

Project Name :- Video Games Sales prediction

I am going to build project of The Machine Learning model to predict future sales of Video Games in global market. We will learn different patterns and relationships between the input Variables or factors that affecting to the video games sales. We analyse data to get meaning information from them and Visualize them to recognise different patterns among them and showing meaning full information that will help us to predict sales data by creating effective machine learning model.

We use Data set which contains the list of video games with sales greater than 16500 games in global market.

Dataset Description

Rank - Ranking of overall sales

Name - The games name

Platform - Platform of the games release (i.e. PC,PS4, etc)

Year - Year of the game's release

Genre - Genre of the game

Publisher - Publisher of the game

NA_Sales - Sales in North America (in millions)

EU_Sales - Sales in Europe (in millions)

JP_Sales - Sales in Japan (in millions)

Other_Sales - Sales in the rest of the world (in millions)

Global_Sales - Total worldwide sales.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("darkgrid")

from sklearn.metrics import mean_absolute_error as mae, mean_squared_error as mse, r2_score

import warnings
warnings.filterwarnings('ignore')
```

Importing Data set

```
In [2]: df = pd.read_csv("vg_sales.csv")
df.head(5)
```

std	4791.853933	5.82898	0.816683	0.505351	0.309291	0.188588	1.555028
min	1.000000	1980.000000	0.000000	0.000000	0.000000	0.000000	0.010000
25%	4151.250000	2003.000000	0.000000	0.000000	0.000000	0.000000	0.060000
50%	8300.500000	2007.000000	0.080000	0.020000	0.000000	0.010000	0.170000
75%	12440.750000	2010.000000	0.240000	0.110000	0.040000	0.040000	0.470000
max	16600.000000	2020.000000	41.490000	29.020000	10.220000	10.570000	82.740000

Handling Missing Values

```
df.isna().sum()
```

Rank	0
Name	0
Platform	0
Year	271
Genre	0
Publisher	58
NA_Sales	0
EU_Sales	0
JP_Sales	0
Other_Sales	0
Global_Sales	0

Exploratory Data Analysis

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16398 entries, 0 to 16597
Data columns (total 11 columns):
# Column Non-Null Count Dtype
---
0 Rank 16598 non-null int64
1 Name 16598 non-null object
2 Platform 16598 non-null object
3 Year 16397 non-null float64
4 Genre 16598 non-null object
5 Publisher 16540 non-null object
6 NA_Sales 16598 non-null float64
7 EU_Sales 16598 non-null float64
8 JP_Sales 16598 non-null float64
9 Other_Sales 16598 non-null float64
10 Global_Sales 16598 non-null float64
dtypes: float64(6), int64(1), object(4)
memory usage: 1.1+ MB
```

As we see, there are 16598 vedio games are there with 11 variables. In the Year and Publisher columns, there will some missing Values

```
In [5]: df.describe()
```

	Rank	Year	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
count	16598.000000	16327.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000
mean	8791.853254	2006.406443	0.264667	0.166652	0.077782	0.048063	0.537441
std	4791.853933	5.828981	0.816683	0.505351	0.309291	0.188588	1.555028
min	1.000000	1980.000000	0.000000	0.000000	0.000000	0.000000	0.010000
25%	4151.250000	2003.000000	0.000000	0.000000	0.000000	0.000000	0.060000
50%	8300.500000	2007.000000	0.080000	0.020000	0.000000	0.010000	0.170000
75%	12449.750000	2010.000000	0.240000	0.110000	0.040000	0.040000	0.470000
max	16600.000000	2020.000000	41.490000	29.020000	10.220000	10.570000	82.740000

Handling Missing Values

```
In [5]: df.isna().sum()
```

```
Out[5]: Rank      0
Name      0
Platform  0
Year     271
Genre     0
Publisher 58
NA_Sales  0
EU_Sales  0
JP_Sales  0
Other_Sales 0
Global_Sales 0
dtype: int64
```

```
In [6]: df.dropna(inplace=True)
```

```
In [7]: df.isna().sum()
```

```
Out[7]: Rank      0
Name      0
Platform  0
Year     0
Genre     0
Publisher 0
NA_Sales  0
EU_Sales  0
JP_Sales  0
Other_Sales 0
Global_Sales 0
dtype: int64
```

Knowing about shape and size of dataset

```
In [10]: df.shape
```

```
Out[10]: (16291, 11)
```

```
In [11]: df.size
```

```
Out[11]: 179201
```

```
In [12]: df.columns
```

```
Out[12]: Index(['Rank', 'Name', 'Platform', 'Year', 'Genre', 'Publisher', 'NA_Sales',
      'EU_Sales', 'JP_Sales', 'Other_Sales', 'Global_Sales'],
      dtype='object')
```

Unique values present in each feature

```
In [13]: print("Total number Categories/Classes present in Each Variable")
for col in df:
    uni=df[col].nunique()
    print(f"{col} :- {uni}")
```

Total number Categories/Classes present in Each Variable
Rank :- 16291
Name :- 1325
Platform :- 31
Year :- 39
Genre :- 12
Publisher :- 576
NA_Sales :- 408
EU_Sales :- 305
JP_Sales :- 244
Other_Sales :- 157
Global_Sales :- 621

Analysis Categorical Features

```
In [14]: # Creating Categorical DataFrame
df_cat=df.select_dtypes("object")
df_cat.head()
```

```
Out[14]:
```

	Rank	Name	Platform	Genre	Publisher
0	1	Wii Sports	Wii	Sports	Nintendo
1	2	Super Mario Bros.	NES	Platform	Nintendo
2	3	Mario Kart Wii	Wii	Racing	Nintendo
3	4	Wii Sports Resort	Wii	Sports	Nintendo
4	5	Pokemon Red/Pokemon Blue	GB	Role-Playing	Nintendo

checking the top (most repetitive) values.

```
In [15]: for cat in df_cat.columns:
    print("-----")
    print(f"cat: {cat}")
    print(df[cat].value_counts().head())
```

cat: Name

Need for Speed: Most Wanted 12
FIFA 14 9
Bakatonville 9
LEGO Marvel Super Heroes 9
The LEGO Movie Videogame 8
Name: Name, dtype: int64

Platform

DS 2131
PS2 2127
PS3 1304
Wii 1290
X360 1234
Name: Platform, dtype: int64

Genre

Action 3251
Sports 2304
Misc 1886
Role-Playing 1470
Shooter 1282
Name: Genre, dtype: int64

Publisher

Electronic Arts 1339
Activision 966
Namco Bandai Games 928
Ubisoft 918
Konami Digital Entertainment 823
Name: Publisher, dtype: int64

The game which has highest global sales

```
In [16]: high = df["Global_Sales"].max()
top_game = df[df["Global_Sales"] == high]
#top_game[["Name", "Global_Sales", "Year"]]
top_game
```

	Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
0	1	Wii Sports	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74

The game which has Lowest global sales

```
In [17]: low = df["Global_Sales"].min()
low_game = df[df["Global_Sales"] == low]
low_game
```

	Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
15980	15983	Turok	PC	2008.0	Action	Touchstone	0.00	0.01	0.00	0.0	0.01
15981	15984	Coven and Labyrinth of Refrain	PSV	2016.0	Action	Nippon Ichi Software	0.00	0.00	0.01	0.0	0.01
15982	15985	Super Battle For Money Sentouchi: Kyuukyoku n...	3DS	2016.0	Action	Namco Bandai Games	0.00	0.00	0.01	0.0	0.01
15983	15986	Dragon Zakura DS	DS	2007.0	Misc	Electronic Arts	0.00	0.00	0.01	0.0	0.01
15984	15987	Chameleon: To Die For	DS	2006.0	Puzzle	SOS Games	0.01	0.00	0.00	0.0	0.01
...
16593	16596	Woody Woodpecker in Black & Blue Alien Case 5	GBA	2002.0	Platform	Kemco	0.01	0.00	0.00	0.0	0.01
16594	16597	Men in Black II: Alien Escape	GC	2003.0	Shooter	Infogrames	0.01	0.00	0.00	0.0	0.01
16595	16598	SCORE International Bga 1000: The Official Game	PS2	2008.0	Racing	Activision	0.00	0.00	0.00	0.0	0.01
16596	16599	Know How 2	DS	2010.0	Puzzle	7G/AMES	0.00	0.01	0.00	0.0	0.01
16597	16600	Spirits & Spells	GBA	2003.0	Platform	Wanadoo	0.01	0.00	0.00	0.0	0.01

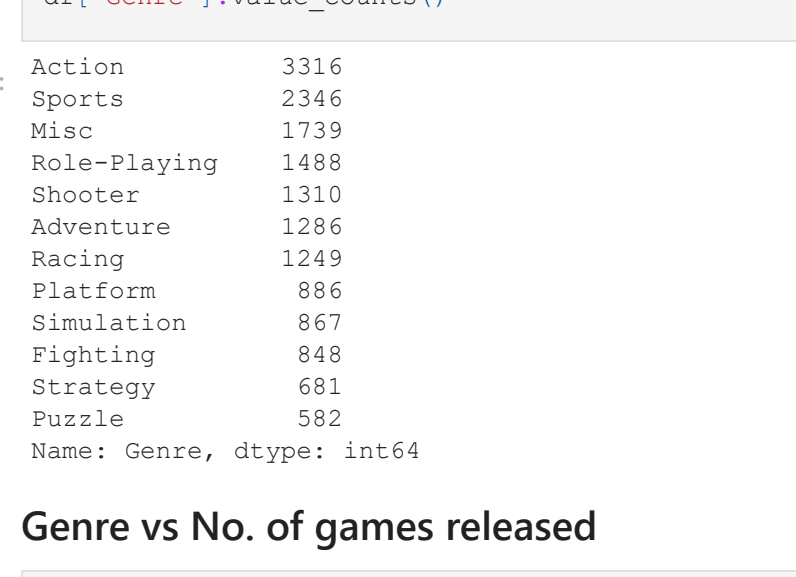
600 rows x 11 columns

Data Visualization

Global sales vs No. of games sold

```
In [3]: plt.hist(df["Global_Sales"], bins=20)
plt.xlabel("Sales in Million", fontsize=15)
plt.xticks(np.arange(0,80,5))
plt.ylabel("No. of games sold")
plt.show()
```

in 0 to 5 million range most games are sold



```
In [4]: # Total number of Platform present
df["Platform"].unique()
```

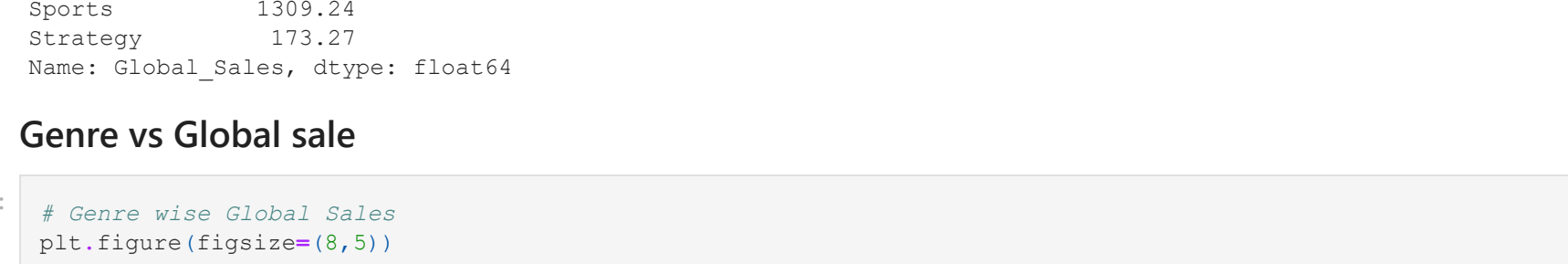
```
Out[4]: array(['WiiU', 'NES', 'GB', 'DS', 'X360', 'PS3', 'PS2', 'SNES', 'GBA',
      'PS4', 'PSV', 'Wii', 'PC', 'Xbox One', 'Xbox 360', 'Xbox', 'GC',
      'WiiU', 'GEN', 'DC', 'PSV', 'SAT', 'SCD', 'MS', 'NG', 'TG16',
      '3DO', 'GG', 'PCFX'], dtype=object)
```

```
In [5]: # Top Platforms who released games
df["Platform"].value_counts()
```

```
Out[5]: DS      2163
PS2     2161
PS3     1329
Wii      1325
X360     1265
PS4     1213
PSV     1196
PC       960
XB       824
GBA      822
GC       556
3DS     509
PSV     413
PS4     336
WiiU     319
SNES     239
Xbox One 213
SAT      173
WiiU     143
2600     133
NES       98
GB        98
DC        52
GEN       27
NG        12
SCD        6
WS         6
3DO        3
TG16       2
OG         1
PCFX       1
Name: Platform, dtype: int64
```

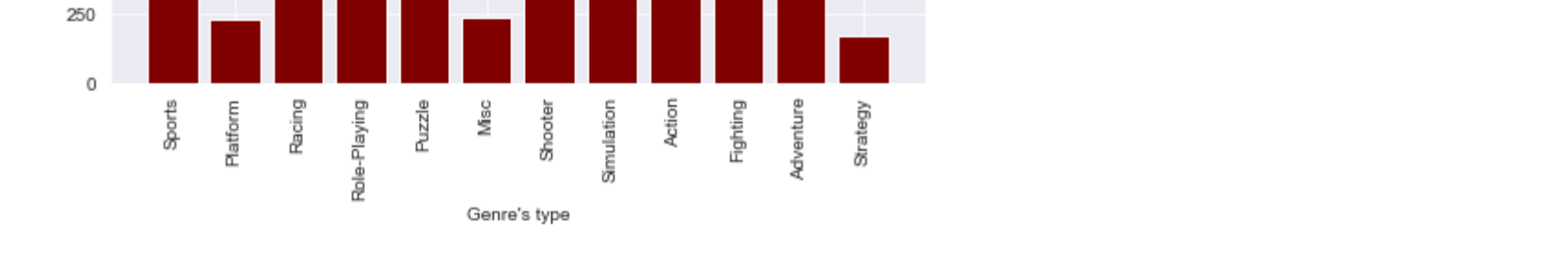
Platform vs No. of games released

```
In [6]: # Platform which released maximum Games
plt.figure(figsize=(15,5))
sns.countplot(x="Platform", data=df, order = df["Platform"].value_counts().index) # unique does not give high to low
plt.xticks(rotation=90, fontsize=12)
plt.grid(True)
```



```
In [7]: # A Platform which contributes more in Global Sales
dfyear = df.groupby(["Year"])["Global_Sales"].sum().reset_index().sort_values("Global_Sales", ascending=False)
dfyear = dfyear.reset_index() #to reset index of a Data Frame. reset_index() method sets a list of integer ranging from 0 to length of data as index.
plt.figure(figsize=(15,5))
sns.barplot(x="Platform", y="Global_Sales", data=dfyear)
plt.xticks(rotation=90, fontsize=12)
plt.xlabel("Platform", fontsize=12)
plt.ylabel("Global Sale", fontsize=12)
plt.grid(True)
```

PS2 platform contributes more in Global Sales



```
In [8]: df["Genre"].unique()
array(['Sports', 'Platform', 'Action', 'Role-Playing', 'Puzzle', 'Misc',
      'Shooter', 'Simulation', 'Racing', 'Role-Playing', 'Puzzle', 'Misc',
      'Strategy'], dtype=object)
```

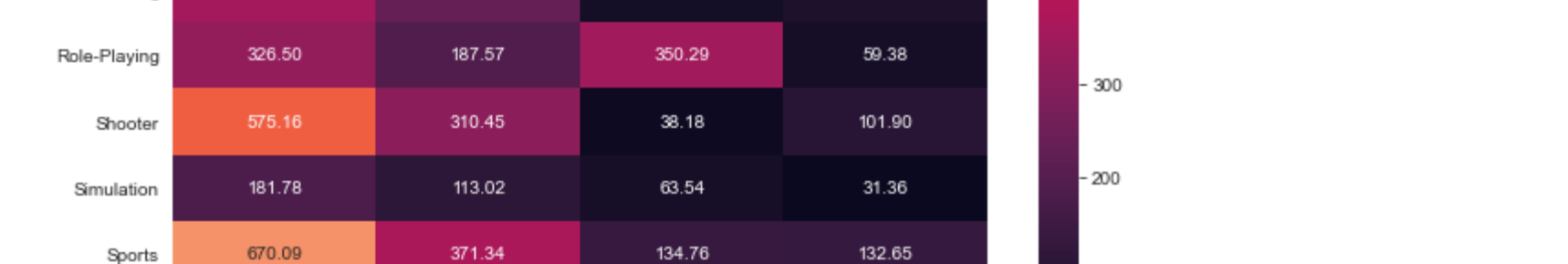
top 10 Publihsers

```
In [9]: # top10 Publishers
top10_pub=df["Publisher"].value_counts().head(10)
top10_pub
```

```
Out[9]: Electronic Arts      1351
Activision      975
Namco Bandai Games      928
Ubisoft      921
Konami Digital Entertainment      832
THQ      822
Nintendo      703
Sony Computer Entertainment      683
Sega      639
Take-Two Interactive      413
Name: Publisher, dtype: int64
```

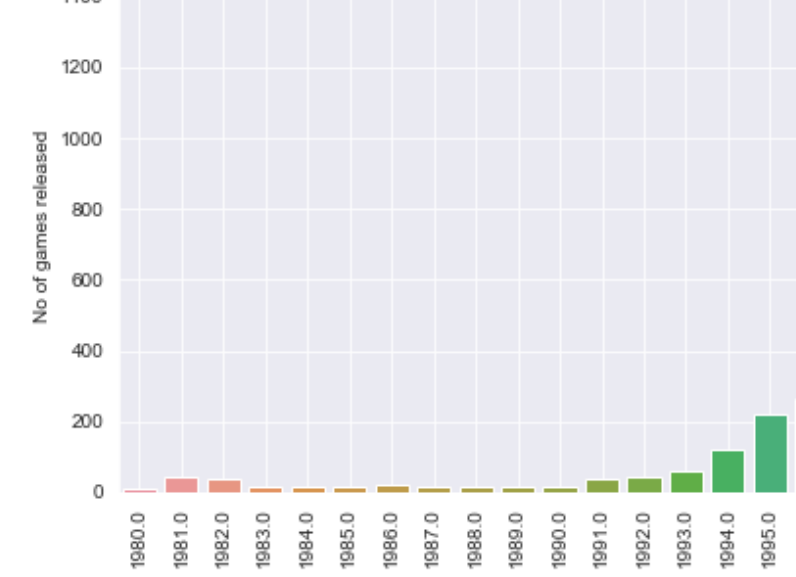
```
In [10]: df3 = df.groupby(["Publisher"]).count().iloc[:10]
df3 = pd.DataFrame(df3.sort_values(ascending=False))[:10]
publishers = df3.index
df3.columns = ["Releases"]
```

```
plt.figure(figsize=(12,8))
sns.barplot(y=publishers, x = "Releases", data=df3)
plt.xlabel("Number of Releases")
plt.ylabel("Publishers")
plt.title("Top 10 Total Publisher Games Released", fontsize=20)
plt.show()
```



```
In [11]: # Years analysis for releasing games
plt.hist(df["Year"], bins=40)
plt.xticks(rotation=90)
plt.xlabel("years")
plt.ylabel("No of games released")
plt.show()
```

2008-10 most games released



```
In [12]: # genres to analysis to how many games have in particular genre
df["Genre"].value_counts()
```

```
Out[12]: Action      3316
Sports      2346
Misc      1739
Role-Playing      1488
Shooter      1310
Adventure      1286
Racing      1249
Platform      886
Simulation      867
Fighting      848
Strategy      681
Puzzle      582
Name: Genre, dtype: int64
```

Genre vs No. of games released

```
In [14]: # Which has released most games
plt.figure(figsize=(9,5))
sns.countplot(x="Genre", data=df, order = df["Genre"].value_counts().index)
plt.xticks(rotation=90, fontsize=12)
plt.xlabel("Genre", fontsize=12)
plt.grid(True)
```

most of the people love action and sports game. in action 3251 and in sports 2304 games release.



```
In [29]: x_val=df["Genre"].unique()
x_val.shape
```

```
Out[29]: (12,)
```

```
In [30]: y_val=df.groupby("Genre")["Global_Sales"].sum()
y_val
```

```
Out[30]: Genre
Action      1722.84
Adventure      234.59
Fighting      444.05
Misc      789.87
Platform      829.13
Puzzle      242.21
Racing      726.76
Role-Playing      923.83
Shooter      1025.20
Simulation      389.98
Sports      1309.24
Strategy      173.27
Name: Global Sales, dtype: float64
```

Genre vs Global Sale

```
In [31]: # Genre wise Global Sales
plt.figure(figsize=(8,5))
plt.bar(x_val,y_val,color="maroon")
plt.xticks(rotation=90)
plt.xlabel("Genre's type")
plt.ylabel("Global Sale")
plt.show()
```


genrewise regions Sales Comparison

```
In [32]: # genrewise regions Sales Comparison
df1 = df[["Genre", "NA_Sales", "EU_Sales", "JP_Sales", "Other_Sales"]]
df1 = df1.groupby("Genre").sum()
df1
```

	NA_Sales	EU_Sales	JP_Sales	Other_Sales
Genre				
Action	861.77	516.48	158.65	184.92
Adventure	101.93	63.74	51.59	16.70
Fighting	209.74	100.00	87.17	36.19
Misc	365.92	211.77	106.67	73.92
Platform	445.99	200.65	190.65	51.51
Puzzle	122.01	50.52	16.68	12.47
Racing	368.93	236.31	96.61	76.68
Role-Playing	326.50	187.57	350.29	59.38
Shooter	575.16	310.45	38.18	101.80
Simulation	181.78	113.02	63.54	31.36
Sports	870.98	373.94	134.76	132.65
Strategy	87.83	44.84	49.10	11.23

```
In [33]: plt.figure(figsize=(10, 8))
sns.heatmap(df1, annot=True, fmt=".2f")
```

```
Out[33]: <axes.SubplotTitled: <Figure>>
```


Years vs No of games released

```
In [34]: # Numbers of games released per year / Year which most games released
dfyear = df.groupby("Year")["Name"].count()
plt.figure(figsize=(15,5))
sns.countplot(x="Year", data=df, order=dfyear.index)
plt.xticks(rotation=90)
```

North America almost Cover 50% of sales generated


```
In [35]: df["Global_Sales"].sum()
Out[35]: 8811.97
```

Year vs Global sale

```
In [36]: # The year which has Highest Global Sale
dfyear = df.groupby("Year")["Global_Sales"].sum() # grouping the data according to the categories and apply +
dfyear = dfyear.reset_index() #to reset index of a Data Frame. reset_index() method sets a list of integer rang
```

```
plt.figure(figsize=(15,5))
sns.barplot(x="Year", y="Global_Sales", data=dfyear)
plt.xticks(rotation=90)
plt.xlabel("Global Sale")
plt.grid(True)
```


Region wise Global Sales analysis

```
In [37]: # Region wise Global Sales analysis
df2 = df[["NA_Sales", "EU_Sales", "JP_Sales", "Other_Sales"]]
df2 = df2.sum().reset_index()
```

```
Out[37]: index      0
0 NA_Sales  4327.65
1 EU_Sales  2406.69
2 JP_Sales  1284.27
3 Other_Sales  788.91
```

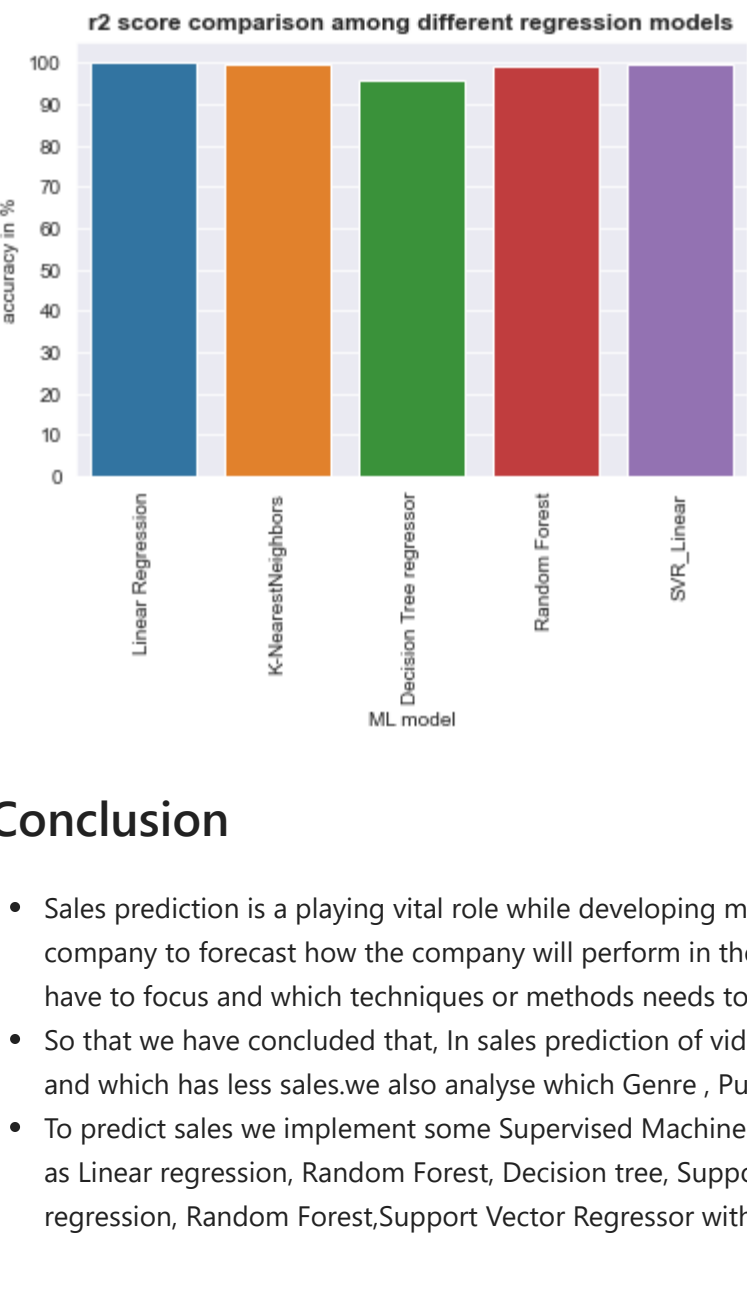
```
In [38]: # By Bar Graph
plt.bar(df2["index"],df2[0])
```

```
Out[38]: <BarContainer object of 4 artists>
```



```
In [39]: # By pie chart
plt.figure(figsize=(10, 8))
plt.pie(df2[0], labels=df2["index"], autopct='%1.2f%%')
plt.show()
```

North America almost Cover 50% of sales generated



Conclusion

- Sales prediction is a playing vital role while developing marketing techniques and and strategic planning by any company. It allows a company to forecast how the company will perform in the future. Sales prediction is als important to make decision on which area we have to focus and which techniques or methods needs to implement to increase sales.
- So that we have concluded that, In sales prediction of video games we observed which game has more sales in the market globally and which has less sales.we also analyse which Genre , Publishers and platforms contributes more and less in global sales.
- To predict sales we implement some Supervised Machine Learning Techniques and Create better model by Different Algorithms such as Linear regression, Random Forest, Decision tree, Support vector regression,kneighbors regressor. Among all these algorithms Linear regression, Random Forest,Support Vector Regressor with Linear gave us the best accurate result with minimum error rate.

Thank You !!!