

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("vgsales.csv")
df.head(5)
```

Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	
0	1	Wii Sports	Wii	2006	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74
1	2	Super Mario Bros.	NES	1985	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24
2	3	Mario Kart Wii	Wii	2008	Racing	Nintendo	15.85	12.88	3.79	3.31	35.82
3	4	Wii Sports Resort	Wii	2009	Sports	Nintendo	15.75	11.01	3.28	2.96	33.00
4	5	Pokemon Red/Pokemon Blue	GB	1996	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37

```
Out[2]: df.columns
```

```
In [3]: index(['Rank', 'Name', 'Platform', 'Year', 'Genre', 'Publisher', 'NA_Sales',
'EU_Sales', 'JP_Sales', 'Other_Sales', 'Global_Sales'],
dtype='object')
```

Exploratory Data Analysis

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

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16598 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                  16327 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 video games are there with 11 variables. In the Year and Publisher columns, there will be some missing values

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

Out[5]:
```

	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
std	830.605254	2006.406443	0.264667	0.146652	0.077782	0.048063	0.537441
mean	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%	8306.500000	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

```
In [6]: df.ndim
print("Dimensions of Data : ",df.shape)
print("Size of Data : ",df.size)
```

Handling Missing Values

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

Out[7]:
Rank          0
Name          0
Platform      0
Year         271
Genre        1470
Publisher     58
NA_Sales      0
EU_Sales      0
JP_Sales      0
Other_Sales   0
Global_Sales   0
dtype: int64

In [8]: df.dropna(inplace=True)

In [9]: df.isna().sum()

Out[9]:
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:
    uval=coll.uniques()
    print("({} : {}".format(col, len(uval)))

total number Categories/Classes present in Each Variable
Rank : 16291
Name : 11325
Platform : 31
Year : 19
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()
```

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

checking the top (most repetitive) values.

```
In [15]: for cat in df_cat.columns:
    print("-----")
    print(cat, df_cat.value_counts().head())
    print("-----")

Name
-----
Need for Speed: Most Wanted    12
FFFA 14                        9
Ratatouille                    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        1686
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	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	0.00	0.01
15981	15984	Coven and Labyrinth of Refrain	PSV	2016.0	Action	Nippon Ichi Software	0.00	0.00	0.01	0.00	0.01
15982	15985	Super Battle For Money	3DS	2016.0	Action	Namco Bandai Games	0.00	0.00	0.01	0.00	0.01
15983	15986	Dragon Zakura DS	DS	2007.0	Misc	Electronic Arts	0.00	0.00	0.01	0.00	0.01
15984	15987	Chameleon To Dye For	DS	2006.0	Puzzle	SOS Games	0.01	0.00	0.00	0.00	0.01
...	
16593	16596	Woody Woodpecker in Crazy Castle 5	GBA	2003.0	Platform	Kemco	0.01	0.00	0.00	0.00	0.01
16594	16599	Men in Black II Alien Escape	GC	2002.0	Shooter	Infogrames	0.01	0.00	0.00	0.00	0.01
16595	16598	SCORE International Baja 1000: The Official Game	PS2	2008.0	Racing	Activision	0.00	0.00	0.00	0.00	0.01
16596	16599	Knock How 2	DS	2010.0	Puzzle	7G/AMES	0.00	0.01	0.00	0.00	0.01
16597	16600	Spirits & Spells	GBA	2003.0	Platform	Vanadoo	0.01	0.00	0.00	0.00	0.01

600 rows × 11 columns

Data Visualization

Global sales vs No. of games sold

```
In [18]: plt.hist(df["Global_Sales"],bins=20)
plt.xlabel("Sales in Million",fontsize=15)
plt.ylabel("No. of games sold")
plt.show()
# In 0 to 5 million range most games are sold
```

```
In [19]: # Total number of Platform present
df["Platform"].unique()

Out[19]: array(['WiiU', 'NES', 'GB', 'DS', 'X360', 'PSP', 'P2P', 'SNES', 'GBA',
'3DS', 'PS4', 'N64', 'PS3', 'XBB', 'PC', '2600', 'PSP', 'XOne', 'GC',
'WiiU', 'GEN', 'DC', 'PSV', 'SAT', 'SCD', 'MS', 'NG', 'TG16',
'3DO', 'GG', 'PCFX'], dtype=object)
```

```
In [20]: # Top Platforms who released games
df["Platform"].value_counts()

Out[20]:
DS      2131
PS2     2127
PS3     1304
Wii     1290
X360    1234
PSP     1197
PS      1189
PC       938
XB       803
GBA      786
GC       542
3DS     499
PSV     410
PSP4     336
N64      316
SNES     239
XOne     213
SAT      173
WiiU     145
2600     116
MS       98
GB        97
DC        52
GEN       27
NG        12
SCD        6
WOS        6
TG16       2
GG         1
GCX        1
PCFX       1
Name: Platform, dtype: int64
```

Platform vs No. of games released

```
In [21]: # Platform which released maximum Games
df.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)
plt.show()
```

Platform vs Global sale

```
In [22]: # A platform which contributes more in Global Sales
dfyear = df.groupby(["Year"])[["Global_Sales"]].sum().reset_index().sort_values("Global_Sales",ascending=False)
#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", data=dfyear, data=dfyear)
plt.xticks(rotation=90,fontsize=12)
plt.xlabel("Platform",fontsize=12)
plt.ylabel("Global Sale",fontsize=12)
plt.grid(True)
plt.show()
# PS2 platform contributes more in Global Sales
```

```
In [23]: df["Genre"].unique()

Out[23]: array(['Sports', 'Platform', 'Racing', 'Role-Playing', 'Puzzle', 'Misc',
'Shooter', 'Simulation', 'Action', 'Fighting', 'Adventure',
'Strategy'], dtype=object)
```

top 10 Pbliherss

```
In [24]: # top 10 Publishers
top10_pub=df["Publisher"].value_counts().head(10)
top10_pub
```

	Electronic Arts	Activision	Namco Bandai Games	Ubisoft	Konami Digital Entertainment	THQ	Nintendo	Sony Computer Entertainment	Sega	Take-Two Interactive
	1339	966	928	918	823	712	696	682	632	412

```
In [25]: df3 = df.groupby(["Publisher"]).count().iloc[:,0]
df3 = pd.DataFrame(df3.sort_values(ascending=False))
publishers = df3.index
df3.columns = ["publishers"]

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 [26]: # Years analysis for releasing games
plt.hist(df["Year"],bins=40)
plt.xticks(rotation=90)
plt.show()
# 2008-10 most games released
```

```
In [27]: # Genres to analysis to how many games have in particular genre
df["Genre"].value_counts()

Out[27]:
Action      3251
Sports      2304
Misc        1686
Role-Playing 1470
Shooter     1282
Adventure    1274
Racing       1225
Platform     875
Simulation    848
Fighting     836
Strategy     670
Puzzle       570
Name: Genre, dtype: int64
```

Genre vs No. of games released

```
In [28]: # Which has released most games
plt.figure(figsize=(12,7))
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)
plt.show()
# most of the people love action and sports game. in action 3316 and in sports 2346 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
```

Genre	Global_Sales
Action	1722.84
Adventure	234.59
Fighting	444.05
Misc	789.87
Platform	922.13
Puzzle	242.21
Racing	726.76
Role-Playing	922.83
Shooter	1026.20
Simulation	389.98
Sports	1309.24
Strategy	173.27

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()

Out[32]:
```

Genre	NA_Sales	EU_Sales	JP_Sales	Other_Sales
Action	861.77	516.48	158.65	184.92
Adventure	101.93	63.74	51.59	16.70
Fighting	220.74	100.00	87.15	36.19
Misc	396.92	211.77	108.67	73.92
Platform	443.99	200.65	130.65	51.51
Puzzle	122.01	50.52	56.68	12.47
Racing	368.93	226.31	96.61	70.68
Role-Playing	328.50	187.57	203.29	59.38
Shooter	672.16	310.45	38.18	101.90
Simulation	181.78	113.02	63.54	31.36
Sports	670.09	371.34	134.76	132.65
Strategy	67.83	44.84	40.10	11.23

Years vs No of games released

```
In [34]: # Numbers of games released per year and Year which most games released
yrgame = df.groupby(["Year"])[["Name"]].count()
plt.figure(figsize=(15,5))
sns.countplot(x="Year", data=yrgame, order=yrgame.index)
plt.xticks(rotation=90)
plt.xlabel("No of games released")
plt.grid(True)
plt.show()
```

```
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 ranging from 0 to length of data as index.

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)
plt.show()
```

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

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
```


