# **Video Games Sales 2019**

Daniel Gutiérrez Rodríguez

1563389

# Índice

# Análisis dataset

· Multiples valores NaN y Null

·Ventas divididas en 2 columnas

```python
N_null = dataset.isnull().sum()
N_null.sort_values(inplace=True, ascending=False)
N_null.plot(kind='bar',stacked=True,figsize=(20,10))
sns.despine(left=True, bottom=True)
print(N_null)
plt.show()
sns.heatmap(dataset.isnull(), cbar=False)
```

```
VGChartz_Score    19862
User_Score        19624
Vgchartzscore     19335
Last_Update       15192
Critic_Score      15156
JP_Sales          13071
PAL_Sales          7737
NA_Sales           7085
ESRB_Rating        5937
Other_Sales        5352
Year                  3
Developer             2
Platform              0
Name                  0
basename              0
Genre                 0
Total_Sales           0
Publisher             0
img_url               0
url                   0
status                0
Rank                  0
dtype: int64
```

| Name | basename | Genre | ESRB_Rating | Platform | Publisher | Developer |
|---|---|---|---|---|---|---|
| Wii Sports | wii-sports | Sports | E | Wii | Nintendo | Nintendo EAD |
| Super Mario Bros. | super-mario-bros | Platform | NaN | NES | Nintendo | Nintendo EAD |
| Mario Kart Wii | mario-kart-wii | Racing | E | Wii | Nintendo | Nintendo EAD |
| PlayerUnknown's Battlegrounds | playerunknowns-battlegrounds | Shooter | NaN | PC | PUBG Corporation | PUBG Corporation |
| Wii Sports Resort | wii-sports-resort | Sports | E | Wii | Nintendo | Nintendo EAD |

# Preparación dataset

· Eliminación de las variables

· Normalización de los datos

```python
dataset = dataset.dropna()
dataset["Year"].fillna(dataset["Year"].mode(), inplace=True)
dataset["Developer"].fillna(dataset["Developer"].mode(), inplace=True)
N_null = dataset.isnull().sum()
N_null.sort_values(inplace=True, ascending=False)
print(N_null)
```

```
Year            3
Developer       2
Total_Sales     0
status          0
Publisher       0
Platform        0
Genre           0
Name            0
Rank            0
dtype: int64
```

```
Total_Sales     0
status          0
Year            0
Developer       0
Publisher       0
Platform        0
Genre           0
Name            0
Rank            0
dtype: int64
```
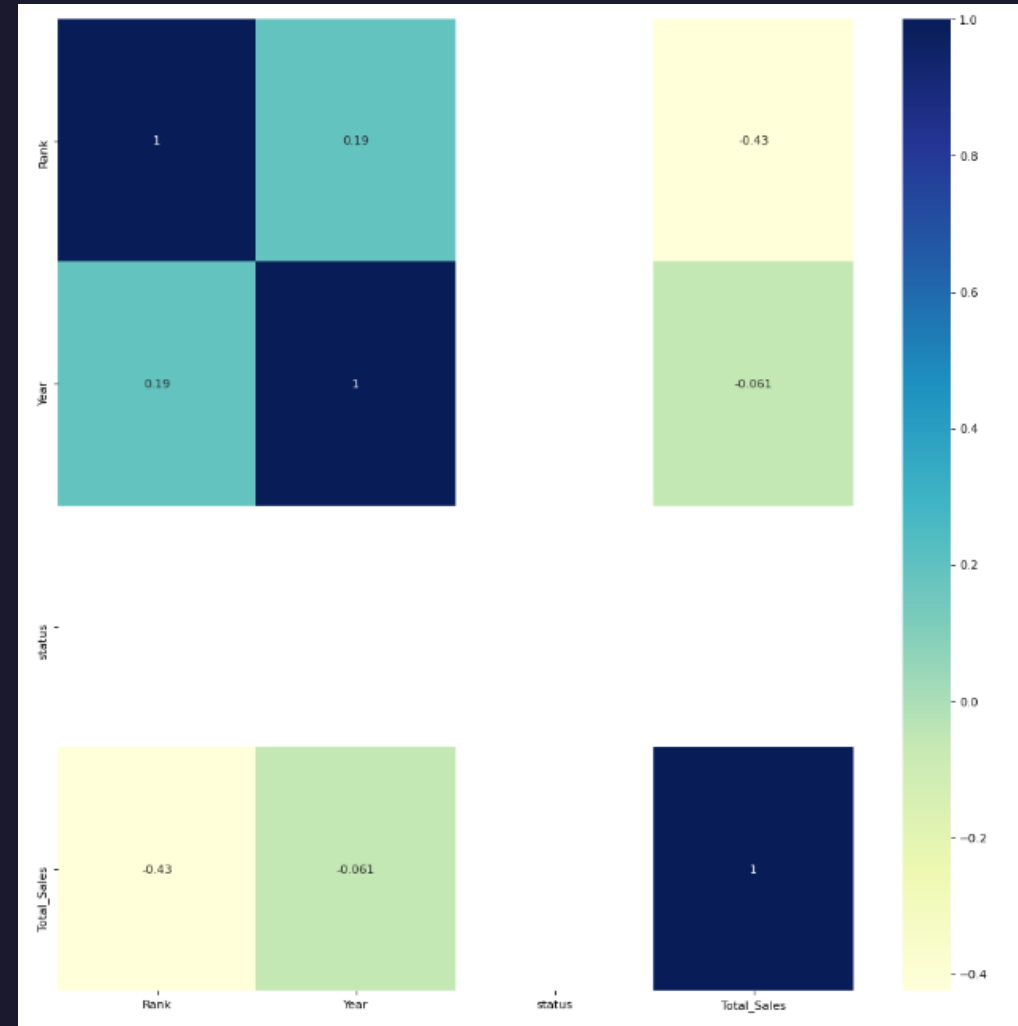
# Preparación dataset

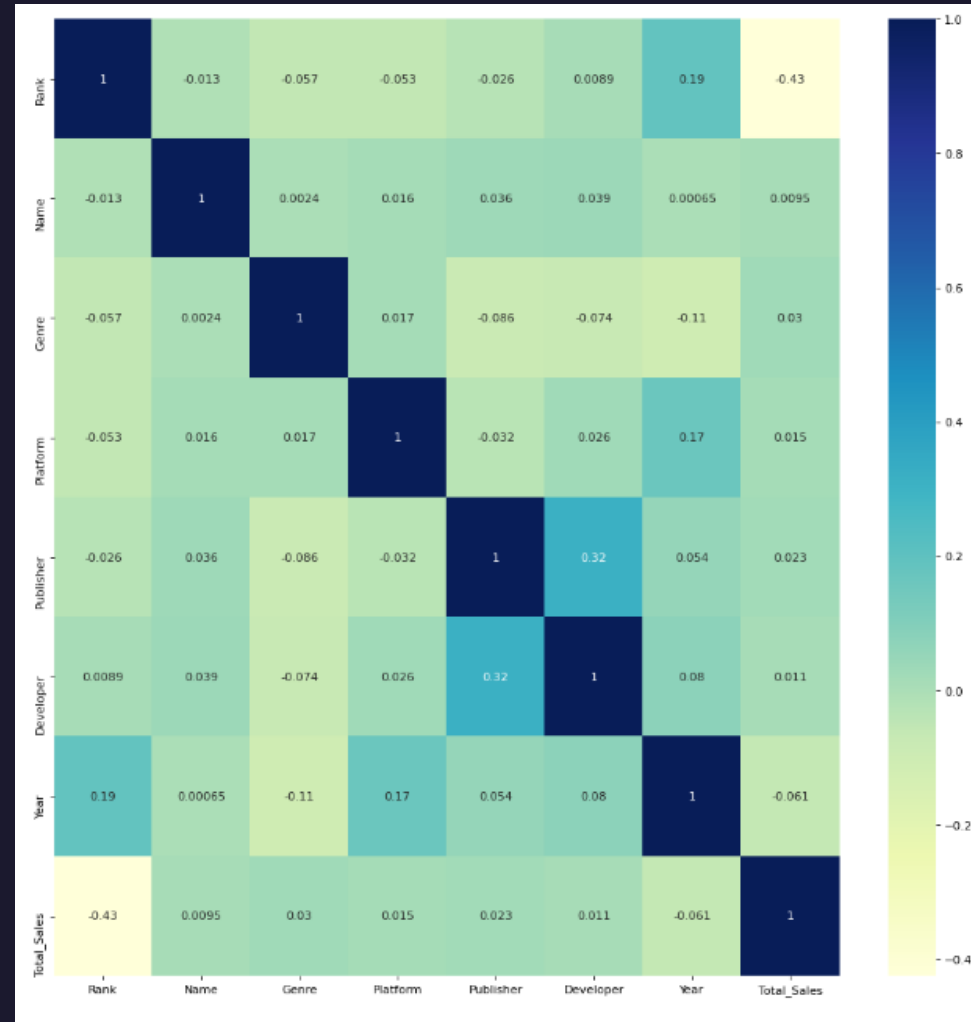- · Falta de algunas variables

- · Analisis variable "Status"

```
LabEncoder = preprocessing.LabelEncoder()
atributes = ['Name', 'Genre', 'Platform','Publisher', 'Developer']
for atribute in atributes:
    LabEncoder.fit(dataset[atribute])
    dataset[atribute] = LabEncoder.transform(dataset[atribute])
print(dataset.head())
```

```
   Rank   Name  Genre  Platform  Publisher  Developer    Year  status  \
0     1  13258     17        34        564       1931  2006.0       1
1     2  11212     10        14        564       1931  1985.0       1
2     3   6713     12        34        564       1931  2008.0       1
3     4   8782     15        18        602       2031  2017.0       1
4     5  13260     17        34        564       1931  2009.0       1

   Total_Sales
0        82.86
1        40.24
2        37.14
3        36.60
4        33.09
```

# Preparación dataset

# Entrenamiento

```
Regresión Logística
F1 score:   0.0705895626921464
C: 0.1

Regresión Logística
F1 score:   0.07159514680622799
C: 10

Regresión Logística
F1 score:   0.07160166804305684
C: 100

Regresión Logística
F1 score: 0.07160166804305684
C: 1000
```

· Decision Tree

· Logistic Regression

```
dt = DecisionTreeClassifier(random_state=0, criterion='gini')
dt.fit(X_train, Y_train)
print ("F1 score: ", f1_score(Y_test, dt.predict(X_test), average='macro'))
print("Criterion:", 'Gini')
print("")

dt = DecisionTreeClassifier(random_state=0, criterion='entropy')
dt.fit(X_train, Y_train)
print ("F1 score: ", f1_score(Y_test, dt.predict(X_test), average='macro'))
print("Criterion:", 'Entropy')
print("")
```

```
F1 score:   0.3609211428383207
Criterion: Gini

F1 score:   0.32092088342272457
Criterion: Entropy
```

# Validación



|   | Modelo | Media | std |
|---|---|---|---|
| 0 | Decision Tree | 0.378563 | 0.033258 |
| 1 | Regresión Logística | 0.075309 | 0.005960 |

```
F1 score:   0.3609211428383207
Criterion: Gini

F1 score:   0.32092088342272457
Criterion: Entropy
```

```
Regresión Logística
F1 score:   0.0705895626921464
C: 0.1

Regresión Logística
F1 score:   0.07159514680622799
C: 10

Regresión Logística
F1 score:   0.07160166804305684
C: 100

Regresión Logística
F1 score: 0.07160166804305684
C: 1000
```

· Decision Tree

· Logistic Regression

# Conclusión

# Gracias

Daniel Gutiérrez

1563389