# Лабораторная работа №2

## Обработка признаков (часть 1).

Рассмотрим исторические данные по выходу и продажам видеоигр из [на 2019 год](#)

## Задание:

[Оригинал](#):

- Выбрать набор данных (датасет), содержащий категориальные и числовые признаки и пропуски в данных. Для выполнения следующих пунктов можно использовать несколько различных наборов данных (один для обработки пропусков, другой для категориальных признаков и т.д.) Просьба не использовать датасет, на котором данная задача решалась в лекции.
- Для выбранного датасета (датасетов) на основе материалов лекций решить следующие задачи:
  - устранение пропусков в данных;
  - кодирование категориальных признаков;
  - нормализацию числовых признаков.

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

```
1  df = pd.read_csv('vgsales-12-4-2019-short.csv')
2  print(f'Total loaded {len(df.index)} video games')
```

```
1  Total loaded 55792 video games
```

```
1  # Объединим продажи в одну колонку
2  df['Sales'] = df[['Total_Shipped', 'Global_Sales']].max(axis=1)
3  df['Sales'].fillna(0, inplace=True)
```

```
1  # Удалим лишнее
2  df = df.drop(columns = ['Total_Shipped', 'Global_Sales', 'NA_Sales', 'PAL_Sales', 'JP_Sales',
   'Other_Sales'])
```

## Добавим выброс

```
1  df.tail(1)
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

| | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 55791 | 55792 | Falcon Age | Action-Adventure | NaN | PS4 | Unknown | Outerloop Games | NaN | NaN | NaN | 0.0 |

```
1  fake_row = (66666, 'FakeGame', 'Music', 'E', 'NS', 'Drool', 'Drool', 29.0, 64.0, 2018.0, 0.0)
2  appendix = pd.DataFrame((fake_row,), columns=df.columns)
3  df = df.append(appendix, ignore_index=True)
4  df.tail(2)
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

| | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 55791 | 55792 | Falcon Age | Action-Adventure | NaN | PS4 | Unknown | Outerloop Games | NaN | NaN | NaN | 0.0 |
| 55792 | 66666 | FakeGame | Music | E | NS | Drool | Drool | 29.0 | 64.0 | 2018.0 | 0.0 |

## Устранение пропусков в данных

```
1  # Подсчитаем NaN-ы по каждой колонке
2  df.isna().sum()
```

```
1  Rank              0
2  Name              0
3  Genre             0
4  ESRB_Rating   32169
5  Platform          0
6  Publisher         0
7  Developer        17
8  Critic_Score  49256
9  User_Score    55457
10 Year            979
11 Sales             0
12 dtype: int64
```

```
1  # Удалим неизвестные игры
2  df.dropna(subset=['Developer', 'Year'], inplace=True)
```

```
1  # Слишком большой перевес NaN-ов относительно настоящих оценок. Лучше удалим
2  df.dropna(subset=['Critic_Score', 'User_Score'], inplace=True)
```

```
1  df['ESRB_Rating'].unique()
```

```
1  array(['E', 'M', 'T', 'E10', nan], dtype=object)
```

```
1  # «E» («Everyone») — «Для всех»
2  # Заменим пропуски на значение по умолчанию
3  df['ESRB_Rating'].fillna('E', inplace=True)
```

```
1  # Заменим оценки критиков и пользователей на средние значения
2  for col in ('Critic_Score', 'User_Score'):
3      df[col].fillna(df[col].mean(), inplace=True)
```

```
1  # Подсчитаем NaN-ы по каждой колонке
2  df.isna().sum()
```

```
1  Rank           0
2  Name           0
3  Genre          0
4  ESRB_Rating    0
5  Platform       0
6  Publisher      0
7  Developer      0
8  Critic_Score   0
9  User_Score     0
10 Year           0
11 Sales          0
12 dtype: int64
```

## Кодирование категориальных признаков

```
1  # One-Hot
2  pd.get_dummies(df[['ESRB_Rating']])
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

|       | ESRB_Rating_E | ESRB_Rating_E10 | ESRB_Rating_M | ESRB_Rating_T |
|-------|---------------|-----------------|---------------|---------------|
| 2     | 1             | 0               | 0             | 0             |
| 4     | 1             | 0               | 0             | 0             |
| 6     | 1             | 0               | 0             | 0             |
| 8     | 1             | 0               | 0             | 0             |
| 11    | 1             | 0               | 0             | 0             |
| ...   | ...           | ...             | ...           | ...           |
| 48498 | 1             | 0               | 0             | 0             |
| 51097 | 0             | 1               | 0             | 0             |
| 54537 | 1             | 0               | 0             | 0             |
| 55423 | 0             | 1               | 0             | 0             |
| 55792 | 1             | 0               | 0             | 0             |

219 rows × 4 columns

```
1  from sklearn.preprocessing import LabelEncoder
2
3  for col in ('Genre', 'Platform', 'Publisher', 'Developer'):
4      le = LabelEncoder()
5      col_le = le.fit_transform(df[col])
6      print(np.unique(df[col]), np.unique(col_le))
7      df[col] = col_le
8
9  df.head()
```

```
1  ['Action' 'Action-Adventure' 'Adventure' 'Fighting' 'MMO' 'Misc' 'Music'
2   'Platform' 'Puzzle' 'Racing' 'Role-Playing' 'Shooter' 'Simulation'
3   'Sports' 'Strategy'] [ 0  1  2  3  4  5  6  7  8  9 10 11 12 13 14]
4  ['3DS' 'DS' 'GC' 'N64' 'NS' 'PC' 'PS' 'PS2' 'PS3' 'PS4' 'PSN' 'PSP' 'WW'
5   'Wii' 'WiiU' 'X360' 'XB' 'XBL' 'XOne'] [ 0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18]
6  ['2K Games' 'Activision' 'Aksys Games' 'Atari' 'Atlus'
7   'Bethesda Softworks' 'Blizzard Entertainment' 'Capcom' 'DotEmu' 'Drool'
8   'EA Sports' 'Electronic Arts' 'FDG Entertainment' 'Feardemic'
9   'Hudson Entertainment' 'Ignition Entertainment' 'Invisible Handlebar'
10  'Konami' 'LucasArts' 'Mediatonic' 'Microprose' 'Microsoft Game Studios'
11  'Midway Games' 'NIS America' 'Namco Bandai' 'Natsume' 'Nicalis'
```

```
12    'Nintendo' 'Paradox Interactive' 'RedOctane' 'Rockstar Games' 'Sega'
13    'Sony Computer Entertainment' 'Sony Interactive Entertainment' 'Square'
14    'Square Enix' 'Strategy First' 'THQ' 'Team17 Software' 'Ubisoft'
15    'Valve Corporation' 'Warner Bros. Interactive' 'Working Designs'
16    'Xseed Games'] [ 0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
17    24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43]
18   ['2K Australia / 2K Boston' '2K Boston / 2K Australia / 2K Marin'
19    '989 Sports' 'Access Games' 'Alfa System' 'Arc System Works' 'Atlus Co.'
20    'Bandai Namco Games' 'BestGameEver.com' 'Bethesda Game Studios'
21    'Bethesda Softworks' 'BioWare Edmonton' 'Blizzard Entertainment'
22    'Blue Tongue' 'Bright Light Productions' 'Bungie Studios' 'Capcom'
23    'Capcom / Dimps Corporation' 'Capcom Production Studio 1'
24    'Capcom Production Studio 1 / TOSE Software' 'Capcom Production Studio 4'
25    'Cavia Inc.' 'Certain Affinity / Valve Software' 'Choice Provisions'
26    'Chris Sawyer' 'Cing' 'Cing Inc. / Town Factory' 'Climax Group'
27    'Climax Studios' 'Clover Studio' 'Crytek' 'Digital Eclipse' 'Drool'
28    'EA Canada' 'EA Digital Illusions CE' 'EA Montreal' 'EA Vancouver'
29    'EXAKT Entertainment' 'Ellipse Studios' 'Ensemble Studios' 'Epic Games'
30    'Evolution Studios' 'Firaxis' 'Flagship' 'From Software' 'GRIN'
31    'Game Arts' 'Game Atelier' 'Game Freak' 'Game Republic' 'Good-Feel'
32    'Grasshopper Manufacture' 'Guerrilla Games' 'HAL Laboratory'
33    'High Voltage Software' 'Hudson Entertainment' 'Infinity Ward'
34    'Insomniac Games' 'Intelligent Systems' 'IronNos Co., Ltd.'
35    'Just Add Water' 'K2 / Kurogane' 'Kojima Productions' 'Konami'
36    'Konami / Eighting' 'Konami Computer Entertainment Japan'
37    'Konami Computer Entertainment Tokyo' 'Krome Studios' 'Level 5'
38    'Lionhead Studios' 'Lizardcube' 'MPS Labs' 'Matrix Software'
39    'Media Molecule' 'Mediatonic Ltd.'
40    'Mistwalker Corporation / Feelplus Inc.' 'Monolith Productions'
41    'Monolith Soft' 'Namco' 'Namco Tales Studio' 'Naughty Dog' 'Neversoft'
42    'Nihon Falcom Corporation' 'Nintendo EAD' 'Nintendo EAD Tokyo'
43    'Nintendo EPD' 'Nintendo Software Technology Corporation'
44    'Nippon Ichi Software' 'Omega Force' 'PlatinumGames' 'Polyphony Digital'
45    'Project Sora' 'Quantic Dream' 'Rainbow Studios' 'Rare Ltd.'
46    'Raven Software' 'Ready at Dawn' 'Ready at Dawn Studios'
47    'Realtime Worlds' 'Retro Studios' 'Rhino Studios' 'Rockstar Leeds/North'
48    'Rockstar North' 'Rockstar San Diego' 'SCEA Bend Studio'
49    'SCEA Santa Monica Studio' 'SCEI' 'SFB Games' 'Secret Level' 'Sega'
50    'Sega WOW Overworks' 'Sonic Team' 'Sonic Team/Dimps Corporation'
51    'Spike / Bandai Namco Games' 'Square Enix' 'SquareSoft'
52    'Sucker Punch Productions' 'Taleworlds' 'Team Bondi' 'Team ICO'
53    'Team Silent' 'Team17 Software' 'The Deep End Games' 'Treyarch'
54    'Turn 10 Studios' 'Ubisoft' 'Ubisoft Montreal' 'Ubisoft Paris'
55    'Ubisoft Romania' 'Ubisoft San Francisco' 'Ubisoft Shanghai'
56    'Ubisoft Toronto' 'Valhalla Game Studios' 'Valve Software' 'Vanillaware'
57    'Vicarious Visions' 'Vigil Games' 'Visceral Games'
58    "Yuke's Future Media Creators" 'Zipper Interactive' 'id Software'
59    'n-Space' 'thatgamecompany' 'tri-Ace'] [  0   1   2   3   4   5   6   7   8   9  10  11  12  13  14  15  16  17
60     18  19  20  21  22  23  24  25  26  27  28  29  30  31  32  33  34  35
61     36  37  38  39  40  41  42  43  44  45  46  47  48  49  50  51  52  53
62     54  55  56  57  58  59  60  61  62  63  64  65  66  67  68  69  70  71
63     72  73  74  75  76  77  78  79  80  81  82  83  84  85  86  87  88  89
64     90  91  92  93  94  95  96  97  98  99 100 101 102 103 104 105 106 107
65    108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125
66    126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143]
```

```
1   .dataframe tbody tr th {
2       vertical-align: top;
3   }
4
5   .dataframe thead th {
6       text-align: right;
7   }
```

|    | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales |
|----|------|------|-------|-------------|----------|-----------|-----------|--------------|------------|------|-------|
| 2  | 3    | Mario Kart Wii | 9 | E | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | 37.14 |
| 4  | 5    | Wii Sports Resort | 13 | E | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | 33.09 |
| 6  | 7    | New Super Mario Bros. | 7 | E | 1 | 27 | 83 | 9.1 | 8.1 | 2006.0 | 30.80 |
| 8  | 9    | New Super Mario Bros. Wii | 7 | E | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | 30.22 |
| 11 | 12   | Wii Play | 5 | E | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | 28.02 |

```
1   # Static label encoding
2   di = {
3       'RP':  0,
4       'EC':  1,
5       'KA':  2,
6       'E':   3,
7       'E10': 4,
8       'T':   5,
9       'M':   6,
10      'AO':  7,
11  }
12
13  df.replace({'ESRB_Rating': di}, inplace=True)
14  df.head()
```

```
1   .dataframe tbody tr th {
2       vertical-align: top;
3   }
4
5   .dataframe thead th {
6       text-align: right;
7   }
```

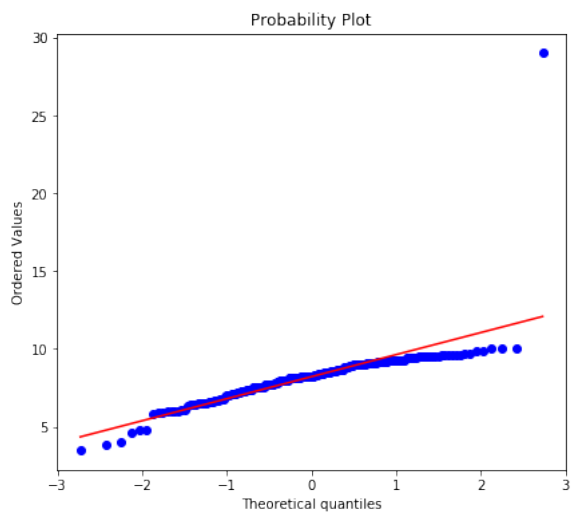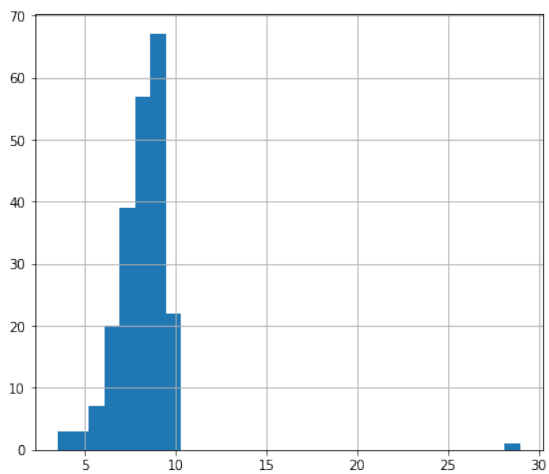|    | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales |
|----|------|------|-------|-------------|----------|-----------|-----------|--------------|------------|------|-------|
| 2  | 3  | Mario Kart Wii        | 9  | 3 | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | 37.14 |
| 4  | 5  | Wii Sports Resort     | 13 | 3 | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | 33.09 |
| 6  | 7  | New Super Mario Bros. | 7  | 3 | 1  | 27 | 83 | 9.1 | 8.1 | 2006.0 | 30.80 |
| 8  | 9  | New Super Mario Bros. Wii | 7 | 3 | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | 30.22 |
| 11 | 12 | Wii Play              | 5  | 3 | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | 28.02 |

## Нормализация числовых признаков

```
1   import scipy.stats as stats
```
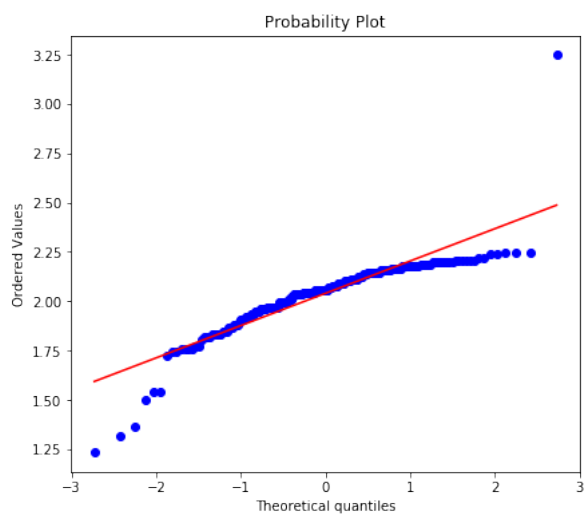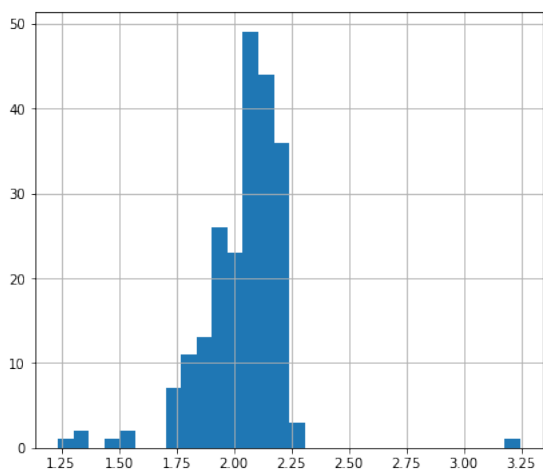
```
1   def diagnostic_plots(df, variable):
2       plt.figure(figsize=(15,6))
3       # гистограмма
4       plt.subplot(1, 2, 1)
5       df[variable].hist(bins=30)
6       ## Q-Q plot
7       plt.subplot(1, 2, 2)
8       stats.probplot(df[variable], dist="norm", plot=plt)
9       plt.show()
```

```
1   def boxcox(data, col):
2       diagnostic_plots(data, col)
3       data[f'{col}_boxcox'], param = stats.boxcox(data[col])
4       print('Оптимальное значение λ = {}'.format(param))
5       diagnostic_plots(data, f'{col}_boxcox')
```
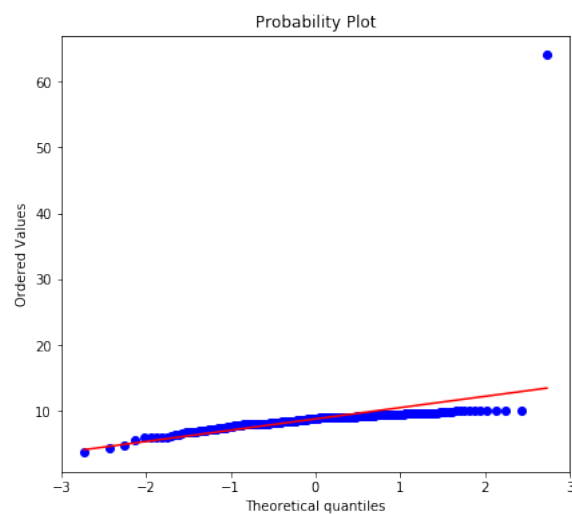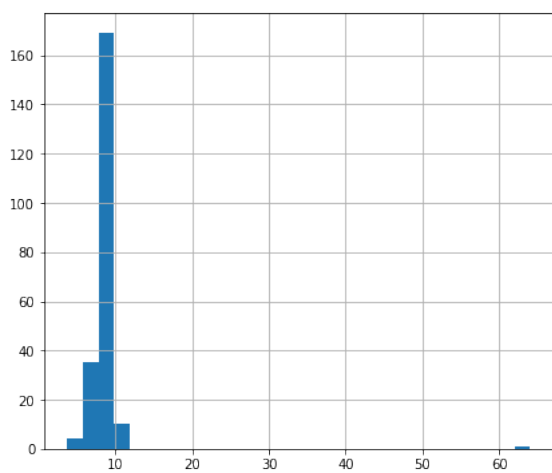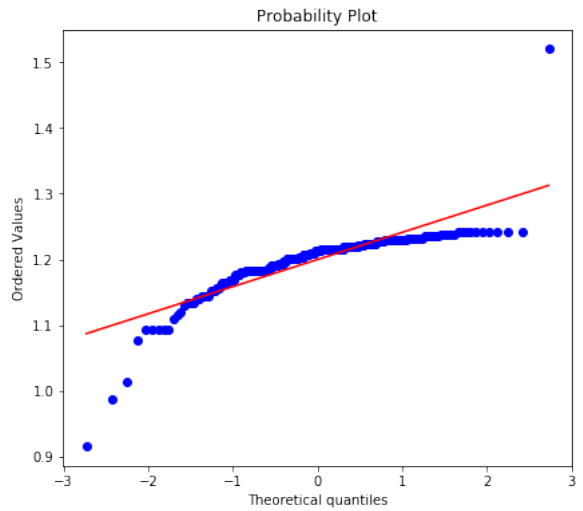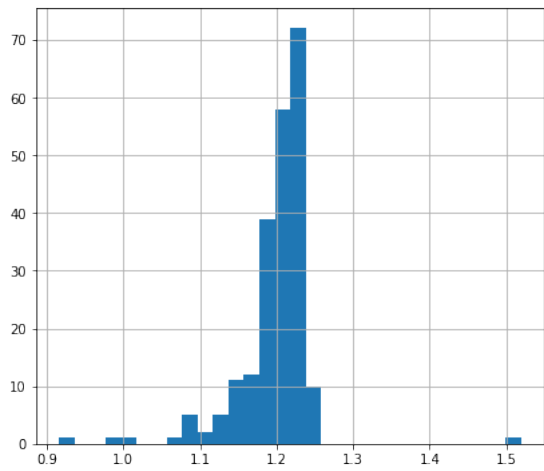
```
1   boxcox(df, 'Critic_Score')
```

Оптимальное значение λ = -0.021821216576747314



```
boxcox(df, 'User_Score')
```



Оптимальное значение λ = -0.6048001579673186

```
1  df.head()
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

|    | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales | Critic_Score_boxcox | Us |
|----|------|------|-------|-------------|----------|-----------|-----------|--------------|------------|------|-------|---------------------|-----|
| 2  | 3  | Mario Kart Wii | 9 | 3 | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | 37.14 | 2.056560 | 1.2 |
| 4  | 5  | Wii Sports Resort | 13 | 3 | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | 33.09 | 2.032969 | 1.2 |
| 6  | 7  | New Super Mario Bros. | 7 | 3 | 1 | 27 | 83 | 9.1 | 8.1 | 2006.0 | 30.80 | 2.155914 | 1.1 |
| 8  | 9  | New Super Mario Bros. Wii | 7 | 3 | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | 30.22 | 2.102027 | 1.2 |
| 11 | 12 | Wii Play | 5 | 3 | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | 28.02 | 1.741018 | 0.9 |

```
1  df.tail(1)
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

|       | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales | Critic_Score_boxco |
|-------|------|------|-------|-------------|----------|-----------|-----------|--------------|------------|------|-------|--------------------|
| 55792 | 66666 | FakeGame | 6 | 3 | 4 | 9 | 32 | 29.0 | 64.0 | 2018.0 | 0.0 | 3.246559 |

## Сохраняем

```
1  save_path = 'video_games_s2.csv'
2  df.to_csv(save_path, index=False)
```

```
1  check_df = pd.read_csv(save_path)
2  check_df.head()
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

| | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales | Critic_Score_boxcox | Use |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | Mario Kart Wii | 9 | 3 | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | 37.14 | 2.056560 | 1.21 |
| 1 | 5 | Wii Sports Resort | 13 | 3 | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | 33.09 | 2.032969 | 1.20 |
| 2 | 7 | New Super Mario Bros. | 7 | 3 | 1 | 27 | 83 | 9.1 | 8.1 | 2006.0 | 30.80 | 2.155914 | 1.18 |
| 3 | 9 | New Super Mario Bros. Wii | 7 | 3 | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | 30.22 | 2.102027 | 1.22 |
| 4 | 12 | Wii Play | 5 | 3 | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | 28.02 | 1.741018 | 0.98 |

```
1
```