**Факультет «Информатика и системы управления»**

**Кафедра ИУ5 «Системы обработки информации и управления»**

Отчёт по лабораторной работе
по дисциплине «Методы Машинного Обучения»

Выполнил:
студент группы № ИУ5-21М
Кучеренко Михаил Александрович
_____, _____

Проверил:
к.т.н., доц., Ю.Е. Гапанюк
_____, _____

2021 г.

# Лабораторная работа №3

## Обработка признаков (часть 2).

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

## Задание:

[Оригинал](#):

- Выбрать один или несколько наборов данных (датасетов) для решения следующих задач. Каждая задача может быть решена на отдельном датасете, или несколько задач могут быть решены на одном датасете. Просьба не использовать датасет, на котором данная задача решалась в лекции.
- Для выбранного датасета (датасетов) на основе материалов лекций решить следующие задачи:
  - масштабирование признаков (не менее чем тремя способами);
  - обработку выбросов для числовых признаков (по одному способу для удаления выбросов и для замены выбросов);
  - обработку по крайней мере одного нестандартного признака (который не является числовым или категориальным);
  - отбор признаков:
    - один метод из группы методов фильтрации (filter methods);
    - один метод из группы методов обертывания (wrapper methods);
    - один метод из группы методов вложений (embedded methods).

```python
import numpy as np
import pandas as pd
import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```python
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import RobustScaler
from sklearn.preprocessing import MaxAbsScaler
```

```python
save_path = 'video_games_s2.csv'
df = pd.read_csv(save_path)
print(f'Loaded {len(df)} games')
```

```
Loaded 219 games
```

```python
df.head()
```

```css
.dataframe tbody tr th {
    vertical-align: top;
}

.dataframe thead th {
    text-align: right;
}
```

| | Rank | Name | Genre | ESRB_Rating | Platform | Publisher | Developer | Critic_Score | User_Score | Year | Sales | Critic |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | Mario Kart Wii | 9 | 3 | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | 37.14 | 2.056 |
| 1 | 5 | Wii Sports Resort | 13 | 3 | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | 33.09 | 2.032 |
| 2 | 7 | New Super Mario Bros. | 7 | 3 | 1 | 27 | 83 | 9.1 | 8.1 | 2006.0 | 30.80 | 2.155 |
| 3 | 9 | New Super Mario Bros. Wii | 7 | 3 | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | 30.22 | 2.102 |
| 4 | 12 | Wii Play | 5 | 3 | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | 28.02 | 1.741 |

```python
# Построение плотности распределения
def draw_kde(col_list, df1, df2, label1, label2):
    fig, (ax1, ax2) = plt.subplots(
        ncols=2, figsize=(12, 5))
    # первый график
    ax1.set_title(label1)
    for col in col_list:
        sns.distplot(df1[col], ax=ax1)
    # второй график
    ax2.set_title(label2)
    for col in col_list:
        sns.distplot(df2[col], ax=ax2)
    plt.show()
```

```python
def draw_data(col_list, df, df_scaled):
    df[col_list].plot()
    df_scaled[col_list].plot()
    plt.show()
```

```python
def get_scaled(df, columns, scaler=StandardScaler()):
    data_scaled = scaler.fit_transform(df[columns])
    df_scaled = pd.DataFrame(data_scaled, columns=columns)
    draw_data(columns, df, df_scaled)
    draw_kde(columns, df, df_scaled, 'Before', 'After')
    return df_scaled
```

```python
def apply_scaled(df, df_scaled, columns):
    for col in columns:
        df[f'{col}_scaled'] = df_scaled[col]
    return df
```

```python
columns = ['Sales', 'ESRB_Rating', 'Genre', 'Platform', 'Publisher', 'Developer', 'Year', 'Critic_Score_boxcox', 'User_Score_boxcox']
```
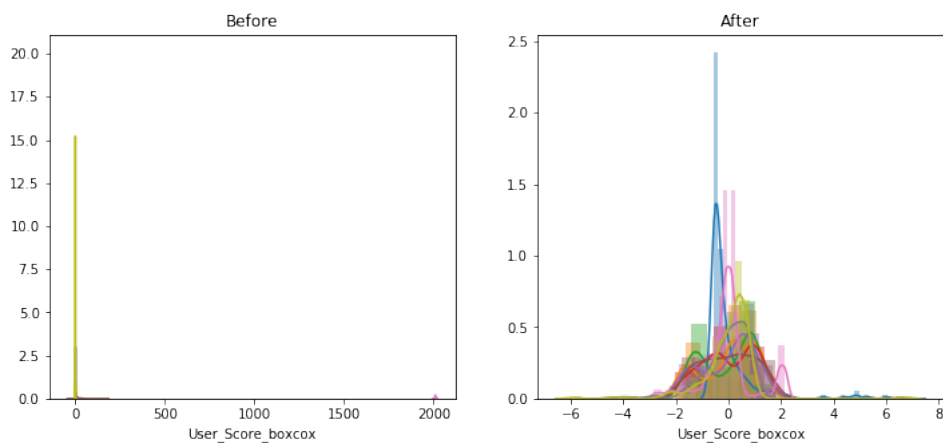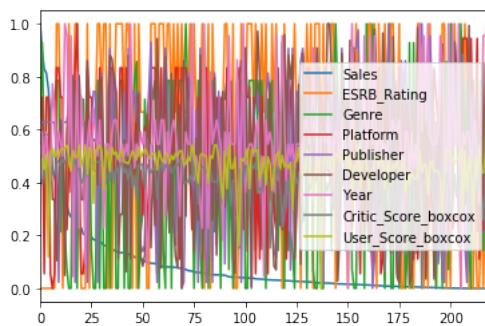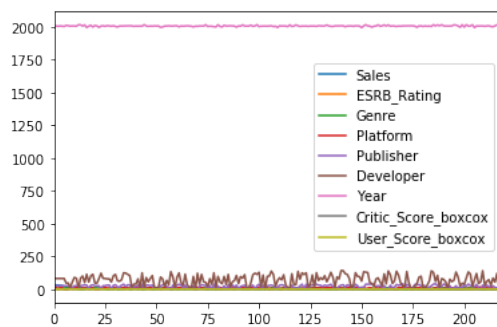
```python
df_scaled_std = get_scaled(df, columns, StandardScaler())
```
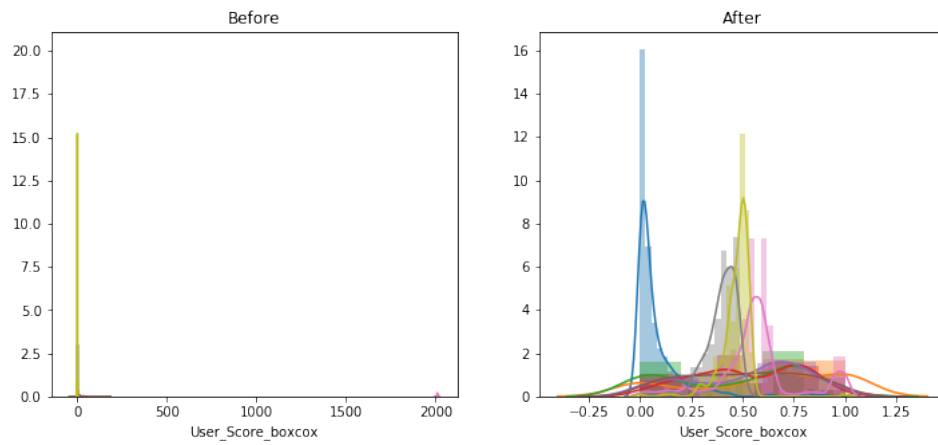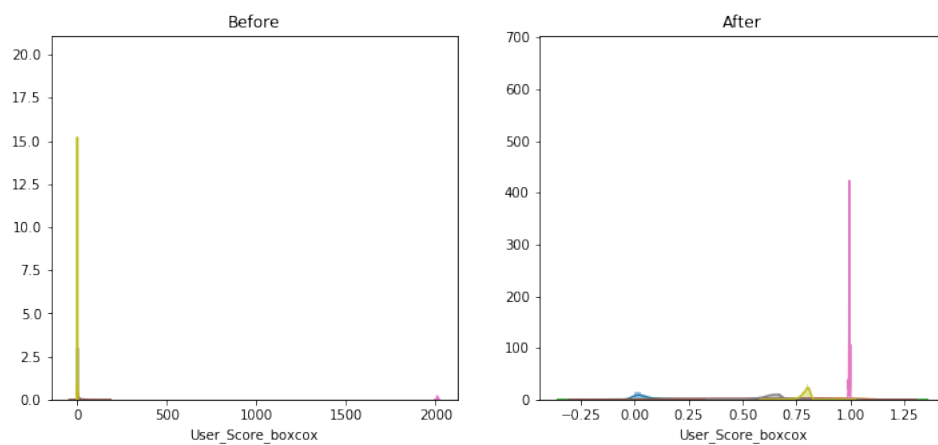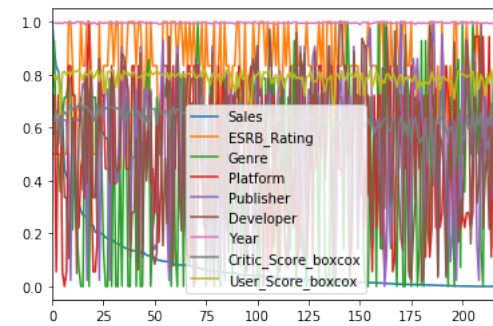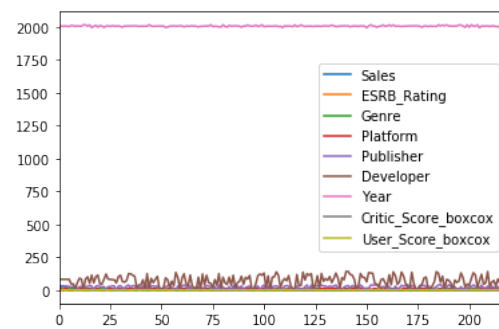
```
1   df_scaled_maxmin = get_scaled(df, columns, MinMaxScaler())
```

**Before** / **After**

```
1  df_scaled_maxabs = get_scaled(df, columns, MaxAbsScaler())
```





**Before** / **After**

```
1  df = apply_scaled(df, df_scaled_maxmin, columns)
2  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 | ... | User_Score_boxcox | Sales_sc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | Mario Kart Wii | 9 | 3 | 13 | 27 | 83 | 8.2 | 9.1 | 2008.0 | ... | 1.218568 | 1.000000 |
| 1 | 5 | Wii Sports Resort | 13 | 3 | 13 | 27 | 83 | 8.0 | 8.8 | 2009.0 | ... | 1.209662 | 0.890953 |
| 2 | 7 | New Super Mario Bros. | 7 | 3 | 1 | 27 | 83 | 9.1 | 8.1 | 2006.0 | ... | 1.186848 | 0.829295 |
| 3 | 9 | New Super Mario Bros. Wii | 7 | 3 | 13 | 27 | 83 | 8.6 | 9.2 | 2009.0 | ... | 1.221433 | 0.813678 |
| 4 | 12 | Wii Play | 5 | 3 | 13 | 27 | 83 | 5.9 | 4.5 | 2007.0 | ... | 0.987667 | 0.754443 |

5 rows × 22 columns

```
1  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 | ... | User_Score_boxcox | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 217 | 55424 | Thumper | 6 | 4 | 4 | 9 | 32 | 9.0 | 9.3 | 2017.0 | ... | 1.224249 | 0 |
| 218 | 66666 | FakeGame | 6 | 3 | 4 | 9 | 32 | 29.0 | 64.0 | 2018.0 | ... | 1.519776 | 0 |

2 rows × 22 columns

## Поиск и устранение выбросов

```
1   from enum import Enum
2   import scipy.stats as stats
3   from sklearn.svm import SVR
4   from sklearn.linear_model import LinearRegression
5   from sklearn.neighbors import KNeighborsRegressor
6   from sklearn.tree import DecisionTreeRegressor
7   from sklearn.ensemble import RandomForestRegressor
8   from sklearn.ensemble import GradientBoostingRegressor
9   from sklearn.metrics import mean_squared_error
10  from sklearn.model_selection import train_test_split
```

```
1   class OutlierBoundaryType(Enum):
2       SIGMA = 1
3       QUANTILE = 2
4       IRQ = 3
5
6   # Функция вычисления верхней и нижней границы выбросов
```

```python
 7   def get_outlier_boundaries(df, col, outlier_boundary_type: OutlierBoundaryType):
 8       if outlier_boundary_type == OutlierBoundaryType.SIGMA:
 9           K1 = 3
10           lower_boundary = df[col].mean() - (K1 * df[col].std())
11           upper_boundary = df[col].mean() + (K1 * df[col].std())
12
13       elif outlier_boundary_type == OutlierBoundaryType.QUANTILE:
14           lower_boundary = df[col].quantile(0.05)
15           upper_boundary = df[col].quantile(0.95)
16
17       elif outlier_boundary_type == OutlierBoundaryType.IRQ:
18           K2 = 1.5
19           IQR = df[col].quantile(0.75) - df[col].quantile(0.25)
20           lower_boundary = df[col].quantile(0.25) - (K2 * IQR)
21           upper_boundary = df[col].quantile(0.75) + (K2 * IQR)
22
23       else:
24           raise NameError('Unknown Outlier Boundary Type')
25
26       return lower_boundary, upper_boundary
```

```python
 1   def diagnostic_plots(df, variable, title):
 2       fig, ax = plt.subplots(figsize=(10,7))
 3       # гистограмма
 4       plt.subplot(2, 2, 1)
 5       df[variable].hist(bins=30)
 6       ## Q-Q plot
 7       plt.subplot(2, 2, 2)
 8       stats.probplot(df[variable], dist="norm", plot=plt)
 9       # ящик с усами
10       plt.subplot(2, 2, 3)
11       sns.violinplot(x=df[variable])
12       # ящик с усами
13       plt.subplot(2, 2, 4)
14       sns.boxplot(x=df[variable])
15       fig.suptitle(title)
16       plt.show()
```

```python
 1   scaled_columns = [f'{x}_scaled' for x in columns]
```

```python
 1   method_list = ['Original']
 2   df_changed = [df]
```

```python
 1   for obt in OutlierBoundaryType:
 2
 3       df2 = df.copy()
 4
 5       for col in columns:
 6           # Вычисление верхней и нижней границы
 7           lower_boundary, upper_boundary = get_outlier_boundaries(df2, col, obt)
 8           # Изменение данных
 9           df2[col] = np.where(df2[col] > upper_boundary, upper_boundary,
10                               np.where(df2[col] < lower_boundary, lower_boundary, df2[col]))
11
12       for col in scaled_columns:
13           # Вычисление верхней и нижней границы
14           lower_boundary, upper_boundary = get_outlier_boundaries(df2, col, obt)
15           # Изменение данных
16           df2[col] = np.where(df2[col] > upper_boundary, upper_boundary,
17                               np.where(df2[col] < lower_boundary, lower_boundary, df2[col]))
18
19       title = '{}-updated'.format(obt)
20       # Сохранение в списки
21       method_list.append(title)
22       df_changed.append(df2)
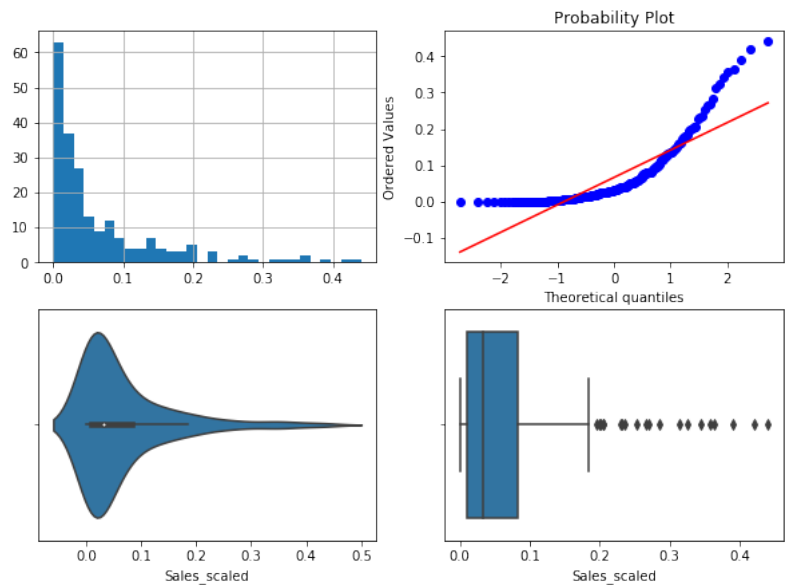```

```python
 1   df_changed[1].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   | ... | User_Score_boxcox | S |
|-----|-------|----------|-------|-------------|----------|-----------|-----------|--------------|------------|--------|-----|-------------------|---|
| 217 | 55424 | Thumper  | 6.0   | 4.0         | 4.0      | 9.0       | 32.0      | 9.0          | 9.3        | 2017.0 | ... | 1.224249          | 0 |
| 218 | 66666 | FakeGame | 6.0   | 3.0         | 4.0      | 9.0       | 32.0      | 29.0         | 64.0       | 2018.0 | ... | 1.344304          | 0 |

2 rows × 22 columns

```
1   for col in scaled_columns:
2       for obt in OutlierBoundaryType:
3           # Вычисление верхней и нижней границы
4           lower_boundary, upper_boundary = get_outlier_boundaries(df, col, obt)
5           # Флаги для удаления выбросов
6           outliers_temp = np.where(df[col] > upper_boundary, True,
7                               np.where(df[col] < lower_boundary, True, False))
8           # Удаление данных на основе флага
9           df_trimmed = df.loc[~(outliers_temp), ]
10          title = 'Поле-{}, метод-{}, строк-{}'.format(col, obt, df_trimmed.shape[0])
11          diagnostic_plots(df_trimmed, col, title)
```



Поле-Sales_scaled, метод-OutlierBoundaryType.SIGMA, строк-213
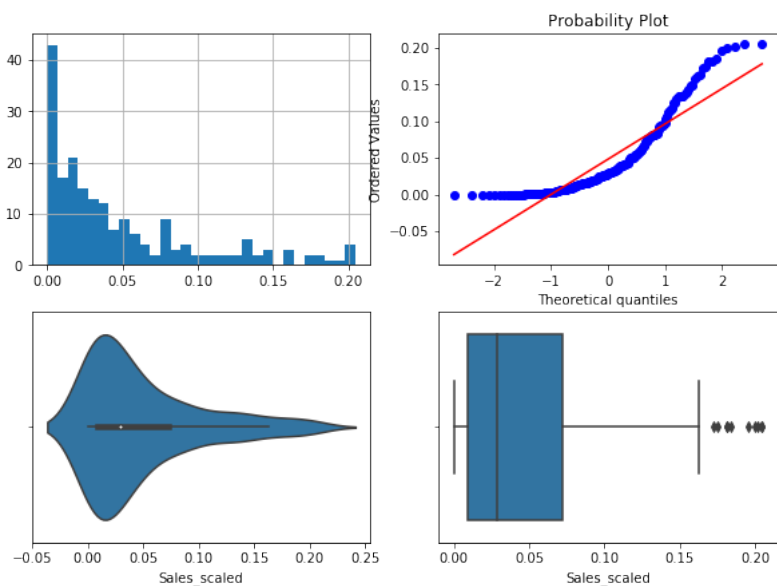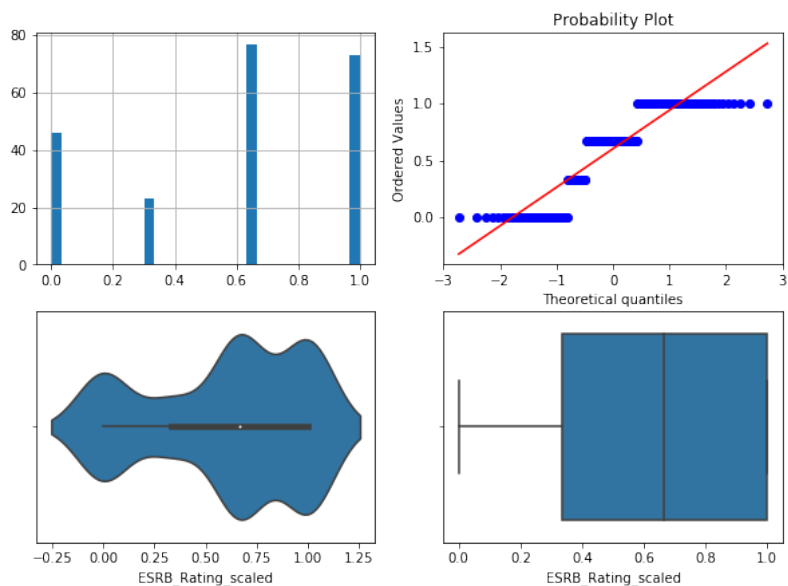
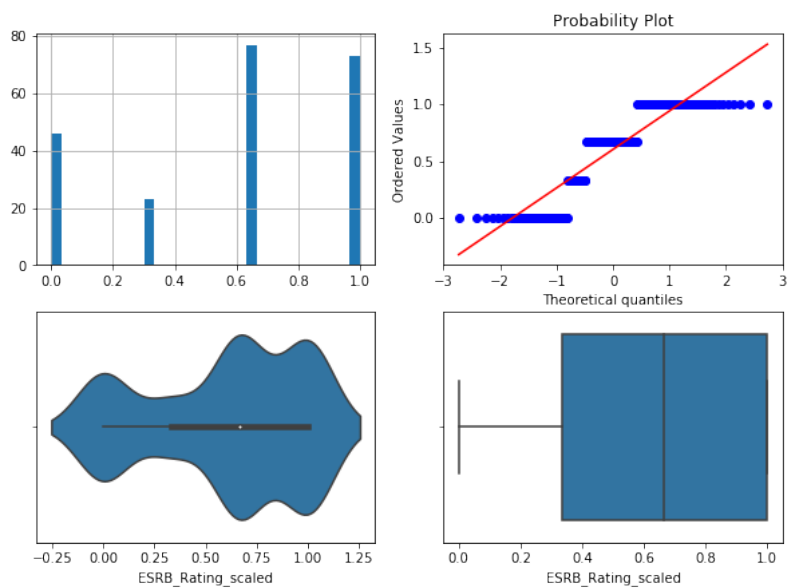Поле-Sales_scaled, метод-OutlierBoundaryType.QUANTILE, строк-208



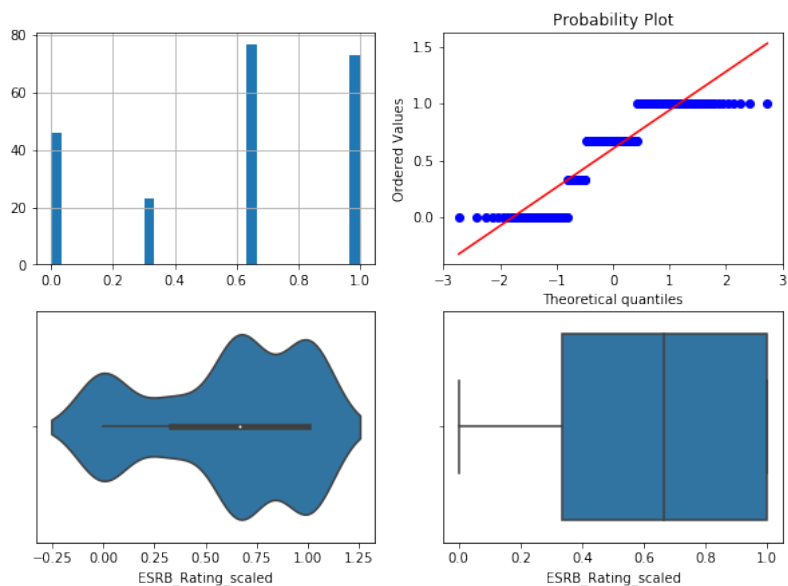Поле-Sales_scaled, метод-OutlierBoundaryType.IRQ, строк-198

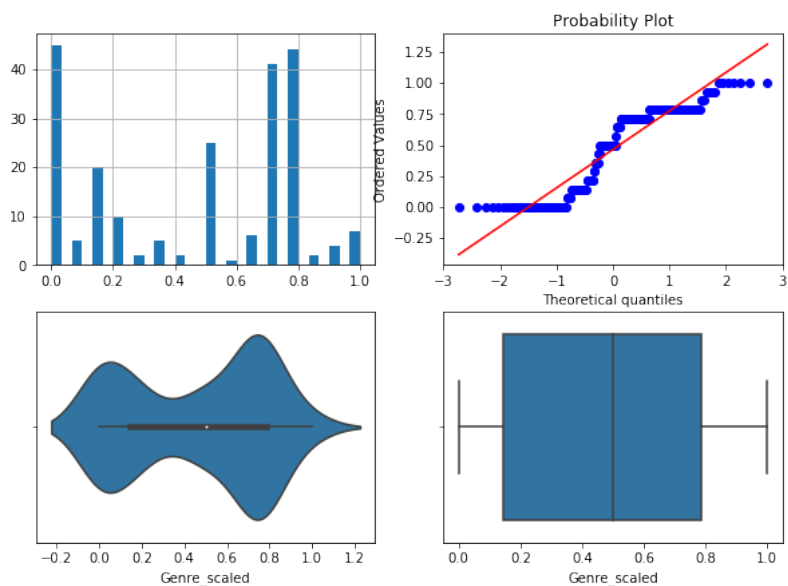Поле-ESRB_Rating_scaled, метод-OutlierBoundaryType.SIGMA, строк-219



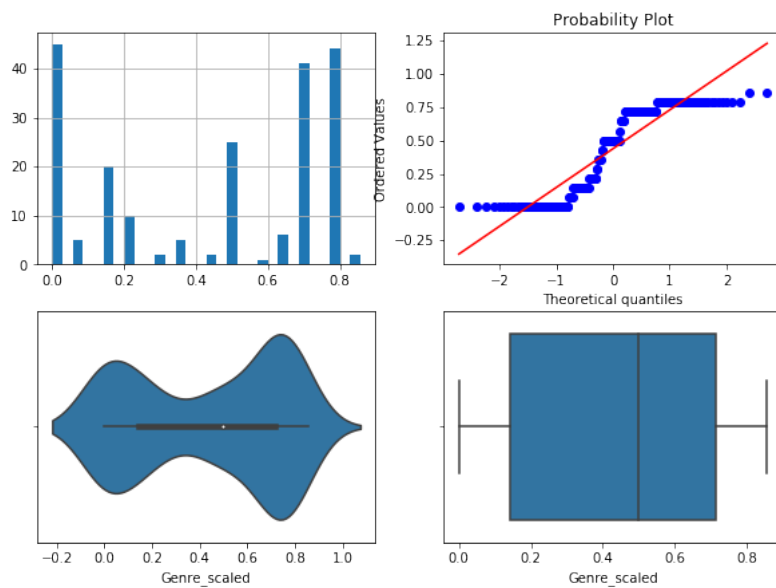Поле-ESRB_Rating_scaled, метод-OutlierBoundaryType.QUANTILE, строк-219

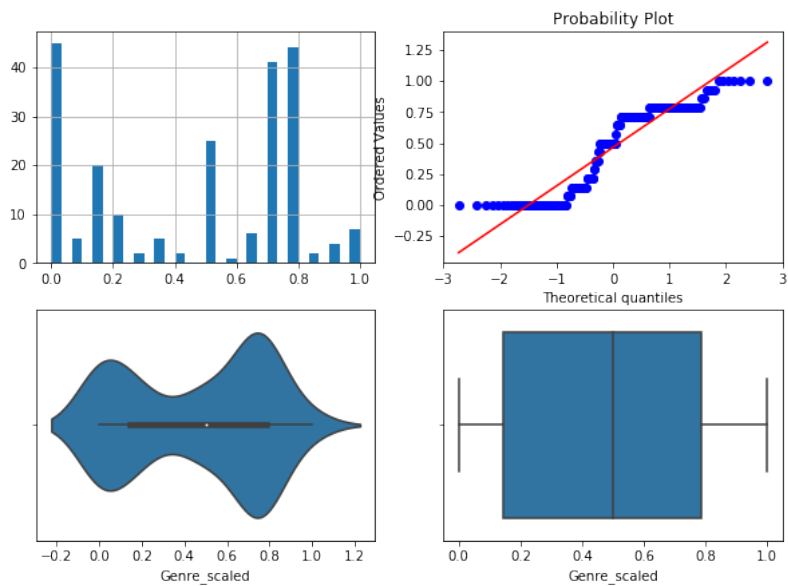Поле-ESRB_Rating_scaled, метод-OutlierBoundaryType.IRQ, строк-219


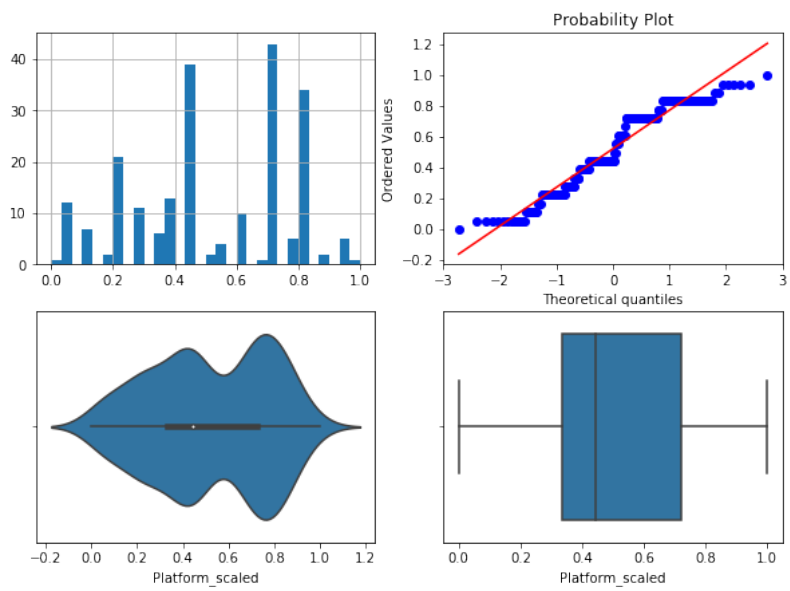Поле-Genre_scaled, метод-OutlierBoundaryType.SIGMA, строк-219

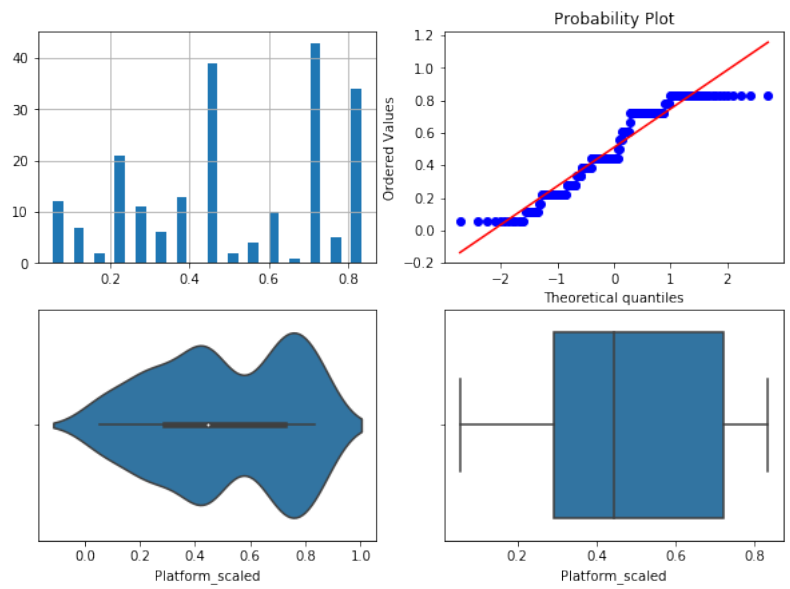Поле-Genre_scaled, метод-OutlierBoundaryType.QUANTILE, строк-208

Probability Plot



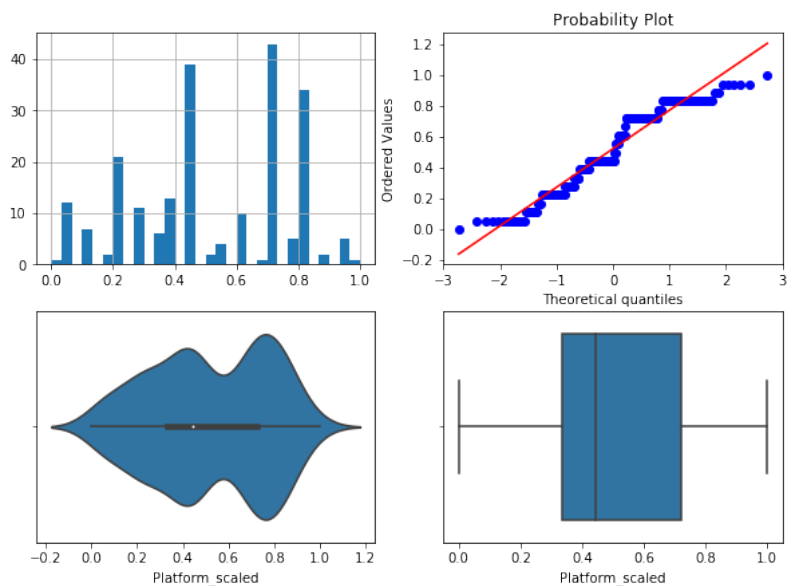Поле-Genre_scaled, метод-OutlierBoundaryType.IRQ, строк-219

Probability Plot

Поле-Platform_scaled, метод-OutlierBoundaryType.SIGMA, строк-219



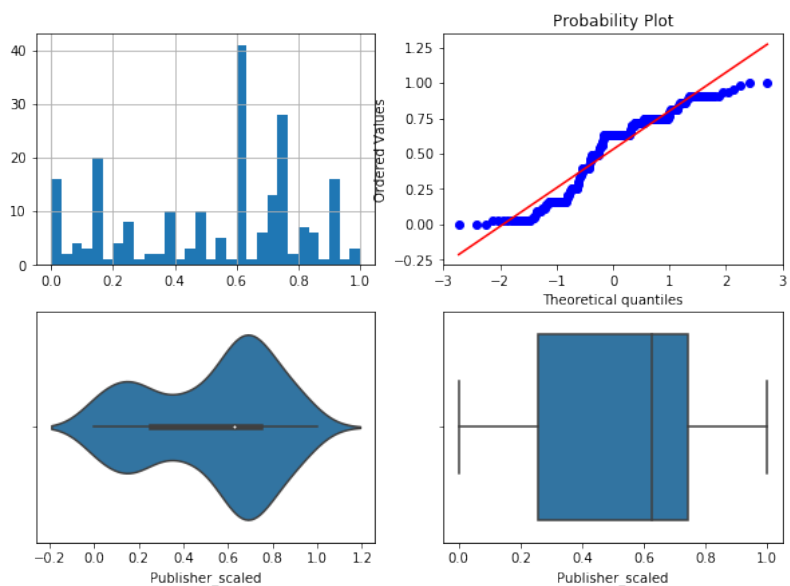Поле-Platform_scaled, метод-OutlierBoundaryType.QUANTILE, строк-210

Поле-Platform_scaled, метод-OutlierBoundaryType.IRQ, строк-219



Probability Plot

Поле-Publisher_scaled, метод-OutlierBoundaryType.SIGMA, строк-219



Probability Plot

Поле-Publisher_scaled, метод-OutlierBoundaryType.QUANTILE, строк-210



Поле-Publisher_scaled, метод-OutlierBoundaryType.IRQ, строк-219

Поле-Developer_scaled, метод-OutlierBoundaryType.SIGMA, строк-219


Поле-Developer_scaled, метод-OutlierBoundaryType.QUANTILE, строк-197

Поле-Developer_scaled, метод-OutlierBoundaryType.IRQ, строк-219

Поле-Year_scaled, метод-OutlierBoundaryType.SIGMA, строк-219

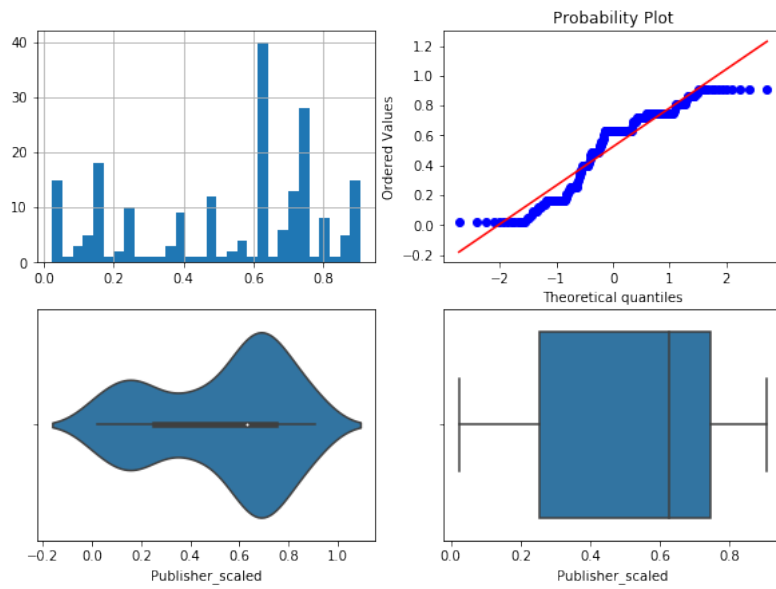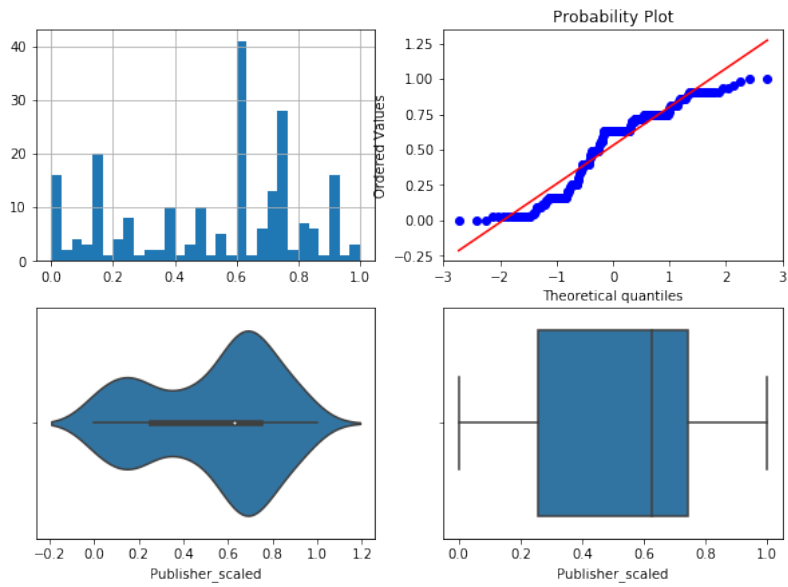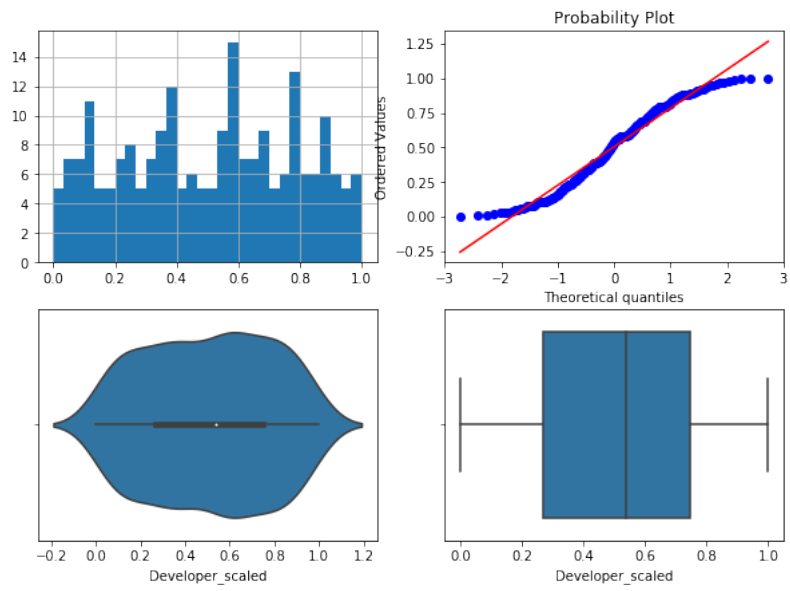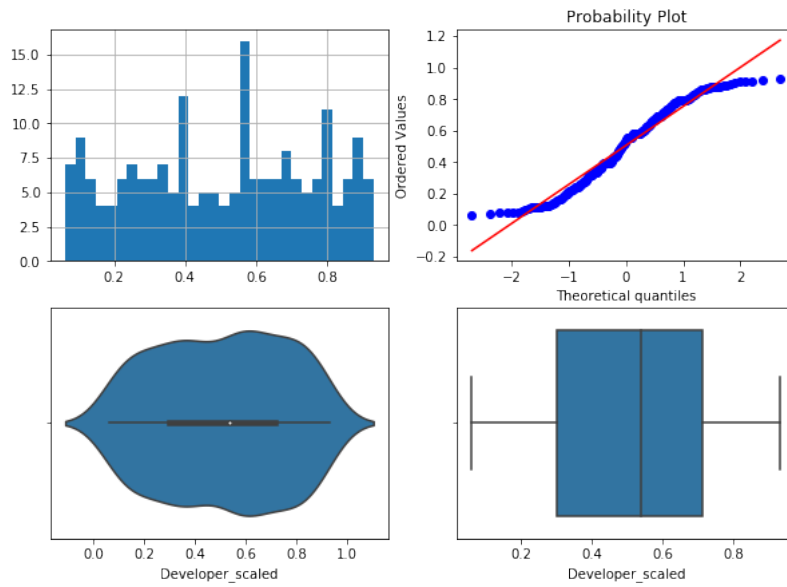Поле-Year_scaled, метод-OutlierBoundaryType.QUANTILE, строк-208



Probability Plot

Поле-Year_scaled, метод-OutlierBoundaryType.IRQ, строк-168



Probability Plot

Поле-Critic_Score_boxcox_scaled, метод-OutlierBoundaryType.SIGMA, строк-215

Поле-Critic_Score_boxcox_scaled, метод-OutlierBoundaryType.QUANTILE, строк-202

Поле-Critic_Score_boxcox_scaled, метод-OutlierBoundaryType.IRQ, строк-212



Поле-User_Score_boxcox_scaled, метод-OutlierBoundaryType.SIGMA, строк-215

Поле-User_Score_boxcox_scaled, метод-OutlierBoundaryType.QUANTILE, строк-197



Поле-User_Score_boxcox_scaled, метод-OutlierBoundaryType.IRQ, строк-206



```python
for col in scaled_columns:
    obt = OutlierBoundaryType.SIGMA
    # Вычисление верхней и нижней границы
    lower_boundary, upper_boundary = get_outlier_boundaries(df, col, obt)
    # Флаги для удаления выбросов
    outliers_temp = np.where(df[col] > upper_boundary, True,
                             np.where(df[col] < lower_boundary, True, False))
    # Удаление данных на основе флага
    df_trimmed = df.loc[~(outliers_temp), ]
    title = 'Поле-{}, метод-{}, строк-{}'.format(col, obt, df_trimmed.shape[0])
    diagnostic_plots(df_trimmed, col, title)
    df = df_trimmed
```

Поле-Sales_scaled, метод-OutlierBoundaryType.SIGMA, строк-213



Поле-ESRB_Rating_scaled, метод-OutlierBoundaryType.SIGMA, строк-213

Поле-Genre_scaled, метод-OutlierBoundaryType.SIGMA, строк-213



Поле-Platform_scaled, метод-OutlierBoundaryType.SIGMA, строк-213

Поле-Publisher_scaled, метод-OutlierBoundaryType.SIGMA, строк-213



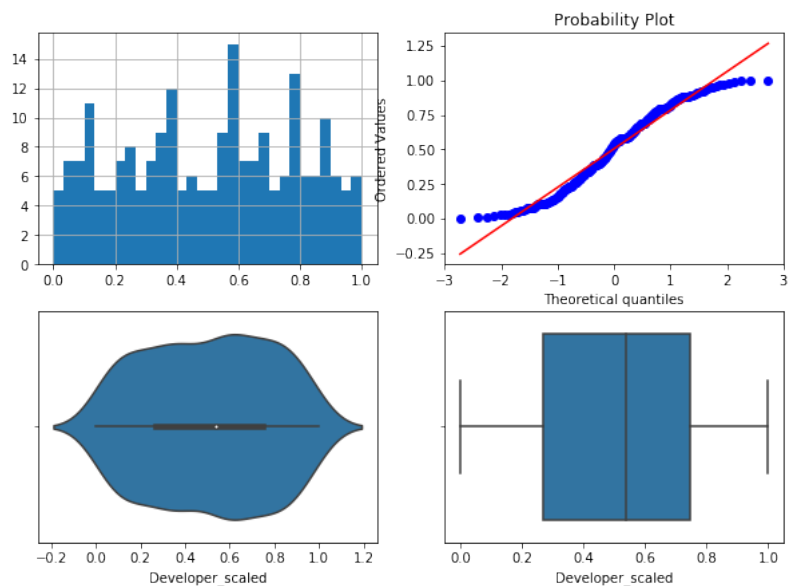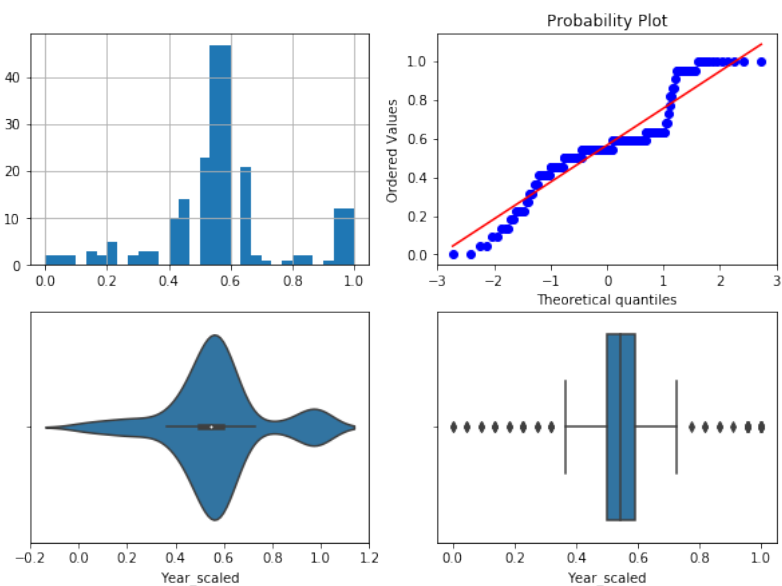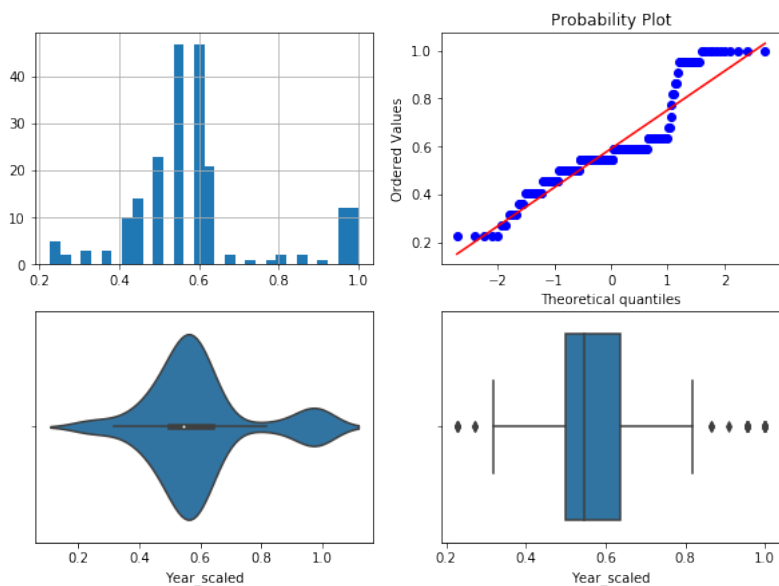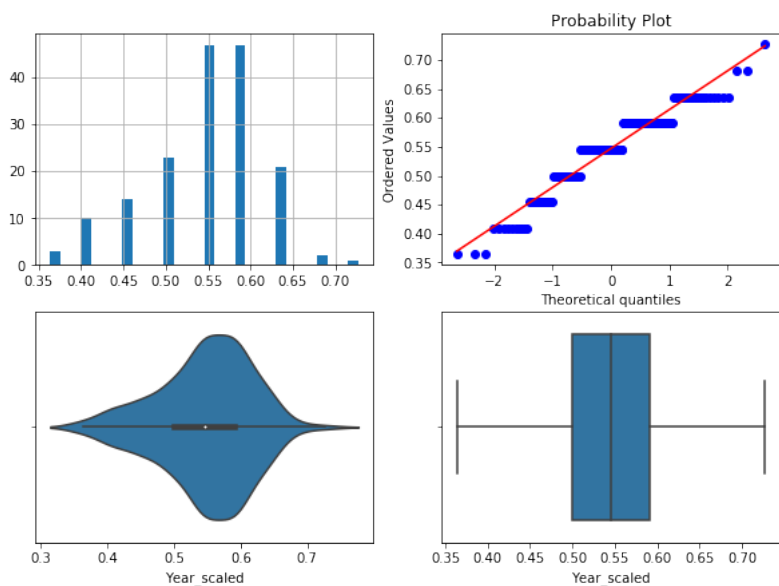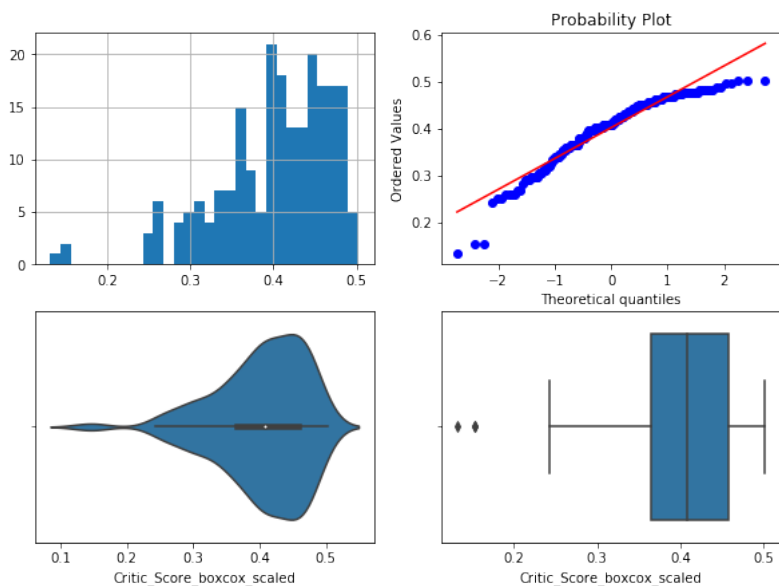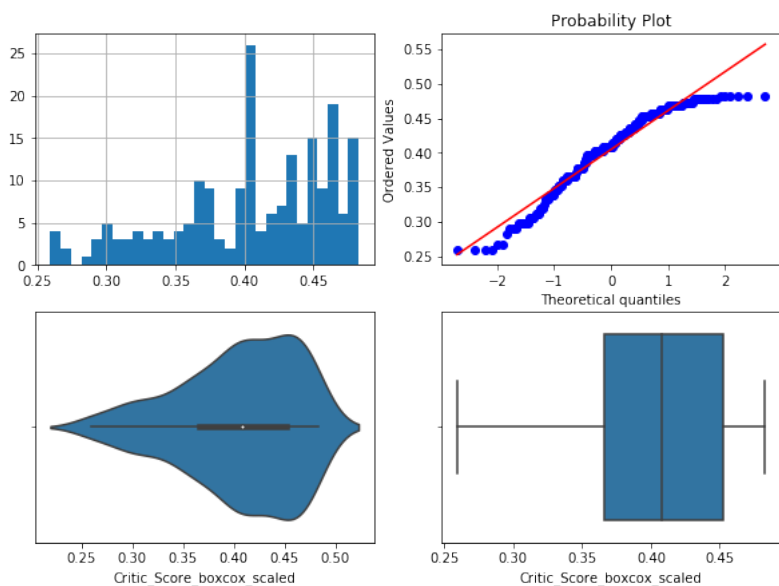Поле-Developer_scaled, метод-OutlierBoundaryType.SIGMA, строк-213
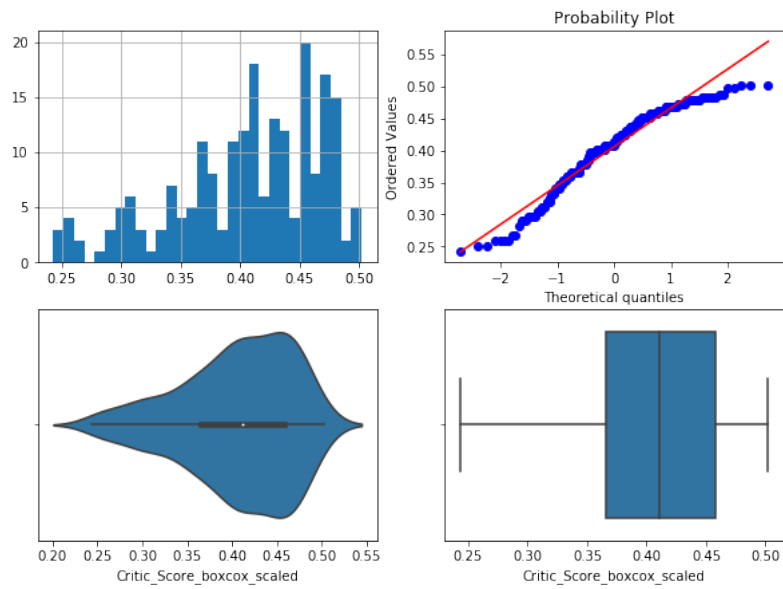
Поле-Year_scaled, метод-OutlierBoundaryType.SIGMA, строк-213



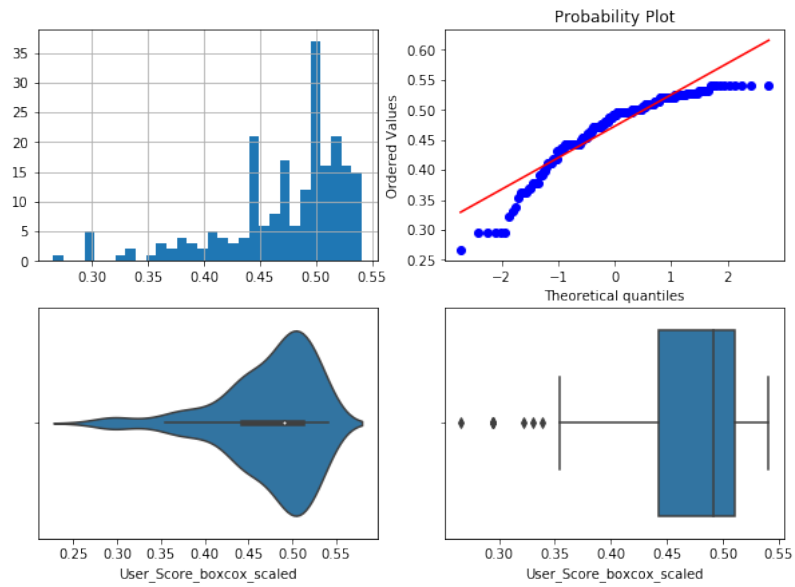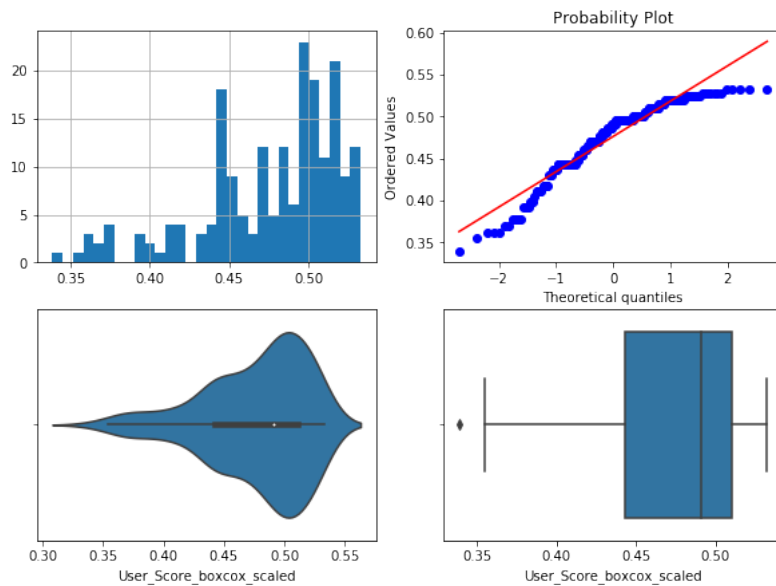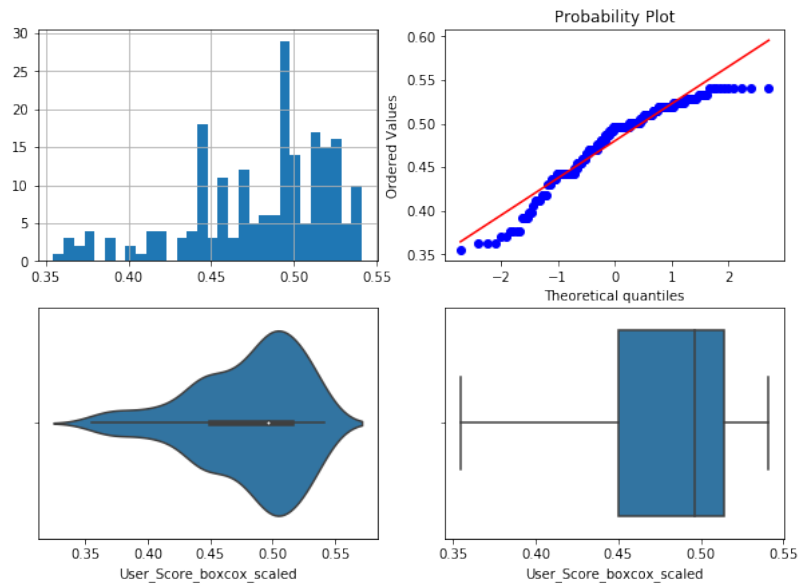Поле-Critic_Score_boxcox_scaled, метод-OutlierBoundaryType.SIGMA, строк-209

Поле-User_Score_boxcox_scaled, метод-OutlierBoundaryType.SIGMA, строк-207



```
1   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 | ... | User_Score_boxcox | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 216 | 54538 | Perception | 2 | 3 | 4 | 13 | 122 | 6.0 | 7.9 | 2017.0 | ... | 1.179739 | 0 |
| 217 | 55424 | Thumper | 6 | 4 | 4 | 9 | 32 | 9.0 | 9.3 | 2017.0 | ... | 1.224249 | 0 |

2 rows × 22 columns

## Отбор признаков

### Методы фильтрации (filter methods)

На основе оценки корреляции

```
1   plt.figure(figsize=(15,15))
2   sns.heatmap(df[scaled_columns].corr(), annot=True, fmt='.3f')
3   plt.show()
```

|  | Sales_scaled | ESRB_Rating_scaled | Genre_scaled | Platform_scaled | Publisher_scaled | Developer_scaled | Year_scaled | Critic_Score_boxcox_scaled | User_Score_boxcox_scaled |
|---|---|---|---|---|---|---|---|---|---|
| Sales_scaled | 1.000 | 0.010 | -0.046 | -0.014 | 0.028 | -0.034 | 0.012 | 0.477 | 0.294 |
| ESRB_Rating_scaled | 0.010 | 1.000 | -0.158 | 0.213 | -0.191 | 0.025 | -0.042 | 0.166 | 0.036 |
| Genre_scaled | -0.046 | -0.158 | 1.000 | 0.012 | -0.077 | -0.123 | -0.114 | -0.097 | 0.061 |
| Platform_scaled | -0.014 | 0.213 | 0.012 | 1.000 | -0.048 | -0.057 | -0.066 | -0.102 | -0.104 |
| Publisher_scaled | 0.028 | -0.191 | -0.077 | -0.048 | 1.000 | 0.467 | -0.035 | -0.029 | -0.033 |
| Developer_scaled | -0.034 | 0.025 | -0.123 | -0.057 | 0.467 | 1.000 | 0.011 | -0.062 | -0.021 |
| Year_scaled | 0.012 | -0.042 | -0.114 | -0.066 | -0.035 | 0.011 | 1.000 | -0.117 | -0.065 |
| Critic_Score_boxcox_scaled | 0.477 | 0.166 | -0.097 | -0.102 | -0.029 | -0.062 | -0.117 | 1.000 | 0.470 |
| User_Score_boxcox_scaled | 0.294 | 0.036 | 0.061 | -0.104 | -0.033 | -0.021 | -0.065 | 0.470 | 1.000 |

```python
def make_corr_df(df, tr=0.6):
    cr = df.corr()
    cr = cr.abs().unstack()
    cr = cr.sort_values(ascending=False)
    cr = cr[cr >= tr]
    cr = cr[cr < 1]
    cr = pd.DataFrame(cr).reset_index()
    cr.columns = ['f1', 'f2', 'corr']
    return cr
```

```python
make_corr_df(df)
```

```css
.dataframe tbody tr th {
    vertical-align: top;
}

.dataframe thead th {
    text-align: right;
}
```

| | f1 | f2 | corr |
|---|---|---|---|

```
1    # Обнаружение групп коррелирующих признаков
2    def corr_groups(cr):
3        grouped_feature_list = []
4        correlated_groups = []
5
6        for feature in cr['f1'].unique():
7            if feature not in grouped_feature_list:
8                # находим коррелирующие признаки
9                correlated_block = cr[cr['f1'] == feature]
10               cur_dups = list(correlated_block['f2'].unique()) + [feature]
11               grouped_feature_list = grouped_feature_list + cur_dups
12               correlated_groups.append(cur_dups)
13       return correlated_groups
```

```
1    # Группы коррелирующих признаков
2    drop_cols = []
3    for g in corr_groups(make_corr_df(df)):
4        for f in g:
5            if '_scaled' not in f:
6                drop_cols.append(f)
7
8    print(drop_cols)
```

```
1    []
```

## Методы обертывания (wrapper methods)

На основе алгоритма полного перебора

```
1    from sklearn.svm import SVR
2    from sklearn.svm import LinearSVC
3    from sklearn.feature_selection import SelectFromModel
4    from sklearn.linear_model import Lasso
5    from sklearn.linear_model import LinearRegression
6    from sklearn.linear_model import LogisticRegression
7    from sklearn.neighbors import KNeighborsClassifier
8    from sklearn.neighbors import KNeighborsRegressor
9    from sklearn.tree import DecisionTreeClassifier
10   from sklearn.ensemble import RandomForestClassifier
11   from sklearn.ensemble import GradientBoostingClassifier
12   from sklearn.tree import DecisionTreeRegressor
13   from sklearn.ensemble import RandomForestRegressor
14   from sklearn.ensemble import GradientBoostingRegressor
15   from sklearn.metrics import mean_squared_error
16   from sklearn.model_selection import train_test_split
17   from sklearn.feature_selection import VarianceThreshold
18   from sklearn.feature_selection import mutual_info_classif, mutual_info_regression
19   from sklearn.feature_selection import SelectKBest, SelectPercentile
```

```
1    from mlxtend.feature_selection import ExhaustiveFeatureSelector as EFS
2
3    lr = LinearRegression()
```

```
1    y_column = 'Sales'
2
3    trash_cols = {y_column, 'Sales_scaled', 'Developer', 'Name', 'Platform', 'Rank', 'Year', 'Publisher'}
4    x_columns = set(df.columns) - trash_cols
5
6    print(y_column, x_columns)
```

```
1    Sales {'Platform_scaled', 'Year_scaled', 'Developer_scaled', 'User_Score_boxcox_scaled', 'Genre_scaled', 'Publisher_scaled',
     'Critic_Score_boxcox_scaled', 'ESRB_Rating_scaled'}
```

```
1    def train_efs(df, x_cols, y_col, min_f=2, max_f=4, cv=5, model=None, scoring='neg_mean_squared_error'):
2        if model is None:
3            model = LinearRegression()
4
5        efs = EFS(model,
6                  min_features=min_f,
7                  max_features=max_f,
```

```
 8                  scoring=scoring,
 9                  print_progress=True,
10                  cv=cv)
11      efs = efs.fit(df[x_cols], pd.DataFrame(df[y_col]))
12
13      print('Best accuracy score: %.2f' % efs.best_score_)
14      print('Best subset (indices):', efs.best_idx_)
15      print('Best subset (corresponding names):', efs.best_feature_names_)
16      return efs1
```

```
 1  efs1 = train_efs(df, x_columns, y_column, 4, 8, 5, lr)
```

```
 1  Features: 163/163
 2
 3  Best accuracy score: -13.71
 4  Best subset (indices): (0, 2, 3, 6)
 5  Best subset (corresponding names): ('Platform_scaled', 'Developer_scaled', 'User_Score_boxcox_scaled', 'Critic_Score_boxcox_scaled')
```

```
 1  efs2 = train_efs(df, x_columns, y_column, 2, 8, 5, lr)
```

```
 1  Features: 247/247
 2
 3  Best accuracy score: -13.61
 4  Best subset (indices): (3, 6)
 5  Best subset (corresponding names): ('User_Score_boxcox_scaled', 'Critic_Score_boxcox_scaled')
```

## Методы вложений (embedded methods)

На основе линейной регрессии

```
 1  e_ls1 = LinearRegression()
 2  e_ls1.fit(df[x_columns], df[y_column])
 3  # Коэффициенты регрессии
 4  list(zip(x_columns, e_ls1.coef_))
```

```
 1  [('Platform_scaled', 0.8135289455413431),
 2   ('Year_scaled', 1.1786807196107814),
 3   ('Developer_scaled', -0.25420000858681396),
 4   ('User_Score_boxcox_scaled', 5.662058101682488),
 5   ('Genre_scaled', -0.10089104471518615),
 6   ('Publisher_scaled', 0.5227810641436722),
 7   ('Critic_Score_boxcox_scaled', 22.695110016228906),
 8   ('ESRB_Rating_scaled', -0.6470729147251617)]
```

```
 1  sel_e_ls1 = SelectFromModel(e_ls1)
 2  sel_e_ls1.fit(df[x_columns], df[y_column])
 3  list(zip(x_columns, sel_e_ls1.get_support()))
```

```
 1  [('Platform_scaled', False),
 2   ('Year_scaled', False),
 3   ('Developer_scaled', False),
 4   ('User_Score_boxcox_scaled', True),
 5   ('Genre_scaled', False),
 6   ('Publisher_scaled', False),
 7   ('Critic_Score_boxcox_scaled', True),
 8   ('ESRB_Rating_scaled', False)]
```

## Очистка фич

```
 1  df = df.drop(columns=drop_cols)
 2  df.head()
```

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

| | Rank | Name | Platform | Publisher | Developer | Year | Sales | Sales_scaled | ESRB_Rating_scaled | Genre_scaled | Platform_scaled | Pub |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 29 | Pokemon X/Y | 0 | 27 | 48 | 2013.0 | 16.37 | 0.440765 | 0.000000 | 0.714286 | 0.000000 | 0.62 |
| 7 | 34 | Pokemon Black / White Version | 1 | 27 | 48 | 2011.0 | 15.64 | 0.421109 | 0.000000 | 0.714286 | 0.055556 | 0.62 |
| 8 | 44 | Halo 3 | 15 | 21 | 15 | 2007.0 | 14.50 | 0.390415 | 1.000000 | 0.785714 | 0.833333 | 0.48 |
| 9 | 50 | Call of Duty: Modern Warfare 2 | 15 | 1 | 56 | 2009.0 | 13.53 | 0.364297 | 1.000000 | 0.785714 | 0.833333 | 0.02 |
| 10 | 53 | Super Smash Bros. Brawl | 13 | 27 | 91 | 2008.0 | 13.29 | 0.357835 | 0.666667 | 0.214286 | 0.722222 | 0.62 |

```
1    # df.drop(columns=['Platform', 'Publisher', 'Developer', 'Year', 'Sales_scaled'])
```

## Сохраняем

```
1    save_path = 'video_games_s3.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 | Platform | Publisher | Developer | Year | Sales | Sales_scaled | ESRB_Rating_scaled | Genre_scaled | Platform_scaled | Publ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 29 | Pokemon X/Y | 0 | 27 | 48 | 2013.0 | 16.37 | 0.440765 | 0.000000 | 0.714286 | 0.000000 | 0.627 |
| 1 | 34 | Pokemon Black / White Version | 1 | 27 | 48 | 2011.0 | 15.64 | 0.421109 | 0.000000 | 0.714286 | 0.055556 | 0.627 |
| 2 | 44 | Halo 3 | 15 | 21 | 15 | 2007.0 | 14.50 | 0.390415 | 1.000000 | 0.785714 | 0.833333 | 0.488 |
| 3 | 50 | Call of Duty: Modern Warfare 2 | 15 | 1 | 56 | 2009.0 | 13.53 | 0.364297 | 1.000000 | 0.785714 | 0.833333 | 0.023 |
| 4 | 53 | Super Smash Bros. Brawl | 13 | 27 | 91 | 2008.0 | 13.29 | 0.357835 | 0.666667 | 0.214286 | 0.722222 | 0.627 |

1