



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ _____ Информатика и системы управления

КАФЕДРА _____ Системы обработки информации и управления

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

К КУРСОВОЙ РАБОТЕ

НА ТЕМУ:

Анализ данных

Студент _____ ИУ5-63
(Группа)

_____ Кашеев М.С.
(Подпись, дата) (И.О.Фамилия)

Руководитель курсовой работы

_____ Гапанюк Ю.Е.
(Подпись, дата) (И.О.Фамилия)

Консультант

_____ (Подпись, дата) (И.О.Фамилия)

2022 г.

Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

УТВЕРЖДАЮ
Заведующий кафедрой _____
(Индекс)

(И.О.Фамилия)
« ____ » _____ 20 ____ г.

ЗАДАНИЕ

на выполнение научно-исследовательской работы

по теме Анализ данных

_____ Студент группы ИУ5-63

Кашеев Максим Сергеевич

(Фамилия, имя, отчество)

Направленность НИР (учебная, исследовательская, практическая, производственная, др.)
учебная

Источник тематики (кафедра, предприятие, НИР) _____ кафедра

График выполнения работы: 25% к ____ нед., 50% к ____ нед., 75% к ____ нед., 100% к ____ нед.

Задание Поиск и выбор набора данных для построения моделей машинного обучения. На основе выбранного набора данных построить модели машинного обучения для решения задачи классификации

Оформление научно-исследовательской работы:

Расчетно-пояснительная записка на 27 листах формата А4.

Дата выдачи задания « ____ » _____ 20 ____ г.

Руководитель курсовой работы

(Подпись, дата) Гапанюк Ю.Е.
(И.О.Фамилия)

Студент

(Подпись, дата) Кашеев М.С.
(И.О.Фамилия)

Примечание: Задание оформляется в двух экземплярах: один выдается студенту, второй хранится на кафедре.

Задание

1	1. Поиск и выбор набора данных для построения моделей машинного обучения. На основе выбранного набора данных студент должен построить модели машинного обучения для решения или задачи классификации, или задачи регрессии.
2	2. Проведение разведочного анализа данных. Построение графиков, необходимых для понимания структуры данных. Анализ и заполнение пропусков в данных.
3	3. Выбор признаков, подходящих для построения моделей. Кодирование категориальных признаков. Масштабирование данных. Формирование вспомогательных признаков, улучшающих качество моделей.
4	4. Проведение корреляционного анализа данных. Формирование промежуточных выводов о возможности построения моделей машинного обучения. В зависимости от набора данных, порядок выполнения пунктов 2, 3, 4 может быть изменен.
5	5. Выбор метрик для последующей оценки качества моделей. Необходимо выбрать не менее трех метрик и обосновать выбор.
6	6. Выбор наиболее подходящих моделей для решения задачи классификации или регрессии. Необходимо использовать не менее пяти моделей, две из которых должны быть ансамблевыми.
7	7. Формирование обучающей и тестовой выборок на основе исходного набора данных.
8	8. Построение базового решения (baseline) для выбранных моделей без подбора гиперпараметров. Производится обучение моделей на основе обучающей выборки и оценка качества моделей на основе тестовой выборки.
9	9. Подбор гиперпараметров для выбранных моделей. Рекомендуется использовать методы кросс-валидации. В зависимости от используемой библиотеки можно применять функцию GridSearchCV, использовать перебор параметров в цикле, или использовать другие методы.
10	10. Повторение пункта 8 для найденных оптимальных значений гиперпараметров. Сравнение качества полученных моделей с качеством baseline-моделей.
11	11. Формирование выводов о качестве построенных моделей на основе выбранных метрик. Результаты сравнения качества рекомендуется отобразить в виде графиков и сделать выводы в форме текстового описания. Рекомендуется построение графиков обучения и валидации, влияния значений гиперпараметров на качество моделей и т.д.

Выполнение работы

В качестве набора данных был выбран набор, который помогает определить контингент, демографические признаки и востребованность рекламной кампании.

Датасет состоит из одного файла. В последствии он будет разбит на обучающий и тестовый датасеты.

Файл содержит включает в себя несколько колонок:

Daily time spent on site - ежедневное время проведенное на сайте

Age – возраст

Area income – доход в выбранной области

Daily internet usage – использование интернета

Ad topic line – содержание рекламы

City – город

Male – пол

Country – страна

Timestap – временная метка

Click on Ad – переход на сайт рекламы

Импорт библиотек

In [1]:

```
1         import numpy as np
2         import pandas as pd
3         import seaborn as sns
4         import matplotlib.pyplot as plt
5         import os
6         from sklearn.preprocessing import MinMaxScaler
7         from sklearn.linear_model import LinearRegression, LogisticRegression
8         from sklearn.model_selection import train_test_split
9         from sklearn.neighbors import KNeighborsRegressor, KNeighborsClassifier
10        from sklearn.metrics import accuracy_score, balanced_accuracy_score
11        from sklearn.metrics import precision_score, recall_score, f1_score, classification_re
12        from sklearn.metrics import confusion_matrix
13        from sklearn.metrics import plot_confusion_matrix
14        from sklearn.model_selection import GridSearchCV
15        from sklearn.metrics import mean_absolute_error, mean_squared_error, mean_squared_log_
16        from sklearn.metrics import roc_curve, roc_auc_score
17        from sklearn.svm import SVC, NuSVC, LinearSVC, OneClassSVM, SVR, NuSVR, LinearSVR
18        from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor, export_graphviz
19        from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
20        from sklearn.ensemble import ExtraTreesClassifier, ExtraTreesRegressor
21        from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor
22        from gmdhpy import gmdh
23        %matplotlib inline
24        sns.set(style="ticks")
```

```

1 # Отрисовка ROC-кривой
2 def draw_roc_curve(y_true, y_score, pos_label=1, average='micro'): fpr, tpr, thresholds =
3     roc_curve(y_true, y_score,
4               pos_label=pos_label)
5     roc_auc_value = roc_auc_score(y_true, y_score, average=average) plt.figure()
6     lw = 2
7     plt.plot(fpr, tpr, color='darkorange',
8             lw=lw, label='ROC curve (area = %0.2f)' % roc_auc_value) plt.plot([0, 1],
9     [0, 1], color='navy', lw=lw, linestyle='--')
10    plt.xlim([0.0, 1.0])
11    plt.ylim([0.0, 1.05])
12    plt.xlabel('False Positive Rate')
13    plt.ylabel('True Positive Rate')
14    plt.title('Receiver operating characteristic')
15    plt.legend(loc="lower right")
16    plt.show()
17

```

Разбиение файла на обучающую и тестовую выборку

Файл разделяется на обучающую и тестовую выборку в соотношении 70:30
соответственно

In [3]:

```

1 def split(filehandler, delimiter=',', row_limit=700,
2           output_name_template='advertising%s.csv', output_path='.', keep_headers=True
3           import csv
4           reader = csv.reader(filehandler, delimiter=delimiter)
5           current_piece = 1
6           current_out_path = os.path.join(
7               output_path,
8               output_name_template % current_piece9
9           )
10          current_out_writer = csv.writer(open(current_out_path, 'w'), delimiter=delimiter)
11          current_limit = row_limit
12          if keep_headers:
13              headers = next(reader)
14              current_out_writer.writerow(headers)
15          for i, row in enumerate(reader):
16              if i + 1 > current_limit:
17                  current_piece += 1
18                  current_limit = row_limit * current_piece
19                  current_out_path = os.path.join(
20                      output_path,
21                      output_name_template % current_piece22
22                  )
23                  current_out_writer = csv.writer(open(current_out_path, 'w'), delimiter=del
24                  if keep_headers:
25                      current_out_writer.writerow(headers)
26                  current_out_writer.writerow(row)27

```

```
1 split(open('advertising.csv', 'r'));
```

Переименуем полученные файлы

In [5]:

```
1 os.rename('advertising1.csv', 'advertising_Train.csv')
2 os.rename('advertising2.csv', 'advertising_Test.csv')
```

Создание обучающей выборки

In [6]:

```
1 # Обучающая выборка
2 train = pd.read_csv('advertising_Train.csv', sep=",")
```

Выведем основную информацию об обучающей выборке

In [7]:

```
1 # Первые 5 строк обучающей выборки
2 train.head()
```

Out[7]:

	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp	Click on
0	68.95	35	61833.90	256.09	Cloned 5thgeneration orchestration	Wrightburgh	0	Tunisia	2016-03-27 00:53:11	
1	80.23	31	68441.85	193.77	Monitored national standardization	West Jodi	1	Nauru	2016-04-04 01:39:02	
2	69.47	26	59785.94	236.50	Organic bottom-line service-desk	Davidton	0	San Marino	2016-03-13 20:35:42	
3	74.15	29	54806.18	245.89	Triple-buffered reciprocal time- frame	West Terrifurt	1	Italy	2016-01-10 02:31:19	
4	68.37	35	73889.99	225.58	Robust logistical utilization	South Manuel	0	Iceland	2016-06-03 03:36:18	

1	<i># Размер обучающего датасета - 700 строк, 10 колонок</i>
2	<code>train.shape</code>

Out[9]:

(700, 10)

In [10]:

1	<i># Набор колонок</i>
2	<code>train.columns</code>

Out[10]:

```
Index(['Daily Time Spent on Site', 'Age', 'Area Income',  
      'Daily Internet Usage', 'Ad Topic Line', 'City', 'Male', 'Country',  
      'Timestamp', 'Clicked on Ad'],  
      dtype='object')
```

In [11]:

1	<i># Набор колонок с типами данных</i>
2	<code>train.dtypes</code>

Out[11]:

Daily Time Spent on Site	float64
Age	int64
Area Income	float64
Daily Internet Usage	float64
Ad Topic Line	object
City	object
Male	int64
Country	object
Timestamp	object
Clicked on Ad	int64
dtype:	object

Типы, отличные от float64 и int64, в последствии не будут использоваться. Однако посчитала нужным их оставить для полноты картины

```
1 # Подробная информация о данных
2 train.info()
```

```
<class
'pandas.core.frame.DataFrame'>
RangeIndex: 700 entries, 0 to 699
Data columns (total 10 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   Daily Time Spent on Site    700 non-null   float64
1   Age                        700 non-null   int64
2   Area Income                700 non-null   float64
3   Daily Internet Usage       700 non-null   float64
4   Ad Topic Line              700 non-null   object
5   City                      700 non-null   object
6   Male                      700 non-null   int64
7   Country                   700 non-null   object
8   Timestamp                 700 non-null   object
9   Clicked on Ad              700 non-null   int64
dtypes: float64(3), int64(3), object(4)
memory usage: 54.8+ KB
```

In [13]:

```
1 # Проверим наличие пустых значений
2 train.isnull().sum()
```

Out[13]:

```
Daily Time Spent on Site    0
Age                        0
Area Income                 0
Daily Internet Usage        0
Ad Topic Line               0
City                       0
Male                       0
Country                    0
Timestamp                  0
Clicked on Ad               0
dtype: int64
```

Создание тестовой выборки

In [14]:

```
1 # Тестовая выборка
2 test = pd.read_csv('advertising_Test.csv', sep=',')
```


Выведем тестовую информацию о созданной выборке

1	# Первые 5 строк тестовой выборки
2	

test.head()

Out[15]:

	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp	Cli o
0	58.60	19	44490.09	197.93	Configurable impactful firmware	West Mariafort	1	Lebanon	2016-06-05 21:38:22	
1	69.77	54	57667.99	132.27	Face-to-face dedicated flexibility	Port Sherrystad	0	Malta	2016-06-01 03:17:50	
2	87.27	30	51824.01	204.27	Fully- configurable 5thgeneration circuit	West Melissashire	1	Christmas Island	2016-03-06 06:51:23	
3	77.65	28	66198.66	208.01	Configurable impactful capacity	Pamelamouth	0	Ukraine	2016-02-26 19:35:54	
4	76.02	40	73174.19	219.55	Distributed leadingedge orchestration	Lesliefort	0	Malta	2016-07-13 14:30:14	

In [17]:

1	# Размер тестового датасета - 300 строк, 10 столбцов
2	

test.shape

Out[17]:

(300, 10)

In [18]:

```
1                                     #Подробная информация о тестовых данных
2                                     test.info()
```

<class

'pandas.core.frame.DataFrame'>

RangeIndex: 300 entries, 0 to 299

Data columns (total 10 columns):

#	Column	Non-Null Count		Dtype
0	Daily Time Spent on Site	300	non-null	float64
1	Age	300	non-null	int64
2	Area Income	300	non-null	float64
3	Daily Internet Usage	300	non-null	float64
4	Ad Topic Line	300	non-null	object
5	City	300	non-null	object
6	Male	300	non-null	int64
7	Country	300	non-null	object
8	Timestamp	300	non-null	object
9	Clicked on Ad	300	non-null	int64

dtypes: float64(3), int64(3), object(4)

memory usage: 23.6+ KB

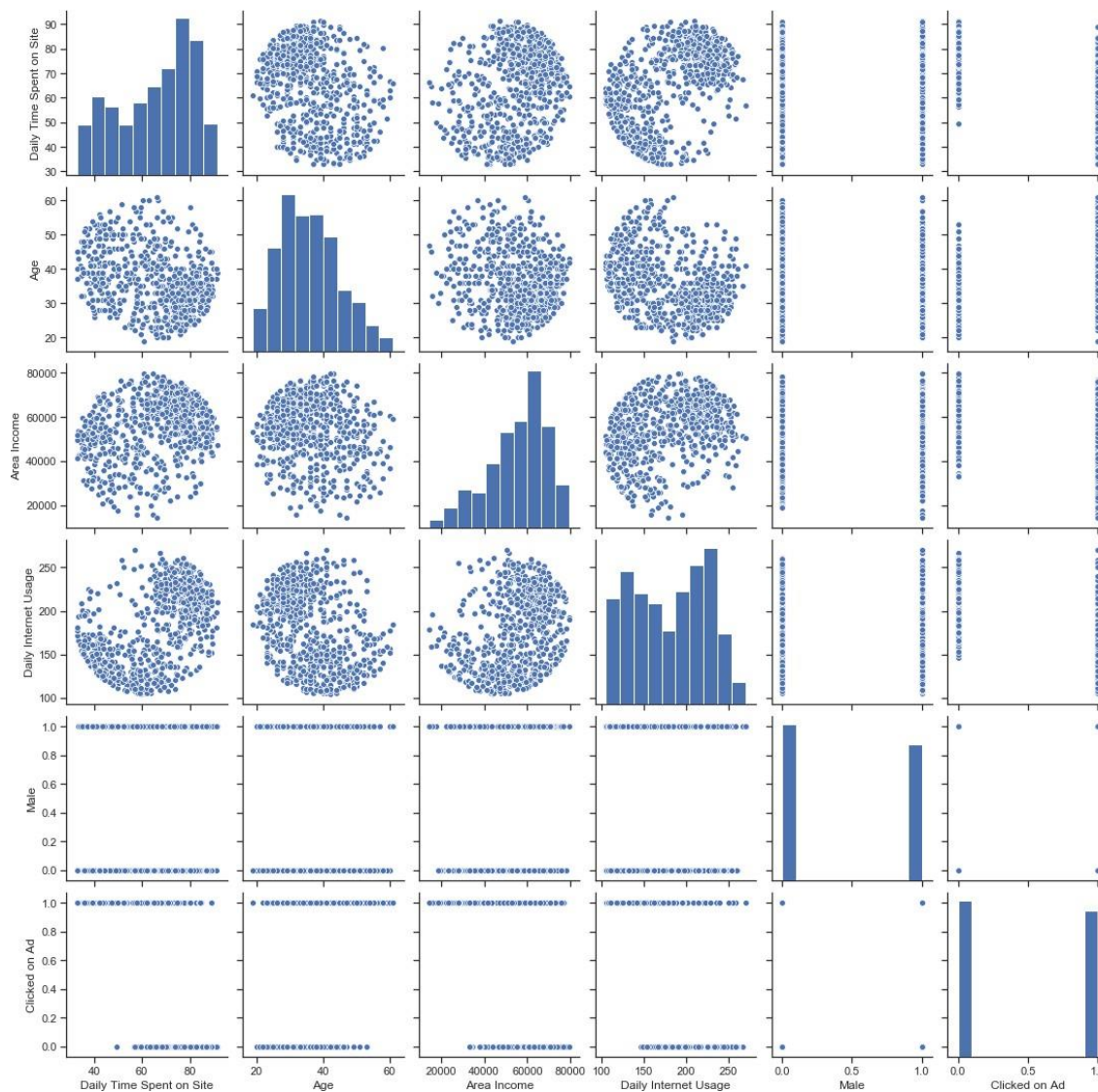
Построим парные диаграммы для обучающей выборки

In [19]:

```
1 # Парные диаграммы
2 sns.pairplot(train)
```

Out[19]:

<seaborn.axisgrid.PairGrid at 0x2132a91aa60>



Построим парные диаграммы, используя целевой признак классификации

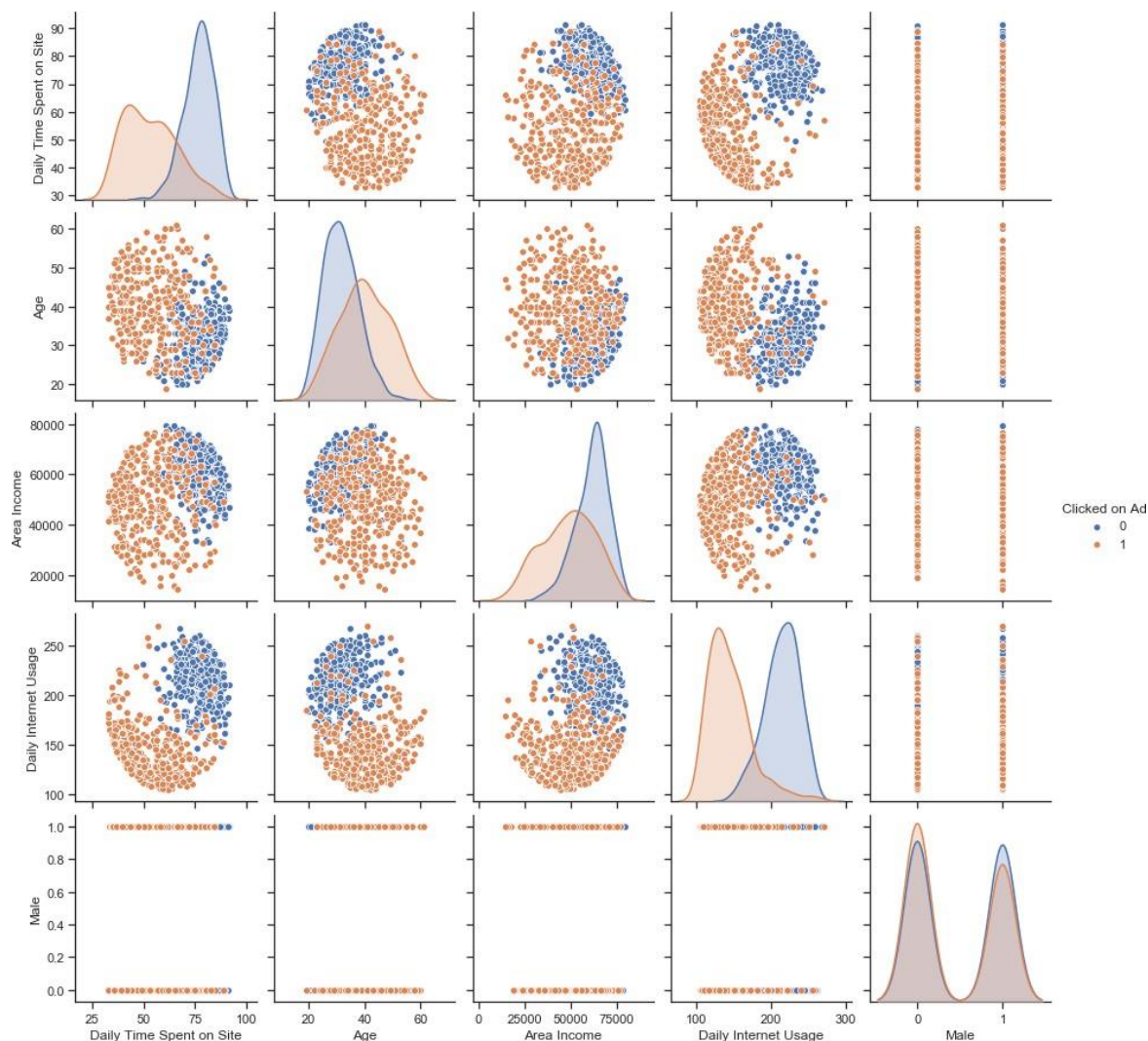
За целевой признак для классификации возьмем "Clicked on Ad" - переход на сайт рекламодателя

In [20]:

```
1 sns.pairplot(train, hue="Clicked on Ad")
```

Out[20]:

<seaborn.axisgrid.PairGrid at 0x2132bbd16a0>



Возможные значения целевого признака классификации в обучающей выборке

In [21]:

```
1 np.sort(train['Clicked on Ad'].unique())
```

Out[21]:

```
array([0, 1], dtype=int64)
```

Возможные значения целевого признака классификации в тестовой выборке

```
1 np.sort(test['Clicked on Ad'].unique())
```

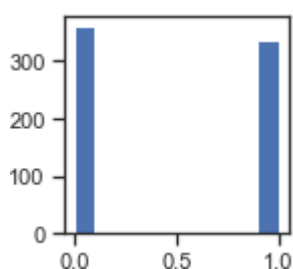
Out[22]:

```
array([0, 1], dtype=int64)
```

Оценка дисбаланса классов для обучающей выборки

In [23]:

```
1 # Оценим дисбаланс классов для Clicked on Ad
2 fig, ax = plt.subplots(figsize=(2,2))
3 plt.hist(train['Clicked on Ad'])
4 plt.show()
```



In [24]:

```
1 # Посмотрим, как часто встречается тот или иной целевой признак
2 train['Clicked on Ad'].value_counts()
```

Out[24]:

```
0    362
```

```
1    338
```

Name: Clicked on Ad, dtype: int64

In [25]:

```
1 # Посчитаем дисбаланс классов
2 total = train.shape[0]
3 class_0, class_1 = train['Clicked on Ad'].value_counts() print('Класс 0
4 составляет { }%, а класс 1 составляет { }%.'
5 .format(round(class_0 / total, 4)*100, round(class_1 / total, 4)*100))
```

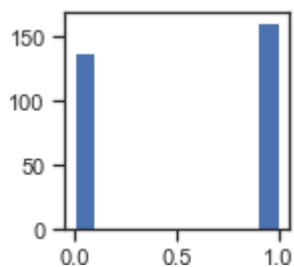
Класс 0 составляет 51.71%, а класс 1 составляет 48.29%.

Оценка дисбаланса классов для тестовой выборки

```

1 # Оценим дисбаланс классов для Clicked on Ad
2 fig, ax = plt.subplots(figsize=(2,2))
3 plt.hist(test['Clicked on Ad'])
4 plt.show()

```



In [27]:

```

1 # Посмотрим, как часто встречается тот или иной целевой признак
2 test['Clicked on Ad'].value_counts()

```

Out[27]:

```

1    162
0    138
Name: Clicked on Ad, dtype: int64

```

In [28]:

```

1 # Посчитаем дисбаланс классов
2 total = test.shape[0]
3 class_0, class_1 = test['Clicked on Ad'].value_counts()
4 print('Класс 0 составляет {}%, а класс 1 составляет {}%.'
5       .format(round(class_0 / total, 4)*100, round(class_1 / total, 4)*100))

```

Класс 0 составляет 54.0%, а класс 1 составляет 46.0%.

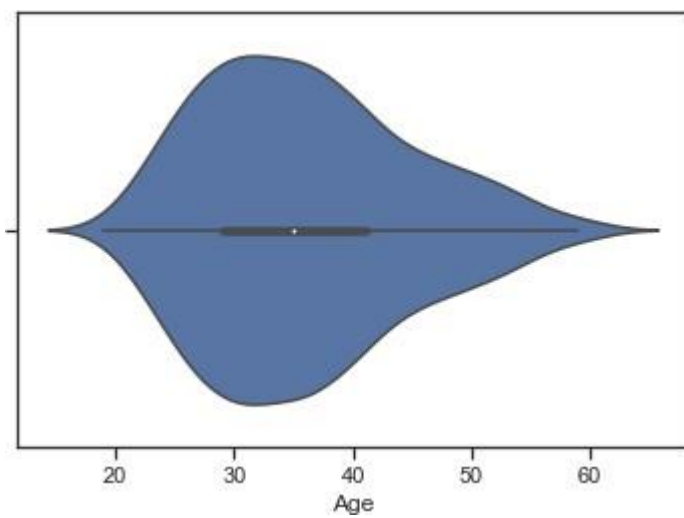
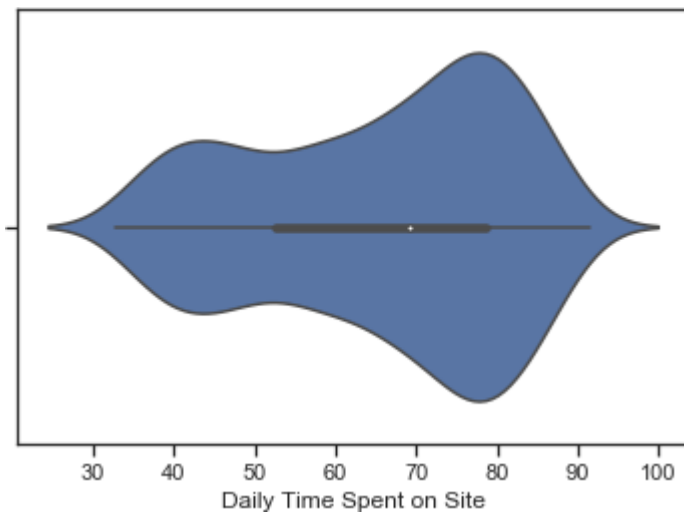
Вывод об оценке дисбаланса классов

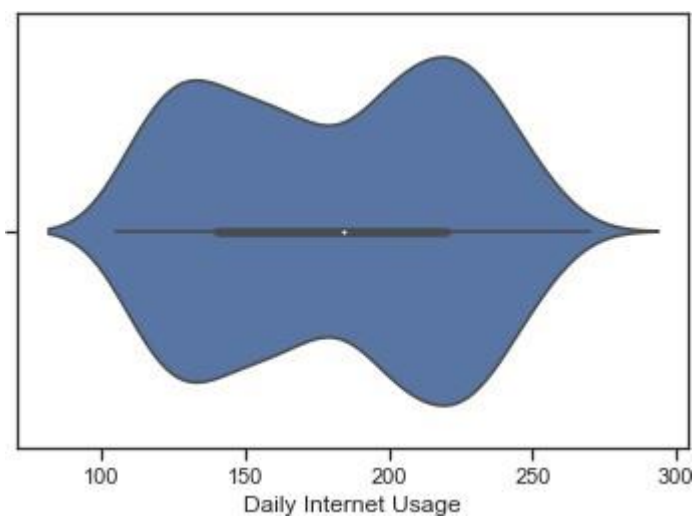
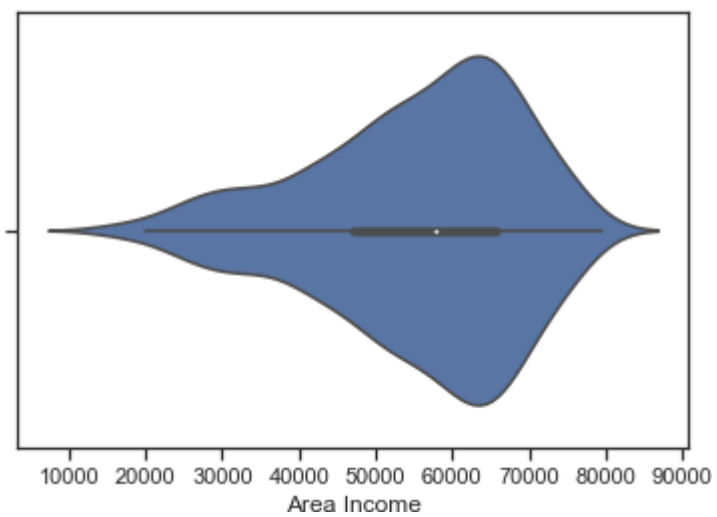
Дисбаланс классов практически отсутствует

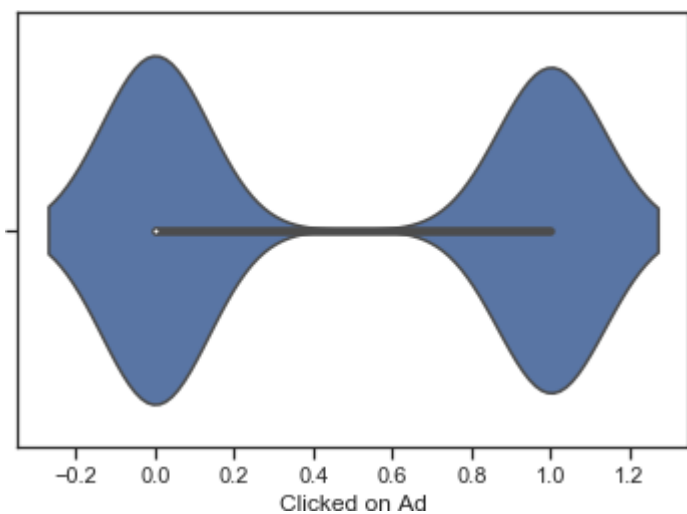
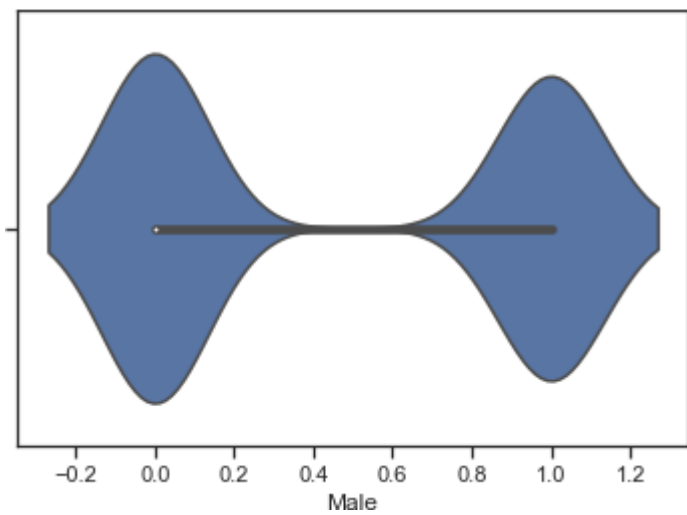
Построение скрипичных диаграмм для обучающей выборки

In [32]:

```
1 # Скрипичные диаграммы для числовых колонок
2 for col in ['Daily Time Spent on Site', 'Age', 'Area Income', 'Daily Internet Usage',sns.violinplot(x=train[col])
3     plt.show()
4
```







In [33]:

```
1 train.dtypes
```

Out[33]:

```
Daily Time Spent on Site    float64
Age                        int64
Area Income                 float64
Daily Internet Usage       float64
Ad Topic Line              object
City                       object
Male                       int64
Country                    object
Timestamp                  object
Clicked on Ad              int64
dtype: object
```

Убираем значения, тип которых не равен float64 и int64

In [34]:

```
1 # Создадим вспомогательные колонки, чтобы наборы данных можно было разделить
2 train['dataset'] = 'TRAIN'
3 test['dataset'] = 'TEST'
```

In [35]:

```
1 # Колонки для объединения
2 join_cols = ['Daily Time Spent on Site', 'Age', 'Area Income', 'Daily Internet Usage',
```

Склеиваем обучающую и тестовую выборку

In [36]:

```
1 data_all = pd.concat([train[join_cols], test[join_cols]])
```

In [38]:

```
1 # Проверим корректность объединения
2 assert data_all.shape[0] == train.shape[0]+test.shape[0]
```

In [39]:

```
1 data_all.head()
```

Out[39]:

	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Male	Clicked on Ad	dataset
0	68.95	35	61833.90	256.09	0	0	TRAIN
1	80.23	31	68441.85	193.77	1	0	TRAIN
2	69.47	26	59785.94	236.50	0	0	TRAIN
3	74.15	29	54806.18	245.89	1	0	TRAIN
4	68.37	35	73889.99	225.58	0	0	TRAIN

Отмасштабируем все признаки кроме целевого

In [46]:

```
1 # Числовые колонки для масштабирования
2 scale_cols = ['Daily Time Spent on Site', 'Age', 'Area Income', 'Daily Internet Usage']
```

In [47]:

```
1 sc1 = MinMaxScaler()
2 sc1_data = sc1.fit_transform(data_all[scale_cols])
```

In [48]:

```
1 for i in range(len(scale_cols)):
2     col = scale_cols[i]
3     new_col_name = col + '_scaled'
4     data_all[new_col_name] = sc1_data[:,i]
```

In [49]:

```
1 data_all.head()
```

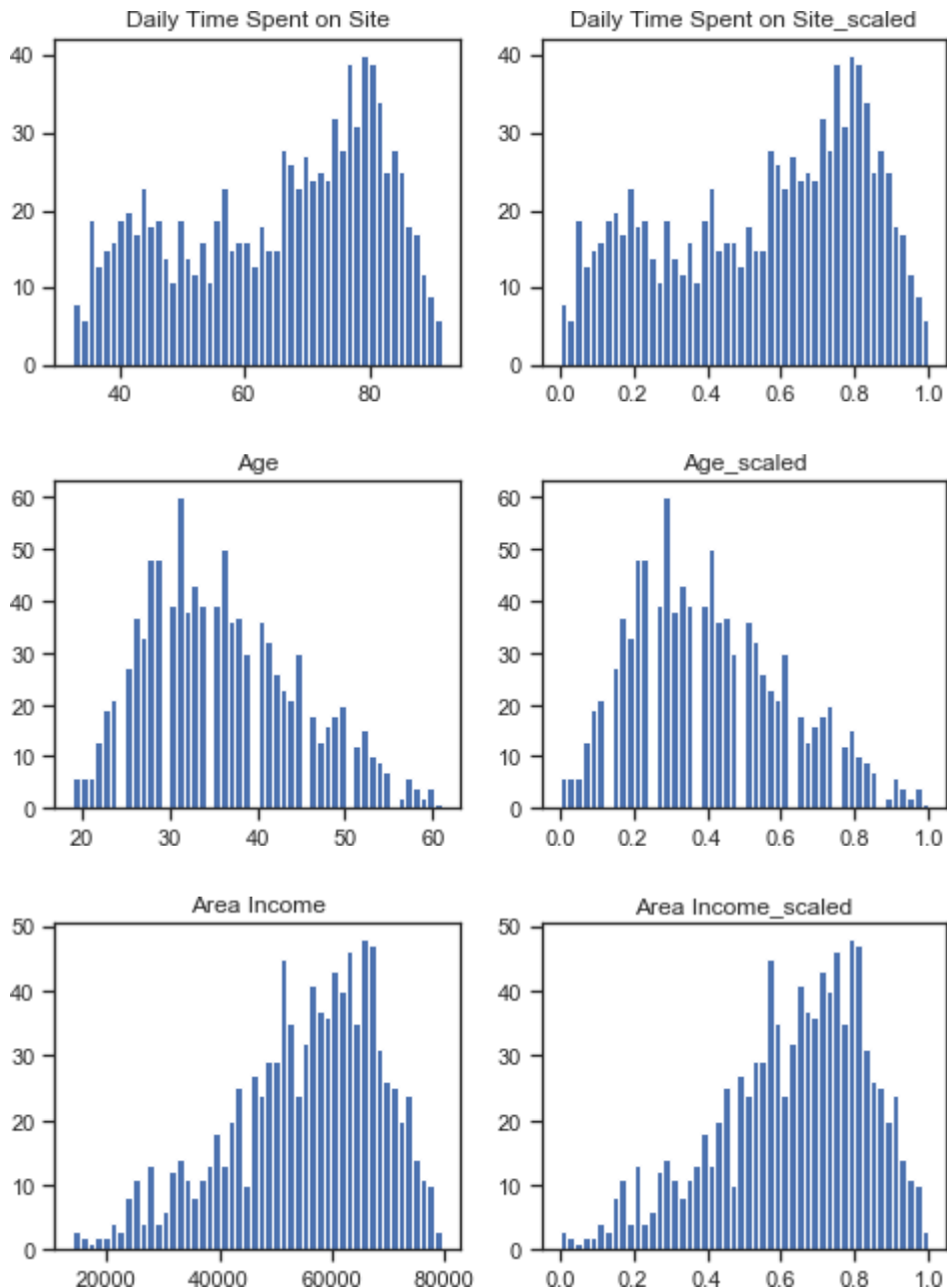
Out[49]:

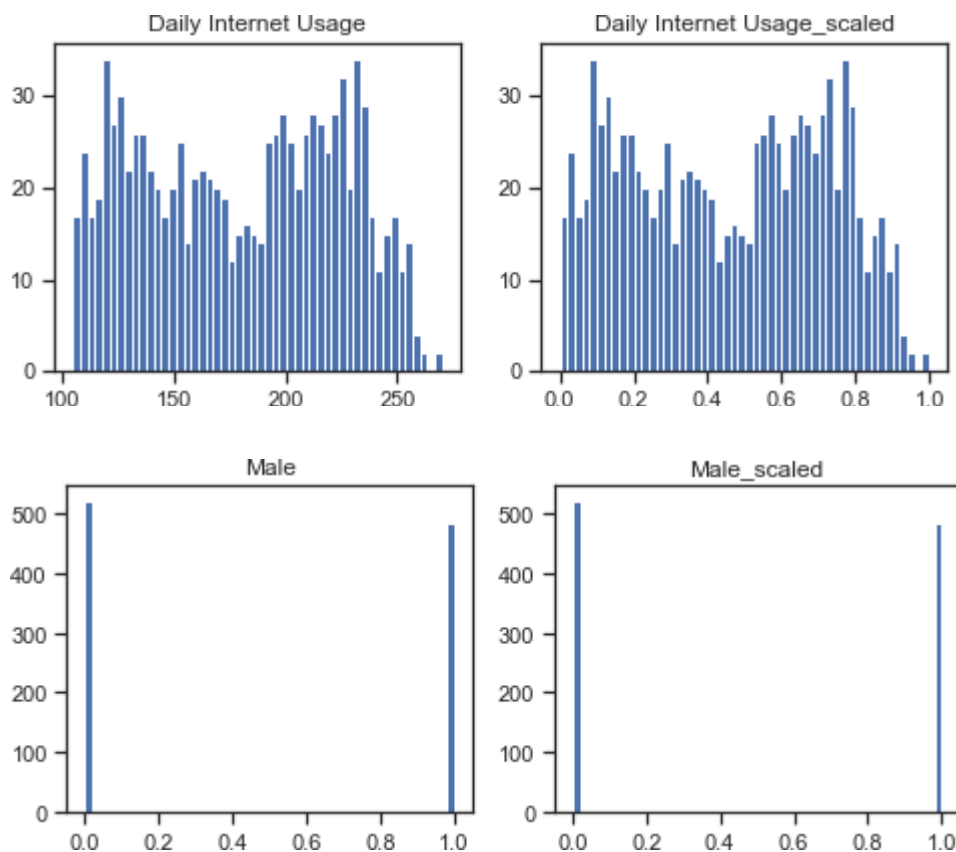
	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Male	Clicked on Ad	dataset	Daily Time Spent on Site_scaled	Age_scaled	Income_sc
0	68.95	35	61833.90	256.09	0	0	TRAIN	0.617882	0.380952	0.730
1	80.23	31	68441.85	193.77	1	0	TRAIN	0.809621	0.285714	0.831
2	69.47	26	59785.94	236.50	0	0	TRAIN	0.626721	0.166667	0.699
3	74.15	29	54806.18	245.89	1	0	TRAIN	0.706272	0.238095	0.623
4	68.37	35	73889.99	225.58	0	0	TRAIN	0.608023	0.380952	0.914

Проверим, что масштабирование не повлияло на распределение данных

In [50]:

```
1      # Проверим, что масштабирование не повлияло на распределение данных
2      for col in scale_cols:
3          col_scaled = col + '_scaled'
4
5      fig, ax = plt.subplots(1, 2, figsize=(8,3))
6      ax[0].hist(data_all[col], 50)
7      ax[1].hist(data_all[col_scaled], 50)
8      ax[0].title.set_text(col)
9      ax[1].title.set_text(col_scaled)
10     plt.show()
```





Вернем в набор данных целевой признак

In [51]:

```
1 corr_cols_1 = scale_cols + ['Clicked on Ad']
2 corr_cols_1
```

Out[51]:

```
['Daily Time Spent on Site',
 'Age',
 'Area Income',
 'Daily Internet Usage',
 'Male',
 'Clicked on Ad']
```

In [52]:

```
1 scale_cols_postfix = [x+'_scaled' for x in scale_cols]
2 corr_cols_2 = scale_cols_postfix + ['Clicked on Ad']
3 corr_cols_2
```

Out[52]:

```
['Daily Time Spent on Site_scaled',
 'Age_scaled',
 'Area Income_scaled',
 'Daily Internet Usage_scaled',
 'Male_scaled',
 'Clicked on Ad']
```

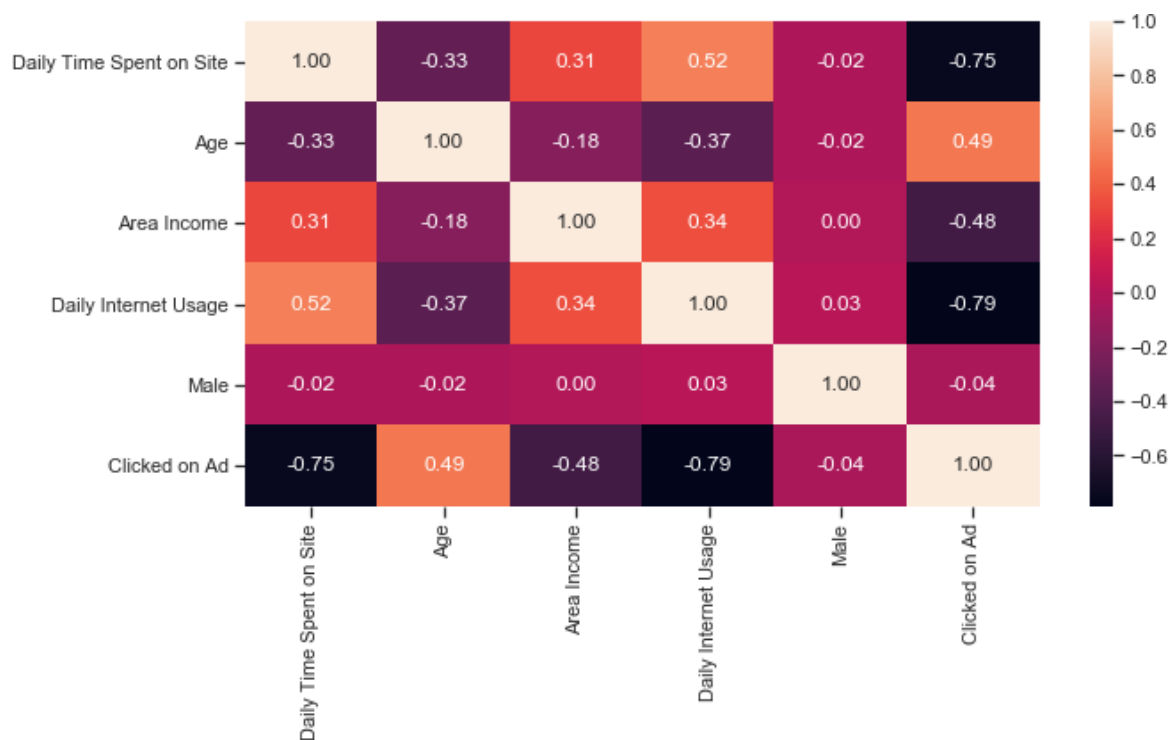
Построим корреляционную матрицу для первичных данных

In [53]:

```
1 fig, ax = plt.subplots(figsize=(10,5))
2 sns.heatmap(data_all[corr_cols_1].corr(), annot=True, fmt='.2f')
```

Out[53]:

<matplotlib.axes._subplots.AxesSubplot at 0x2132ea26910>



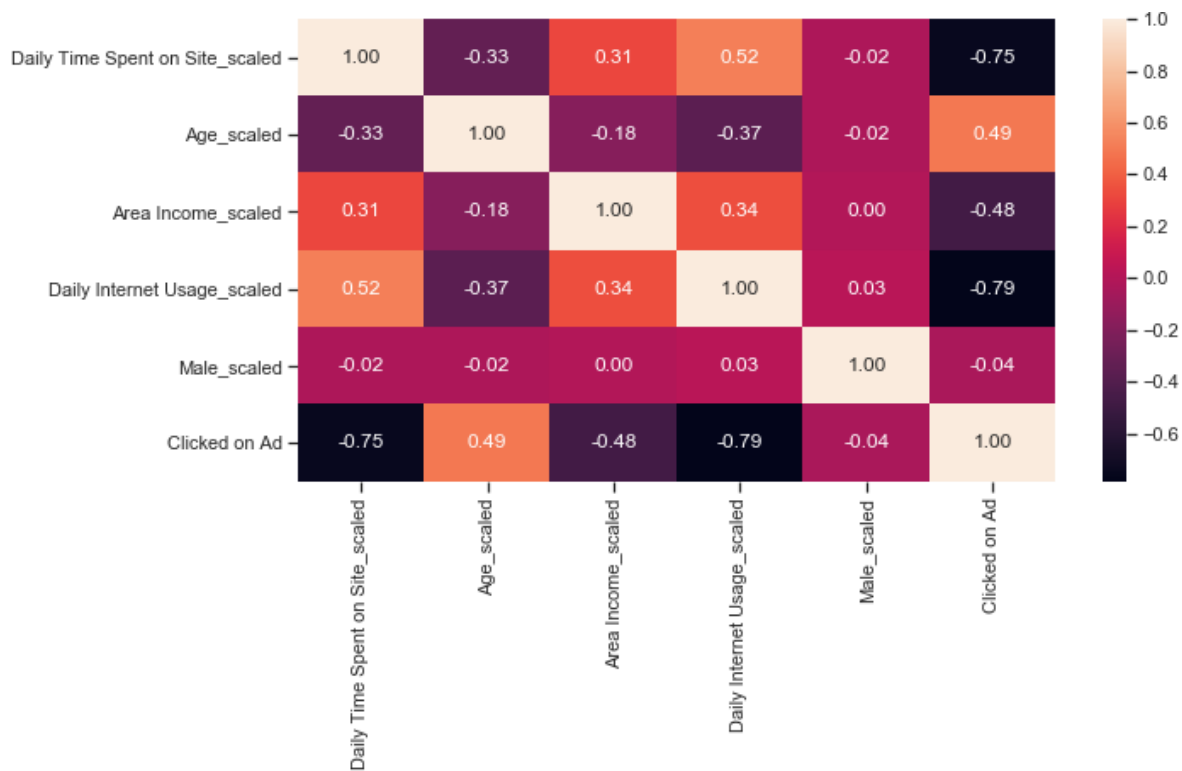
Построим корреляционную матрицу для масштабированных данных

In [54]:

```
1 fig, ax = plt.subplots(figsize=(10,5))
2 sns.heatmap(data_all[corr_cols_2].corr(), annot=True, fmt='.2f')
```

Out[54]:

<matplotlib.axes._subplots.AxesSubplot at 0x2132eb61730>



Выводы о коррелирующих признаках

1. Корреляционные матрицы для исходных и масштабированных данных совпадают

2. Целевой признак классификации ("Clicked on AD") наиболее сильно коррелирует с "Daily Time Spent on Site" (0.75), "Daily Internet Usage" (0.79), "Age" (0.49), "Area Income" (0.48)
3. Между собой признаки коррелируют не достаточно сильно, чтобы их исключать
4. На основании корреляционной матрицы можно сделать вывод о том, что данные позволяют построить модель машинного обучения.

Разработаем класс для визуализации и сохранения меток

In [55]:

```
1 class MetricLogger:
2
3     def __init__(self):
4         self.df = pd.DataFrame(
5             {'metric': pd.Series([], dtype='str'),
6              'alg': pd.Series([], dtype='str'),
7              'value': pd.Series([], dtype='float')})
8
9     def add(self, metric, alg, value):
10         """
11         Добавление значения
12         """
13         # Удаление значения если оно уже было ранее добавлено
14         self.df.drop(self.df[(self.df['metric']==metric)&(self.df['alg']==alg)].index,
15                     # Добавление нового значения
16                     temp = [{'metric':metric, 'alg':alg, 'value':value}]
17                     self.df = self.df.append(temp, ignore_index=True)
18
19     def get_data_for_metric(self, metric, ascending=True):
20         """
21         Формирование данных с фильтром по метрике
22         """
23         temp_data = self.df[self.df['metric']==metric]
24         temp_data_2 = temp_data.sort_values(by='value', ascending=ascending)
25         return temp_data_2['alg'].values, temp_data_2['value'].values
26
27     def plot(self, str_header, metric, ascending=True, figsize=(5, 5)):
28         """
29         Вывод графика
30         """
31         array_labels, array_metric = self.get_data_for_metric(metric, ascending)
32         fig, ax1 = plt.subplots(figsize=figsize)
33         pos = np.arange(len(array_metric))
34         rects = ax1.barh(pos, array_metric,
35                          align='center',
36                          height=0.5,
37                          tick_label=array_labels)
38         ax1.set_title(str_header)
39         for a,b in zip(pos, array_metric):
40             plt.text(0.5, a-0.05, str(round(b,3)), color='white')
41         plt.show()
```

Разделим обучающую и тестовую выборку

Заметим, что данные в этих выборках масштабируемые

```
1 # На основе масштабированных данных выделим обучающую и тестовую выборки с помощью фил
2 train_data_all = data_all[data_all['dataset']=="TRAIN"]
3 test_data_all = data_all[data_all['dataset']=="TEST"]
4 train_data_all.shape, test_data_all.shape
```

Out[56]:

((700, 13), (300, 13))

Задачи классификации и регрессии

Выделим признак для задачи классификации

In [60]:

```
1 # Признаки для задачи классификации
2 task_clas_cols = ['Daily Time Spent on Site_scaled', 'Age_scaled',
3                  'Area Income_scaled', 'Daily Internet Usage_scaled']
```

Определим выборки для задачи классификации

In [61]:

```
1 clas_X_train = train_data_all[task_clas_cols]
2 clas_X_test = test_data_all[task_clas_cols]
3 clas_Y_train = train_data_all['Clicked on Ad']
4 clas_Y_test = test_data_all['Clicked on Ad']
5 clas_X_train.shape, clas_X_test.shape, clas_Y_train.shape, clas_Y_test.shape
```

Out[61]:

((700, 4), (300, 4), (700,), (300,))

Выделим признак для задачи регрессии

Целевым признаком будет "Daily Time Spent on Site".

Для решения задачи регрессии возьмем параметры, которые наиболее сильно коррелируют с целевым. К таким признакам относятся: "Clicked on Ad" (0.75), "Daily Internet Usage_scaled" (0.52), "Age_scaled" (0.33), "Area Income_scaled" (0.31)

In [62]:

```
1 # Признаки для задачи регрессии
2 task_regr_cols = ['Age_scaled', 'Area Income_scaled',
3                  'Daily Internet Usage_scaled', 'Clicked on Ad']
```

Определим выборки для задачи регрессии

```
1 # Выборки для задачи регрессии
2 regr_X_train = train_data_all[task_regr_cols]
3 regr_X_test = test_data_all[task_regr_cols]
4 regr_Y_train = train_data_all['Daily Time Spent on Site_scaled']
5 regr_Y_test = test_data_all['Daily Time Spent on Site_scaled']
6 regr_X_train.shape, regr_X_test.shape, regr_Y_train.shape, regr_Y_test.shape
```

Out[63]:

((700, 4), (300, 4), (700,), (300,))

Создадим словарь моделей, которые будем строить

In [64]:

```
1 # Модели
2 clas_models = {'LogR': LogisticRegression(),
3               'KNN_5': KNeighborsClassifier(n_neighbors=5),
4               'SVC': SVC(),
5               'Tree': DecisionTreeClassifier(),
6               'RF': RandomForestClassifier(),
7               'GB': GradientBoostingClassifier() }
```

In [65]:

```
1 # Сохранение метрик
2 clasMetricLogger = MetricLogger()
```

Построим модель класса

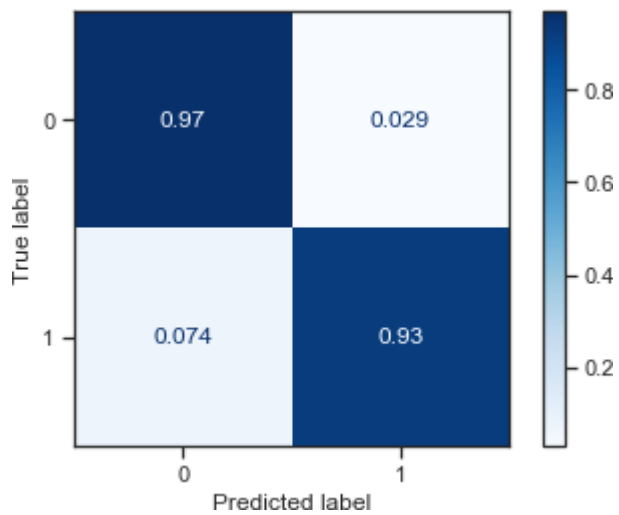
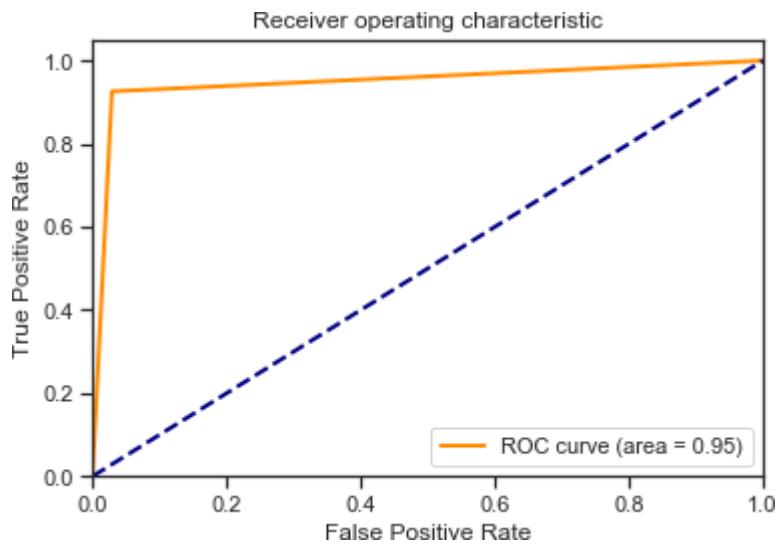
In [66]:

```
1 def clas_train_model(model_name, model, clasMetricLogger):
2     model.fit(clas_X_train, clas_Y_train)
3     Y_pred = model.predict(clas_X_test)
4     precision = precision_score(clas_Y_test.values, Y_pred)
5     recall = recall_score(clas_Y_test.values, Y_pred)
6     f1 = f1_score(clas_Y_test.values, Y_pred)
7     roc_auc = roc_auc_score(clas_Y_test.values, Y_pred)
8
9     clasMetricLogger.add('precision', model_name, precision)
10    clasMetricLogger.add('recall', model_name, recall)
11    clasMetricLogger.add('f1', model_name, f1)
12    clasMetricLogger.add('roc_auc', model_name, roc_auc)
13
14    print('*****')
15    print(model)
16    print('*****')
17    draw_roc_curve(clas_Y_test.values, Y_pred)
18
19    plot_confusion_matrix(model, clas_X_test, clas_Y_test.values,
20                          display_labels=['0', '1'],
21                          cmap=plt.cm.Blues, normalize='true')
22    plt.show()
```

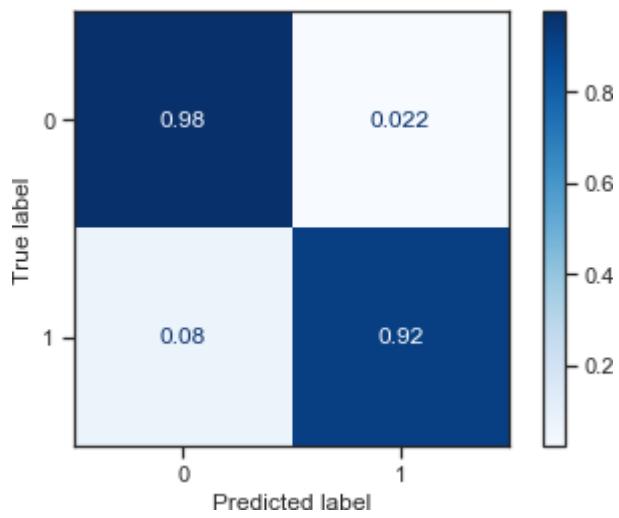
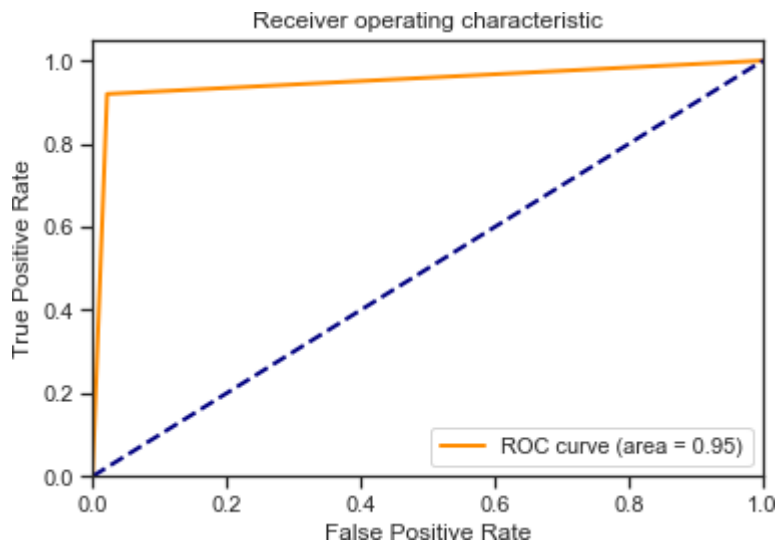
In [67]:

```
1 for model_name, model in clas_models.items():
2     clas_train_model(model_name, model, clasMetricLogger)
```

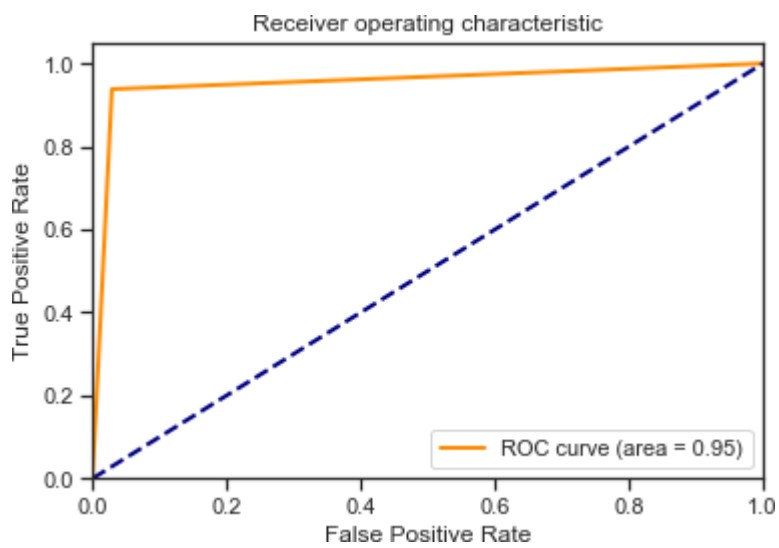
LogisticRegression()

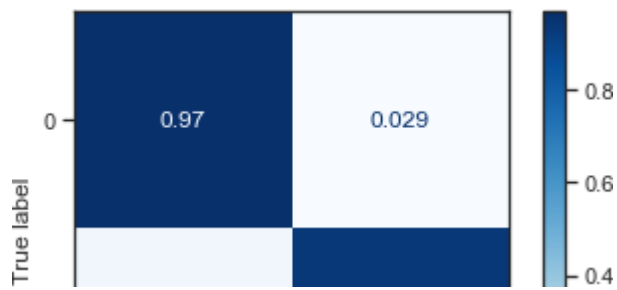


KNeighborsClassifier()

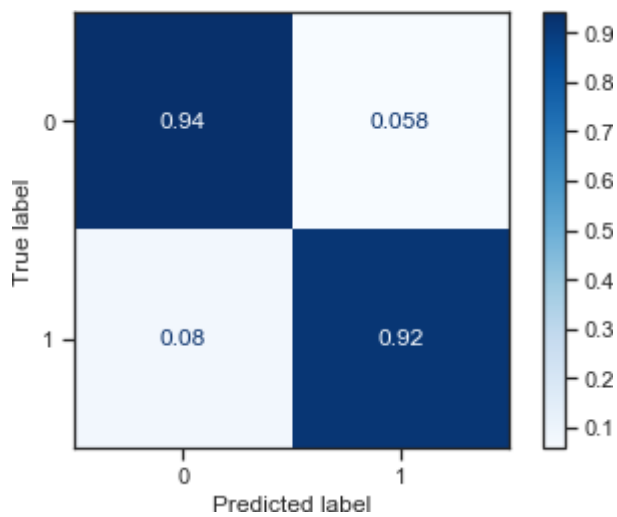
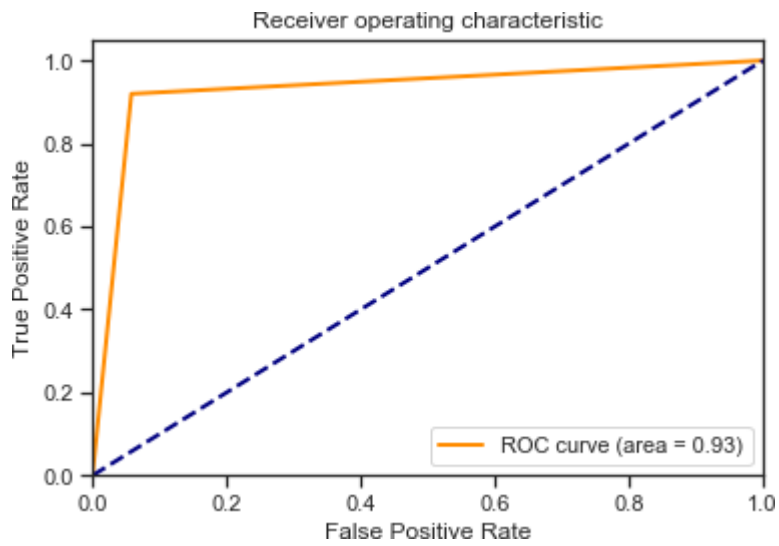


SVC()

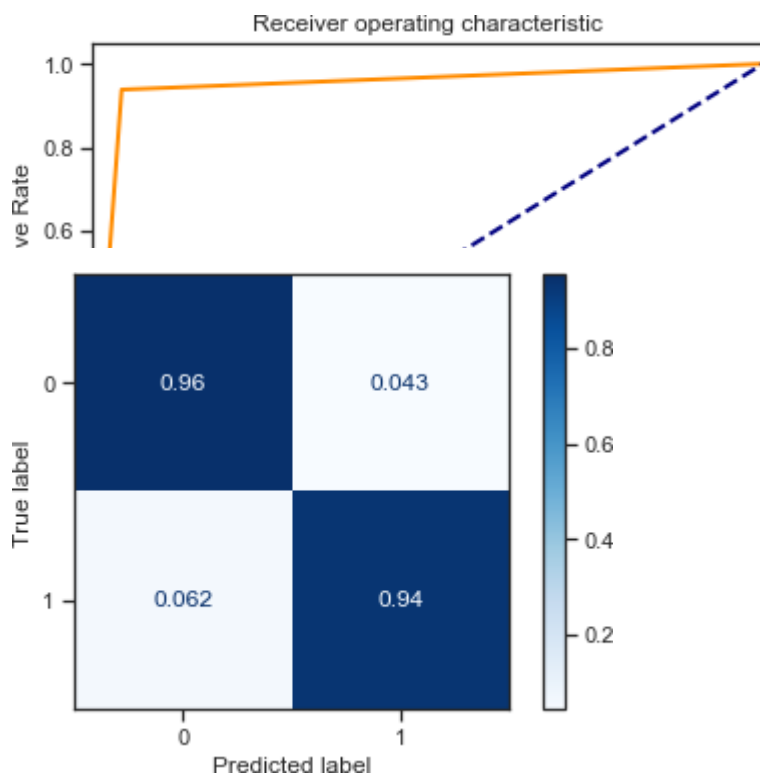




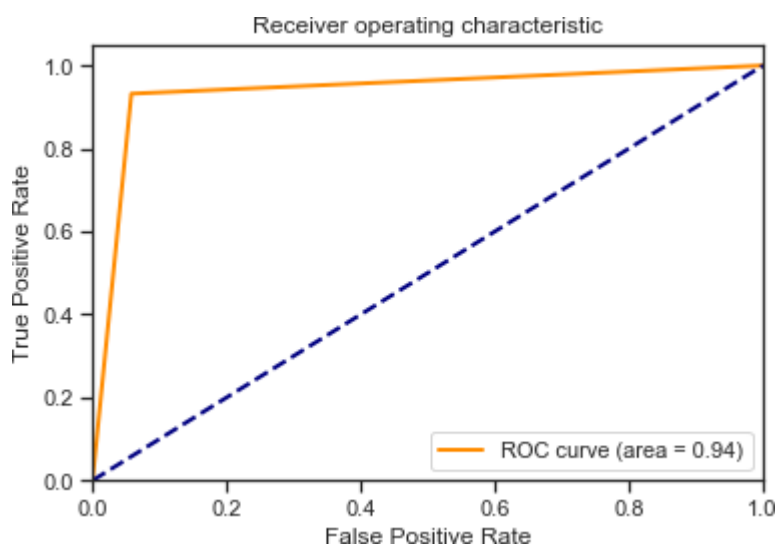
DecisionTreeClassifier()

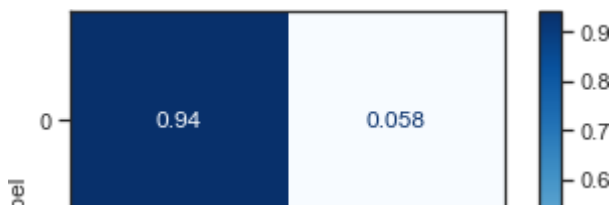


RandomForestClassifier()



GradientBoostingClassifier()





Создадим словарь моделей, которые будем строить

In [68]:

```
1 # Модели
2 regr_models = {'LR': LinearRegression(),
3               'KNN_5': KNeighborsRegressor(n_neighbors=5),
4               'SVR': SVR(),
5               'Tree': DecisionTreeRegressor(),
6               'RF': RandomForestRegressor(),
7               'GB': GradientBoostingRegressor() }
```

In [69]:

```
1 # Сохранение метрик
2 regrMetricLogger = MetricLogger()
```

In [70]:

```
1 def regr_train_model(model_name, model, regrMetricLogger):
2     model.fit(regr_X_train, regr_Y_train)
3     Y_pred = model.predict(regr_X_test)
4
5     mae = mean_absolute_error(regr_Y_test, Y_pred)
6     mse = mean_squared_error(regr_Y_test, Y_pred)
7     r2 = r2_score(regr_Y_test, Y_pred)
8
9     regrMetricLogger.add('MAE', model_name, mae)
10    regrMetricLogger.add('MSE', model_name, mse)
11    regrMetricLogger.add('R2', model_name, r2)
12
13    print('*****')
14        print(model)
15        print()
16        print('MAE={ }, MSE={ }, R2={ }'.format(
17            round(mae, 3), round(mse, 3), round(r2, 3)))
18    print('*****')
```

```
1 for model_name, model in regr_models.items():
2     regr_train_model(model_name, model, regrMetricLogger)
```

*****LinearRegression()

MAE=0.147, MSE=0.036, R2=0.525

*****KNeighborsRegressor()

MAE=0.098, MSE=0.026, R2=0.661

*****SVR()

MAE=0.097, MSE=0.022, R2=0.708

*****DecisionTreeRegressor()

MAE=0.112, MSE=0.042, R2=0.453

*****RandomForestRegressor()

MAE=0.086, MSE=0.024, R2=0.69

*****GradientBoostingRegressor()

MAE=0.112, MSE=0.026, R2=0.653

Подбор гиперпараметров для выбранных моделей

Пример для задачи классификации

In [72]:

```
1 clas_X_train.shape
```

Out[72]:

(700, 4)

```

1 n_range = np.array(range(1,450,10))
2 tuned_parameters = [{'n_neighbors': n_range}]
3 tuned_parameters

```

Out[130]:

```

[{'n_neighbors': array([ 1, 11, 21, 31, 41, 51, 61, 71, 81, 91, 101, 111, 121,
131, 141, 151, 161, 171, 181, 191, 201, 211, 221, 231, 241, 251,
261, 271, 281, 291, 301, 311, 321, 331, 341, 351, 361, 371, 381,
391, 401, 411, 421, 431, 441])}]

```

In [131]:

```

1 %%time
2 clf_gs = GridSearchCV(KNeighborsClassifier(), tuned_parameters, cv=5, scoring='roc_auc')
3 clf_gs.fit(clas_X_train, clas_Y_train)

```

Wall time: 6 s

Out[131]:

```

GridSearchCV(cv=5, estimator=KNeighborsClassifier(),
param_grid=[{'n_neighbors': array([ 1, 11, 21, 31, 41, 51, 61, 71, 81, 91, 101, 111, 121,
131, 141, 151, 161, 171, 181, 191, 201, 211, 221, 231, 241, 251,
261, 271, 281, 291, 301, 311, 321, 331, 341, 351, 361, 371, 381,
391, 401, 411, 421, 431, 441])}],
scoring='roc_auc')

```

In [132]:

```

1 # Лучшая модель
2 clf_gs.best_estimator_

```

Out[132]:

```
KNeighborsClassifier(n_neighbors=271)
```

In [133]:

```

1 # Лучшее значение параметров
2 clf_gs.best_params_

```

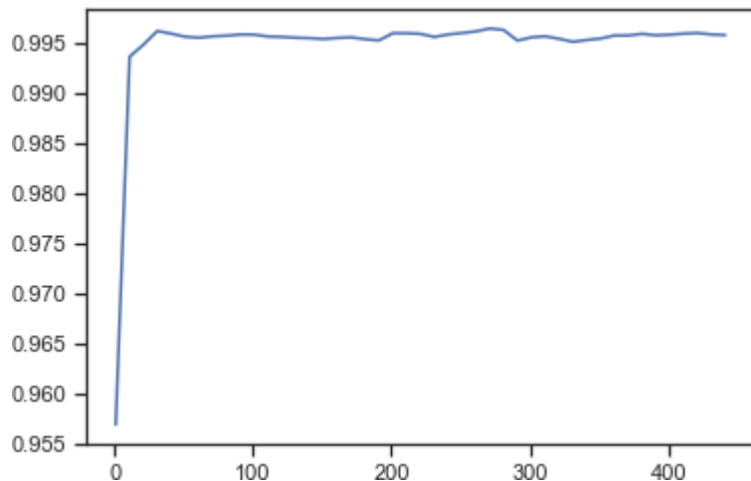
Out[133]:

```
{'n_neighbors': 271}
```

```
1 # Изменение качества на тестовой выборке в зависимости от K-соседей
2 plt.plot(n_range, clf_gs.cv_results_['mean_test_score'])
```

Out[134]:

[<matplotlib.lines.Line2D at 0x2133003d730>]



Пример для задачи регрессии

In [135]:

```
1 n_range = np.array(range(1,450,25))
2 tuned_parameters = [{'n_neighbors': n_range}]
3 tuned_parameters
```

Out[135]:

```
[{'n_neighbors': array([ 1, 26, 51, 76, 101, 126, 151, 176, 201, 226, 251, 276, 301, 326, 351, 376, 401, 426])}]
```

In [136]:

```
1 %%time
2 regr_gs = GridSearchCV(KNeighborsRegressor(), tuned_parameters, cv=5, scoring='neg_me
3 regr_gs.fit(regr_X_train, regr_Y_train)
```

Wall time: 1.68 s

Out[136]:

```
GridSearchCV(cv=5, estimator=KNeighborsRegressor(),
             param_grid=[{'n_neighbors': array([ 1, 26, 51, 76, 101, 12
6, 151, 176, 201, 226, 251, 276, 301,
326, 351, 376, 401, 426])}],
             scoring='neg_mean_squared_error')
```

In [137]:

```
1 # Лучшая модель
2 regr_gs.best_estimator_
```

Out[137]:

```
KNeighborsRegressor(n_neighbors=26)
```

In [103]:

```
1 # Лучшее значение параметров
2 regr_gs.best_params_
```

Out[103]:

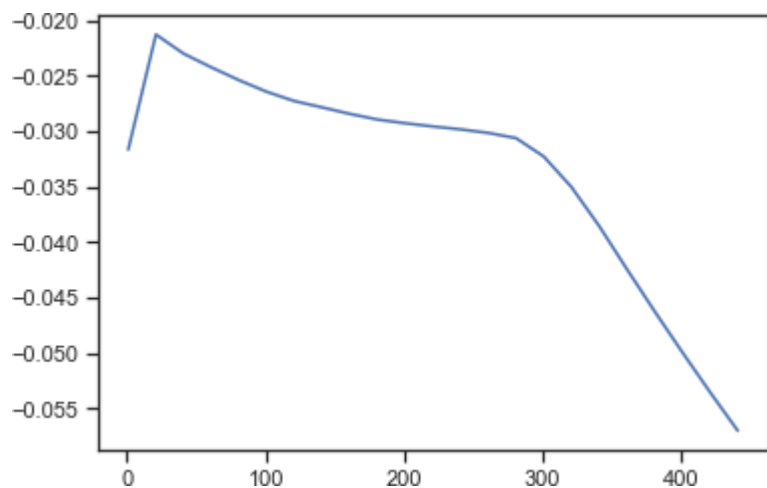
```
{'n_neighbors': 21}
```

In [104]:

```
1 # Изменение качества на тестовой выборке в зависимости от K-соседей
2 plt.plot(n_range, regr_gs.cv_results_['mean_test_score'])
```

Out[104]:

```
[<matplotlib.lines.Line2D at 0x2132ff68a30>]
```



Сравнение качества моделей с гиперпараметрами

Для задачи классификации

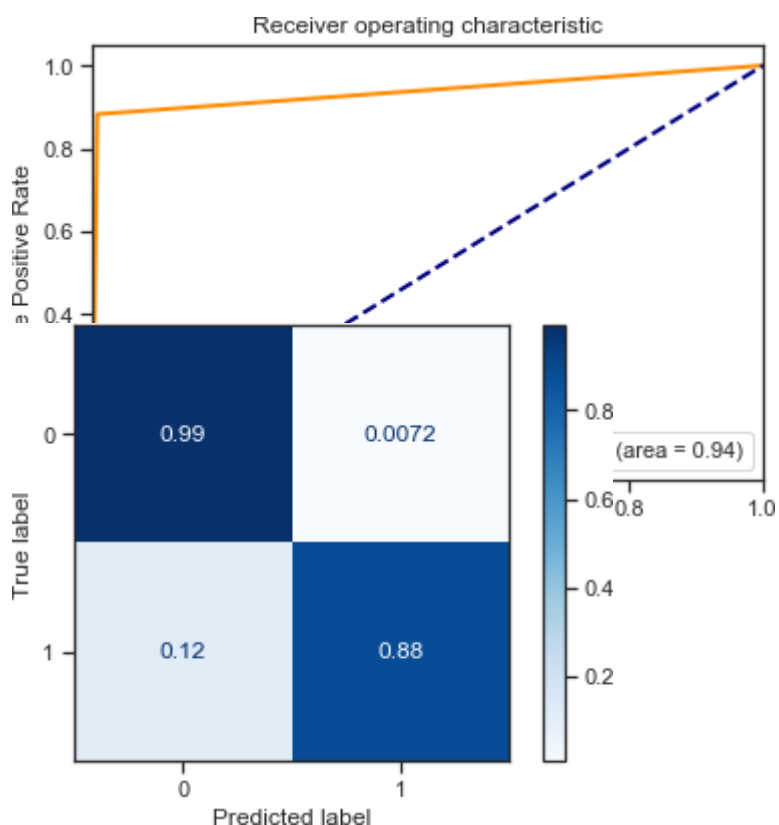
In [116]:

```
1 clas_models_grid = {'KNN_271':clf_gs.best_estimator_}
```

In [117]:

```
1 for model_name, model in clas_models_grid.items():
2     clas_train_model(model_name, model, clasMetricLogger)
```

*****KNeighborsClassifier(n_neighbors=271)



Для задачи регрессии

In [91]:

```
1 regr_models_grid = {'KNN_26':regr_gs.best_estimator_}
```

```
1 for model_name, model in regr_models_grid.items():
2     regr_train_model(model_name, model, regrMetricLogger)
```

*****KNeighborsRegressor(n_neighbors=26)

MAE=0.116, MSE=0.027, R2=0.639

Выводы о качестве моделей на основе выбранных метрик

Задача классификации

In [93]:

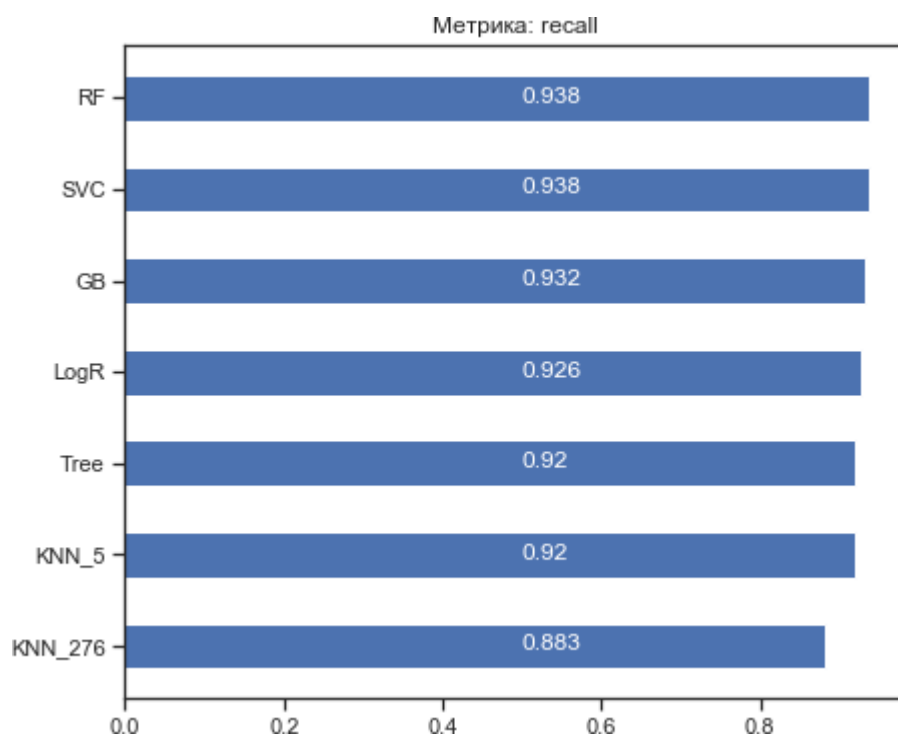
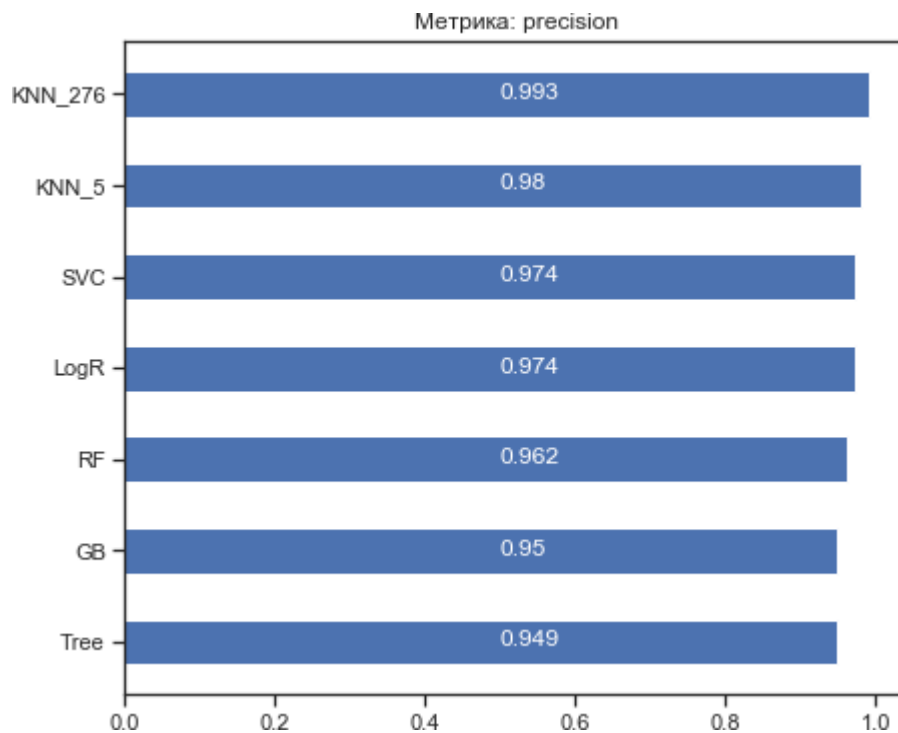
```
1                                     # Метрики качества модели
2     clas_metrics = clasMetricLogger.df['metric'].unique()
3     clas_metrics
```

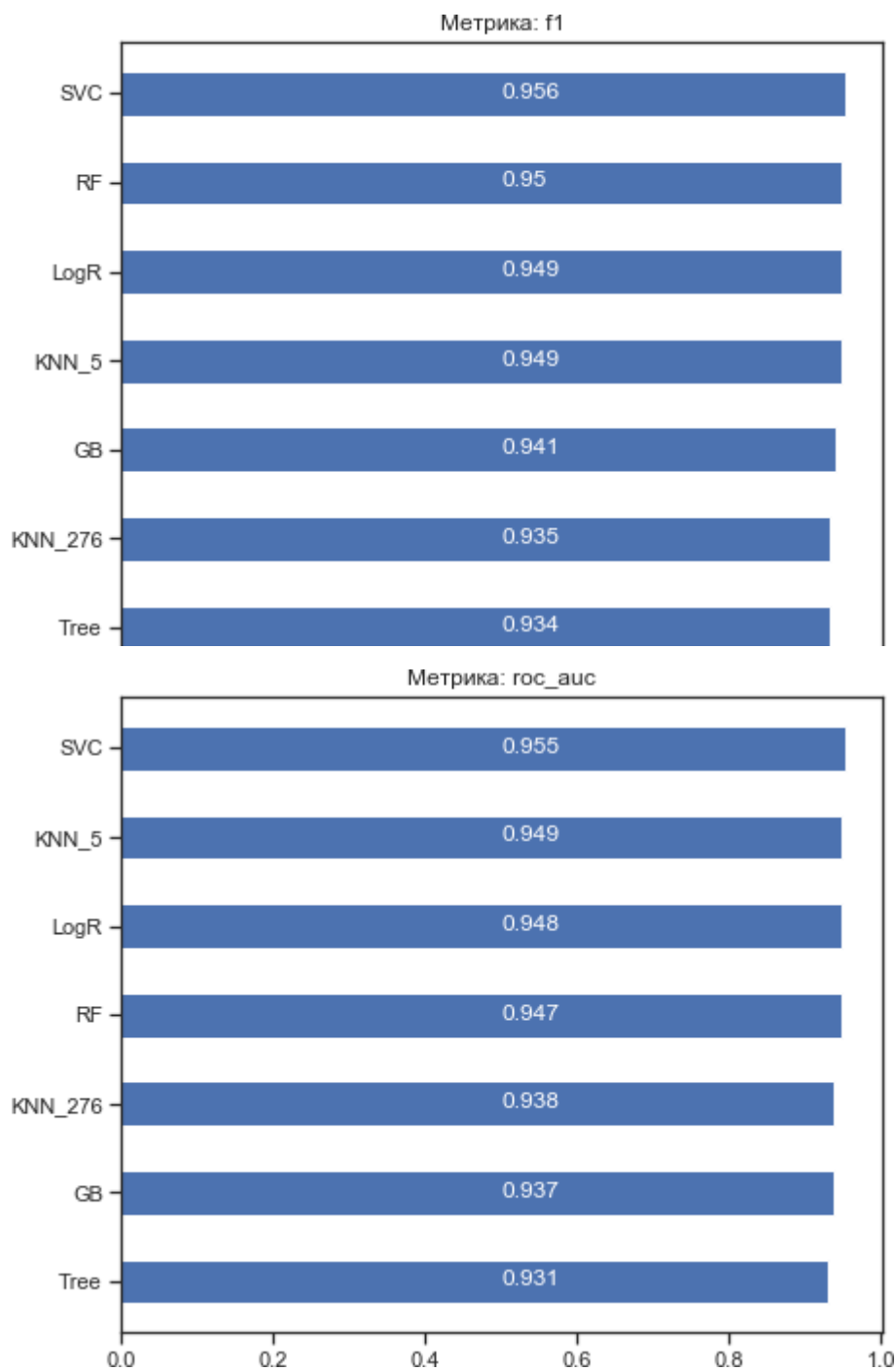
Out[93]:

array(['precision', 'recall', 'f1', 'roc_auc'], dtype=object)

In [94]:

```
1 # Построим графики метрик качества модели
2 for metric in clas_metrics:
3     clasMetricLogger.plot('Метрика: ' + metric, metric, figsize=(7, 6))
```





Вывод

Две из четырех метрики показали, что лучшим является метод опорных векторов (SVC). Однако не для всех метрик данный метод является идеальным.

Задача регрессии

In [95]:

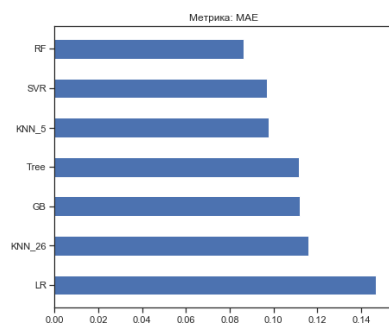
```
1 # Метрики качества модели
2 regr_metrics = regrMetricLogger.df['metric'].unique()
3 regr_metrics
```

Out[95]:

```
array(['MAE', 'MSE', 'R2'], dtype=object)
```

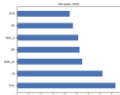
In [96]:

```
1 regrMetricLogger.plot('Метрика: ' + 'MAE', 'MAE', ascending=False, figsize=(7, 6))
```



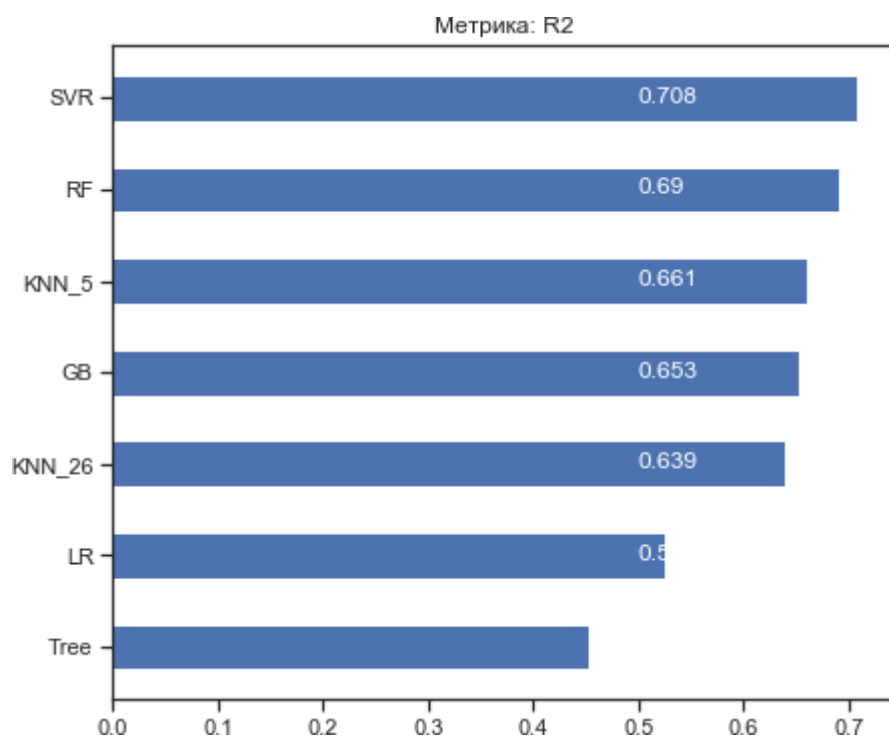
In [98]:

```
1 regrMetricLogger.plot('Метрика: ' + 'MSE', 'MSE', ascending=False, figsize=(7, 6))
```



In [99]:

```
1 regrMetricLogger.plot('Метрика: ' + 'R2', 'R2', ascending=True, figsize=(7, 6))
```



Вывод

1	Очень сложно определить наилучший метод, так как для какой-то метрики данный метод идеален, для другой же он совершенно не подходит. Одним из самых подходящих методов является метод к ближайших соседей.
---	--