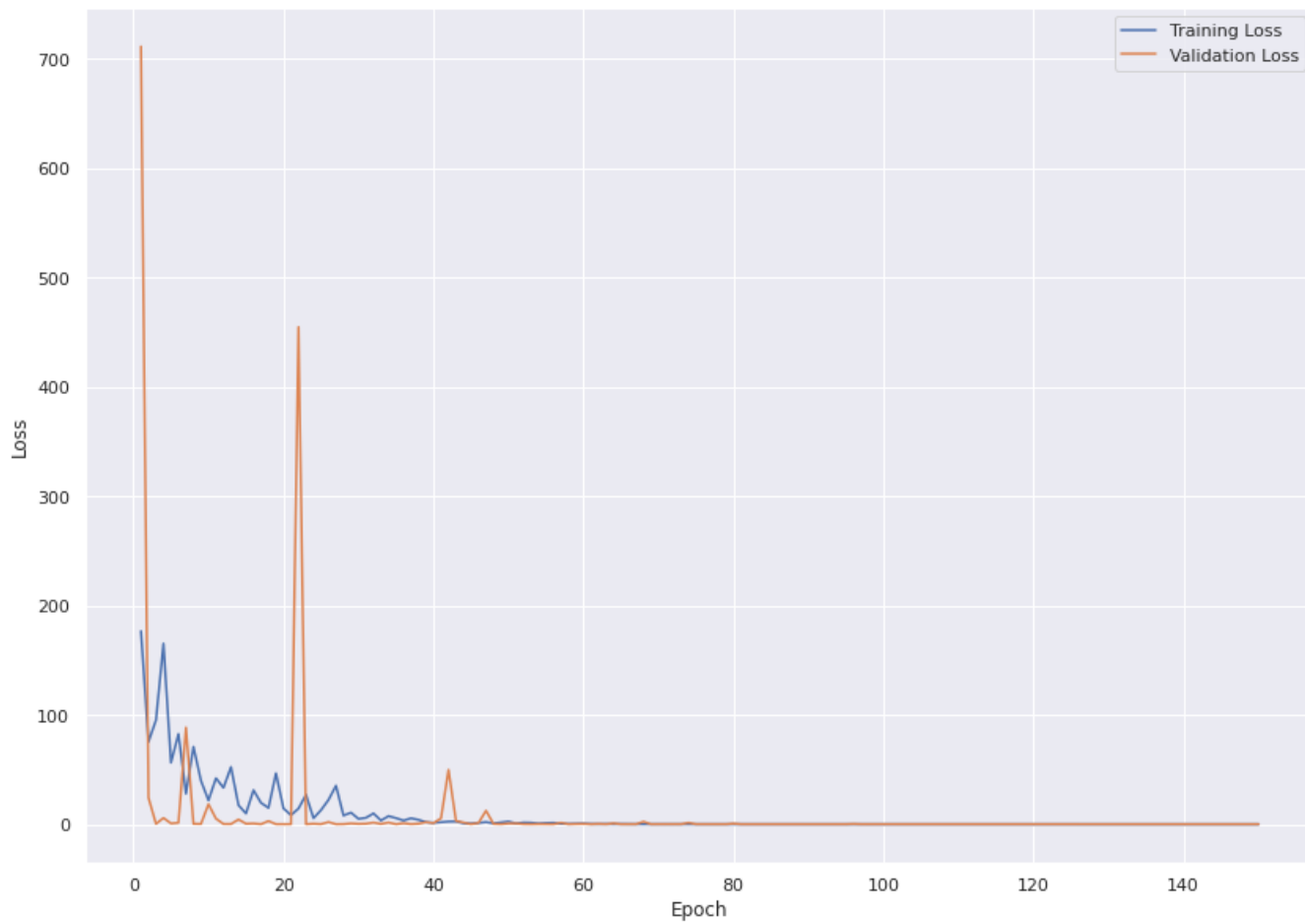
```
 1 import numpy as np
 2 from keras.models import Sequential
 3 from keras.layers import Dense, Dropout, Activation
 4 from keras.utils import np_utils
 5 from sklearn import metrics
 6 import tensorflow as tf
 7
 8 inputs = tf.keras.Input(shape=(7,))
 9 x = tf.keras.layers.Dense(128, activation='relu')(inputs)
10 x = tf.keras.layers.Dense(128, activation='relu')(x)
11 outputs = tf.keras.layers.Dense(1)(x)
12
13 model = tf.keras.Model(inputs=inputs, outputs=outputs)
14
15
16 optimizer = tf.keras.optimizers.RMSprop(0.001)
17
18 model.compile(
19     optimizer=optimizer,
20     loss='mse'
21 )
22
23
24 batch_size = 64
25 epochs = 150
26
27 history = model.fit(
28     x2_train,
29     y2_train,
30     validation_split=0.2,
31     batch_size=batch_size,
32     epochs=epochs,
33     verbose=0
34 )
```

```
 1 plt.figure(figsize=(14, 10))
 2
 3 epochs_range = range(1, epochs + 1)
 4 train_loss = history.history['loss']
 5 val_loss = history.history['val_loss']
 6
 7 plt.plot(epochs_range, train_loss, label="Training Loss")
```

```
 8 plt.plot(epochs_range, val_loss, label="Validation Loss")
 9
10 plt.xlabel("Epoch")
11 plt.ylabel("Loss")
12 plt.legend()
13
14 plt.show()
```



```
1 np.argmin(val_loss)
```

    93

```
1 !pip install tensorflow_addons
```

    Collecting tensorflow_addons
      Downloading tensorflow_addons-0.16.1-cp37-cp37m-manylinux_2_12_x86_64.manylinux2010_x8
         |████████████████████████████████| 1.1 MB 8.0 MB/s
    Requirement already satisfied: typeguard>=2.7 in /usr/local/lib/python3.7/dist-packages

```
Installing collected packages: tensorflow-addons
Successfully installed tensorflow-addons-0.16.1
```

◄ ▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐▐ ►

```
1 from tensorflow_addons.metrics import RSquare
2 y_pred = np.squeeze(model.predict(x_test))
3
4 result = RSquare()
5 result.update_state(y_test, y_pred)
6
7 print("R^2 Score:", result.result())
```

```
R^2 Score: tf.Tensor(0.038393497, shape=(), dtype=float32)
```

```
1 model.evaluate(x_test, y_test)
```

```
43/43 [==============================] - 0s 2ms/step - loss: 4.0662
4.066182613372803
```

Colab paid products  ·  Cancel contracts here

● ✕