

Title: Integrated Assessment (HDip for DA)

- Student Name: Eoin Cantwell
- Student Number: sba21230
- Submission Date: 02/06/2021

Assignment Introduction

The data set was found on Kaggle and it can be accessed through the following link:
<https://www.kaggle.com/rgregout/video-gamesales>

- The data set is on the sales of worldwide video games that each sold 100k copies or more. Overall, there are 11 feature columns that have both quantitative/qualitative data, there are a total of 16,598 observations and this data set has not been used by me before. All of which confirms the assignments specifications and requirements at the outset are met.
- This integrated assignment is based on an video game sales data set. The assignment covers two modules 1) Statistics 2) Data Preparation and there are a set of questions for each module. Overlap may occur but each question heading will specify what module/question is included for discussion.
- I chose this data set as when I was younger I was into playing computer games and so I found it very interesting to dig into what insights could be extracted from the data set.

Q1. Statistics / Q2 Data Preparation: Characterization of the Data Set

```
In [2]: #Importing required libraries
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
import warnings
warnings.filterwarnings("ignore")

In [5]: #Reading in the data set from where it is stored on my PC
#The data is housed in a .csv file (comma separated values)
df = pd.read_csv("data/vgsales.csv")

In [6]: #Shows us the first 5 rows, this can be increased by inserting a number inside brackets
df.head()
```

	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	84.66
1	2	Super Mario Bros.	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	39.24
2	3	Mario Kart Wii	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	33.83
3	4	Wii Sports Resort	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	32.99
4	5	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37

```
In [5]: #Renaming columns into lower and snake case to simplify code writing and aid accuracy
new_names = {'Rank':'rank', 'Name':'name', 'Platform':'platform',
              'Year':'year', 'Genre':'genre', 'Publisher':'publisher',
              'NA_Sales':'north_american_sales', 'EU_Sales':'european_sales',
              'JP_Sales':'japanese_sales', 'Other_Sales':'other_sales',
              'Global_Sales':'global_sales'}

In [6]: #Passing through the renaming of the columns so it is in place
df.rename(columns=new_names, inplace=True)

In [7]: #Increasing the frame width to the max width, this spaces columns nicely and data can
pd.set_option('display.max_colwidth', None)
df.head()
```

	rank	name	platform	year	genre	publisher	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
0	1	Wii Sports	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	84.66
1	2	Super Mario Bros.	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	39.24
2	3	Mario Kart Wii	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	33.83
3	4	Wii Sports Resort	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	32.99
4	5	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37

The above output confirms the columns were successfully renamed

```
In [8]: #Shows the number of observations and features the data set has
df.shape
```

```
Out[8]: (16598, 11)
```

Data Dictionary

- A data dictionary is a practical depiction of the metadata connected with a data object, such as a database, table, or column (OCTOPAL, 2019). The data dictionary helps give an understanding as to what each feature means and is very useful at the start of any data analysis project, for both the user and any person using/reading the data at a later date.

Feature	Description
rank	The order of each video game based on total sales
name	The name of the video game
platform	The machine the game can be played on
year	The year of release
genre	The type of the video game
publisher	The company that created the video game
north_american_sales	The total of North American video game sales
european_sales	The total of European video game sales
japanese_sales	The total of Japanese video game sales
other_sales	The total of other regions/countries video game sales
global_sales	The overall global total of video games sales

```
In [9]: 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 north_american_sales 16598 non-null    float64
7 european_sales       16598 non-null    float64
8 japanese_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: 1.4+ MB

The above code shows the number of observations (rows), the number of columns (variables/features), what data type each feature is and finally the memory usage the data set demands from your pc.
```

```
In [10]: #Creating new df to enable numeric values be used in the statistical function below
stats_df1 = df[['north_american_sales', 'european_sales', 'japanese_sales', 'other_sales', 'global_sales']]

stats_df1.describe()
```

	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
count	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000
mean	0.264667	0.146652	0.077782	0.048063	0.537441
std	0.816683	0.505351	0.309291	0.188588	1.555028
min	0.000000	0.000000	0.000000	0.000000	0.010000
25%	0.000000	0.000000	0.000000	0.000000	0.060000
50%	0.080000	0.020000	0.000000	0.010000	0.170000
75%	0.240000	0.110000	0.040000	0.040000	0.470000
max	41.490000	29.020000	10.220000	10.570000	82.740000

The above table shows key statistical information regarding our quantitative data columns.

```
In [156]: #Show data types that are objects
df.describe(include='object')
```

	name	platform	genre	publisher
count	16598	16598	16598	16540
unique	11493	31	12	578
top	Need for Speed: Most Wanted	DS	Action	Electronic Arts
freq	12	2163	3316	1351

The above code gives an understanding of the categorical data in the dataset and is useful for future analysis. Insights gained, 31 different platforms/devices the video games are played on. There are 12 different genre types. Need for Speed: Most Wanted is the game that appears most.

Mode

- The following mode functions confirm what we is seen in the above table

```
In [10]: trend = df["platform"].mode()
print(trend)

0 DS
dtype: object

In [11]: trend = df["year"].mode()
print(trend)

0 2009.0
dtype: float64

In [12]: trend = df["genre"].mode()
print(trend)

0 Action
dtype: object

In [13]: trend = df["publisher"].mode()
print(trend)

0 Electronic Arts
dtype: object
```

The above functions show the mode of each categorical data type, the results are that Nintendo DS is the most frequent platform for the video games. The year where most video games are sold is 2009. The favourite genre is action and the most frequent publisher of video games is the company Electronic Arts.

```
In [158]: #This code reveals if there is any correlation between numerical variables
stats_df2 = df[['north_american_sales', 'european_sales', 'japanese_sales', 'other_sales', 'global_sales']]
stats_df2.corr()
```

	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
north_american_sales	1.000000	0.767727	0.449787	0.634737	0.941047
european_sales	0.767727	1.000000	0.435584	0.726385	0.902836
japanese_sales	0.449787	0.435584	1.000000	0.290186	0.618186
other_sales	0.634737	0.726385	0.290186	1.000000	0.748331
global_sales	0.941047	0.902836	0.618186	0.748331	1.000000

```
In [159]: #Shows how spread out the data is
skewness = df.skew()
print(skewness)
```

	rank	year	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
rank	0	0.000066	-1.002560	18.799627	18.875533	3.58	6.81
year	0	0	11.206458	24.233923	17.400645
north_american_sales	0	0.01	0.00	0.00	0.00	0.00	0.00
european_sales	0	0.01	0.00	0.00	0.00	0.00	0.00
japanese_sales	0	0.00	0.01	0.00	0.00	0.00	0.00
other_sales	0	0.00	0.01	0.00	0.00	0.00	0.00
global_sales	0	0.01	0.00	0.00	0.00	0.00	0.00
dtype:	float64

```
In [160]: df['year'].min()

Out[160]: 1980.0

In [161]: df['year'].max()

Out[161]: 2020.0
```

The min/max function reveals the dataset covers 40 years of video game sales which supports why there are so many observations.

```
In [162]: #Z Score
from scipy import stats
z = df[['north_american_sales', 'european_sales', 'japanese_sales', 'other_sales', 'global_sales']]
print(z)
```

	rank	year	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
0	41.49	29.02	3.77	8.46
1	29.08	3.58	6.81	0.77
2	15.85	12.88	3.79	3.31
3	15.75	11.01	3.28	2.96
4	11.27	8.89	10.22	1.00
...
16593	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16594	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16595	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16596	0.00	0.01	0.00	0.00	0.00	0.00	0.00
16597	0.01	0.00	0.00	0.00	0.00	0.00	0.00
...
global_sales	82.74	40.24	35.82	33.00	31.37
0	40.24	35.82	33.00	31.37
1	35.82	33.00	31.37
2	33.00	31.37
3	31.37
4
...
16593	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16594	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16595	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16596	0.01	0.00	0.00	0.00	0.00	0.00	0.00
16597	0.01	0.00	0.00	0.00	0.00	0.00	0.00

[16598 rows x 5 columns]

A Z-score is a numerical measurement that calls out a value's relationship to the mean of a group of values. Z-score is calculated in terms of standard deviations from the mean.

A Z-score equaling to 0 indicates that the data point's score is identical to the mean score. A Z-score of 1 means a value that is one standard deviation from the mean.

Z-scores may be positive or negative, with a positive value indicating the score is above the mean and a negative score indicating it is below the mean (Hayes, 2021).

Q1. Statistics / Q2 Data Preparation: Exploratory Data Analysis

```
In [18]: #Investigate if there is any null values present
df.isnull().sum()
```

	rank	name	platform	year	genre	publisher	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
rank	0	0	0	0	0	0	0	0	0	0	0
name	0	271	0	0	0	0	0	0	0	0	0
platform	0	0	0	0	0	0	0	0	0	0	0
year	0	0	0	0	0	0	0	0	0	0	0
genre	0	0	0	0	0	0	0	0	0	0	0
publisher	58	0	0	0	0	0	0	0	0	0	0
north_american_sales	0	0	0	0	0	0	0	0	0	0	0
european_sales	0	0	0	0	0	0	0	0	0	0	0
japanese_sales	0	0	0	0	0	0	0	0	0	0	0
other_sales	0	0	0	0	0	0	0	0	0	0	0
global_sales	0	0	0	0	0	0	0	0	0	0	0
dtype:	int64	object	object	object	object	object	float64	float64	float64	float64	float64

Two features have null values, year and publisher, there are multiple ways to deal with null values, I am going to drop the rows that contain them.

```
In [163]: #Dropping Null values
df.dropna(subset=['year', 'publisher'], inplace=True)

In [164]: #Confirming null values have been dropped
df.isnull().sum()
```

	rank	name	platform	year	genre	publisher	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
rank	0	0	0	0	0	0	0	0	0	0	0
name	0	271	0	0	0	0	0	0	0	0	0
platform	0	0	0	0	0	0	0	0	0	0	0
year	0	0	0	0	0	0	0	0	0	0	0
genre	0	0	0	0	0	0	0	0	0	0	0
publisher	0	0	0	0	0	0	0	0	0	0	0
north_american_sales	0	0	0	0	0	0	0	0	0	0	0
european_sales	0	0	0	0	0	0	0	0	0	0	0
japanese_sales	0	0	0	0	0	0	0	0	0	0	0
other_sales	0	0	0	0	0	0	0	0	0	0	0
global_sales	0	0	0	0	0	0	0	0	0	0	0
dtype:	int64	object	object	object	object	object	float64	float64	float64	float64	float64

```
In [166]: #Verifying the new shape of the dataset after null values and rows were dropped
df.shape

Out[166]: (16291, 11)
```

```
In [167]: #Check for duplicate data
duplicate = df.duplicated()
print(duplicate.sum())
df[duplicate]
```

	rank	name	platform	year	genre	publisher	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
0	0	0	0	0	0	0	0	0	0	0	0

There are no duplicates recorded, therefore nothing further needs to be done

Outlier Detection and Handling

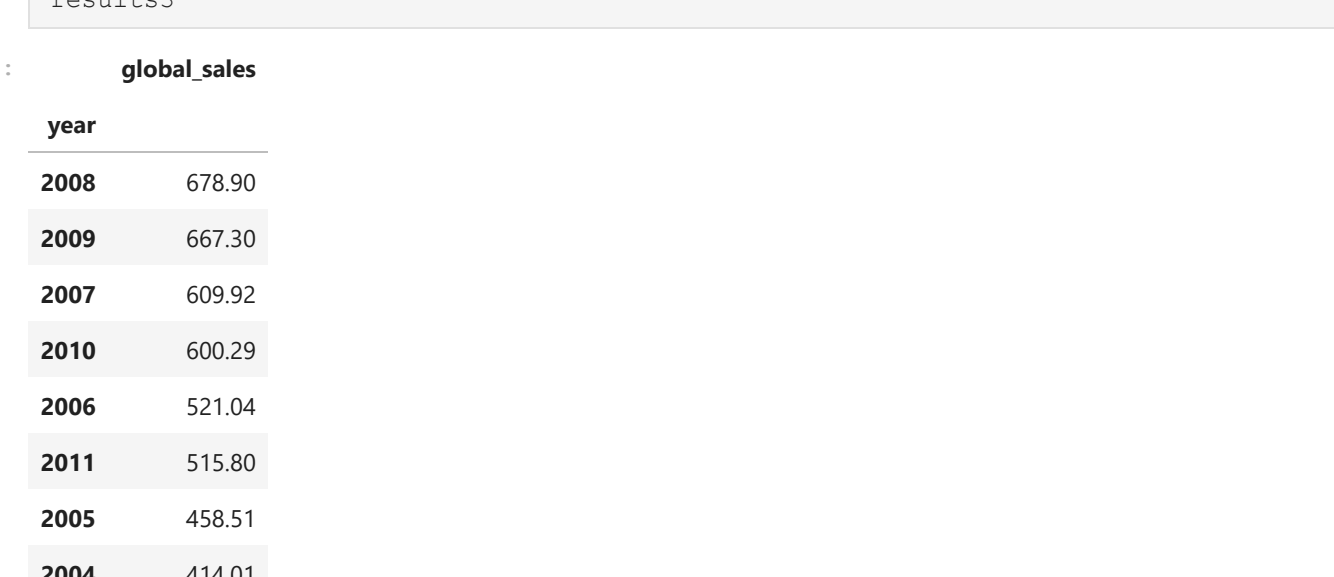
- Feature Engineering

Feature engineering is the initial stage in a machine learning channel and it comprises all the techniques utilised to clean existing datasets, improve their signal-noise ratio, and reduce their dimensionality (Bonaccorso, 2017).

```
In [22]: # Create a list for numerical columns that are to be visualized
Column_List = ['north_american_sales', 'european_sales', 'japanese_sales', 'other_sales', 'global_sales']

In [23]: # Generate whisker plots to detect the presence of any outliers
fig, ax = plt.subplots(len(Column_List), figsize=(10, 20))

for i, col_list in enumerate(Column_List):
    sns.boxplot(df[col_list], ax=ax[i], palette="winter", orient="h")
    ax[i].set_title(f'Whisker Plot for Outlier Detection on {col_list}')
    ax[i].set_ylabel(f'Number of Sales')
    ax[i].set_xlabel(f'Number of Sales')
    fig.tight_layout(pad=1.1)
```



A Whisker Plot function was used to examine the presence of outliers in the categorical data. The findings show there are numerous outliers present in each variable. The Z score method can only be used with normally distributed data, a check will have to be done to see if the data is so.

```
In [168]: # In this scenario, the outliers are removed using Z-Score due to the variability in
df = df[(np.abs(stats.zscore(df[['north_american_sales', 'european_sales',
                                'japanese_sales', 'other_sales',
                                'global_sales']))) < 3).all(axis = 1)] # all
df = df.reset_index() # Due to elimination of rows, index has to be reset
```

```
In [169]: #Compare the new dimension with the old one, the old shape of data set was (16598, 11)
#so we see 777 values have been removed.
df.shape

In [26]: #Change dates to remove decimal points, maybe truncate is better and then to pd datet.
df['year'] = df['year'].astype(int)
```

```
In [27]: #Placing Rank as the Index
df.reset_index(inplace=True)
df.head()
```

	rank	name	platform	year	genre	publisher	north_american_sales	european_sales	japanese_sales	other_sales	global_sales
0	1	Wii Sports	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	84.66
1	2	Super Mario Bros.	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	39.24
2	3	Mario Kart Wii	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	33.83
3	4	Wii Sports Resort	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	32.99
4	5	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo	11.27	8.89			

[illegible]

