



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 \_\_\_\_\_ Радиотехнический (РТ) \_\_\_\_\_

КАФЕДРА \_\_\_\_\_ Системы обработки информации и управления (ИУ-5) \_\_\_\_\_

**РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА  
К НАУЧНО-ИССЛЕДОВАТЕЛЬСКОЙ РАБОТЕ**

**НА ТЕМУ:**

***Прогнозирование оформления клиентом  
срочного депозита***

Студент \_\_\_\_\_  
\_\_\_\_\_ (Группа)

\_\_\_\_\_  
(Подпись, дата) \_\_\_\_\_  
Т. М. Алиев  
(И.О.Фамилия)

Руководитель

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

Консультант

\_\_\_\_\_  
(Подпись, дата) \_\_\_\_\_  
(И.О.Фамилия)

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)  
« \_\_\_\_ » \_\_\_\_\_ 2022 \_\_\_\_ г.

**ЗАДАНИЕ**  
**на выполнение научно-исследовательской работы**

по теме Прогнозирование оформления клиентом срочного депозита

Студент группы РТ5-61Б

Алиев Тимур Муратович  
(Фамилия, имя, отчество)

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

Источник тематики (кафедра, предприятие, НИР) кафедра

График выполнения НИР: 25% к \_\_\_\_ нед., 50% к \_\_\_\_ нед., 75% к \_\_\_\_ нед., 100% к \_\_\_\_ нед.

Техническое задание решить задачу классификации на основе материалов дисциплины по  
выбранной предметной области

**Оформление научно-исследовательской работы:**

Расчетно-пояснительная записка на 23 листах формата А4.

Перечень графического (иллюстративного) материала (чертежи, плакаты, слайды и т.п.)

Дата выдачи задания « 15 » февраля 2022 г.

**Руководитель НИР**

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

**Студент**

Т. М. Алиев  
(Подпись, дата) (И.О.Фамилия)

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

## Содержание

1) Введение.....	4
2) Описание датасета.....	4
3) Импорт библиотек.....	5
4) Загрузка данных .....	5
5) Проведение разведочного анализа данных .....	5
6) Построение графиков для понимания структуры данных.....	6
7) Выбор признаков, подходящих для построения моделей. Кодирование категориальных признаков. Масштабирование данных. Формирование вспомогательных признаков, улучшающих качество моделей .....	10
8) Проведение корреляционного анализа данных. Формирование промежуточных выводов о возможности построения моделей машинного обучения .....	14
9) Выбор метрик для последующей оценки качества моделей .....	15
10) Сохранение и визуализация метрик.....	16
11) Выбор наиболее подходящих моделей для решения задачи классификации или регрессии .....	17
12) Формирование обучающей и тестовой выборок на основе исходного набора данных.....	17
13) 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 .....	17
14) Подбор гиперпараметров для выбранных моделей.....	21
15) Сравнение качества полученных моделей с качеством baseline-моделей .....	22
16) Формирование выводов о качестве построенных моделей на основе выбранных метрик.....	22
17) Заключение .....	23
18) Список использованных источников информации .....	23

## 1) Введение

В качестве предметной области был выбран датасет с информацией о клиентах банка. В исследовании будет решаться задача бинарной классификации.

Данные связаны с кампаниями прямого маркетинга португальского банковского учреждения. Маркетинговые кампании были основаны на телефонных звонках. Часто требовалось более одного контакта с одним и тем же клиентом, чтобы узнать, будет ли продукт (срочный банковский депозит) подписан («да») или нет («нет»).

Цель классификации — предсказать, подпишется ли клиент на срочный депозит (переменная  $y$ ).

## 2) Описание датасета

В качестве набора данных мы будем использовать набор данных прогнозирования подпишется ли клиент на срочный депозит: <https://www.kaggle.com/datasets/hariharanpavan/bank-marketing-dataset-analysis-classification>

Данные связаны с кампаниями прямого маркетинга (телефонными звонками) португальского банковского учреждения. Цель классификации — предсказать, подпишется ли клиент на срочный депозит (переменная  $y$ ).

- Age Возраст клиента
- Job Работа заказчика
- Martial Семейное положение клиента
- Education Уровень образования клиента
- Default Имеет ли кредит по умолчанию?
- Housing Если у клиента есть жилищный кредит
- Loan Имеет личный кредит
- Balance Индивидуальный баланс клиента
- Contact Тип связи
- Month Последний контактный месяц года
- Day Последний контактный день недели
- Duration Длительность последнего контакта, в секундах
- Campaign Количество контактов, выполненных во время этой кампании и для этого клиента
- Pdays Количество дней, прошедших с момента последнего контакта с клиентом из предыдущей кампании
- Previous Количество контактов, выполненных до этой кампании и для этого клиента
- Poutcome результат предыдущей маркетинговой кампании
- Y целевой признак (подпишется ли клиент на срочный депозит)

### 3) Импорт библиотек

```
: 1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.preprocessing import MinMaxScaler
6 from sklearn.linear_model import LinearRegression, LogisticRegression
7 from sklearn.model_selection import train_test_split
8 from sklearn.neighbors import KNeighborsRegressor, KNeighborsClassifier
9 from sklearn.metrics import accuracy_score, balanced_accuracy_score
10 from sklearn.metrics import precision_score, recall_score, f1_score, classification_report
11 from sklearn.metrics import confusion_matrix
12 from sklearn.metrics import plot_confusion_matrix
13 from sklearn.model_selection import GridSearchCV
14 from sklearn.metrics import mean_absolute_error, mean_squared_error, mean_squared_log_error, median_absolute_error, r2_score
15 from sklearn.metrics import roc_curve, roc_auc_score
16 from sklearn.svm import SVC, NuSVC, LinearSVC, OneClassSVM, SVR, NuSVR, LinearSVR
17 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor, export_graphviz
18 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
19 from sklearn.ensemble import ExtraTreesClassifier, ExtraTreesRegressor
20 from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor
21 from sklearn.preprocessing import LabelEncoder
22 %matplotlib inline
23 sns.set(style="ticks")
24 import warnings
25 warnings.filterwarnings('ignore')
```

### 4) Загрузка данных

```
1 first_data = pd.read_csv('bank-full.csv')
```

```
1 # Удалим дубликаты записей, если они присутствуют
2 data = first_data.drop_duplicates()
```

### 5) Проведение разведочного анализа данных

Основные характеристики датасета

```
1 data.head()
```

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
0	58	management	married	tertiary	no	2143	yes	no	unknown	5	may	261	1	-1	0	unknown	no
1	44	technician	single	secondary	no	29	yes	no	unknown	5	may	151	1	-1	0	unknown	no
2	33	entrepreneur	married	secondary	no	2	yes	yes	unknown	5	may	76	1	-1	0	unknown	no
3	47	blue-collar	married	unknown	no	1506	yes	no	unknown	5	may	92	1	-1	0	unknown	no
4	33	unknown	single	unknown	no	1	no	no	unknown	5	may	198	1	-1	0	unknown	no

```
1 data.shape
```

(45211, 17)

```
1 # Список колонок
2 data.columns
```

```
Index(['age', 'job', 'marital', 'education', 'default', 'balance', 'housing',
      'loan', 'contact', 'day', 'month', 'duration', 'campaign', 'pdays',
      'previous', 'poutcome', 'y'],
      dtype='object')
```

```

1 # Список колонок с типами данных
2 data.dtypes

```

```

age          int64
job          object
marital      object
education    object
default      object
balance      int64
housing      object
loan         object
contact      object
day          int64
month        object
duration     int64
campaign     int64
pdays      int64
previous     int64
poutcome    object
y            object
dtype: object

```

```

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

```

```

age          0
job          0
marital      0
education    0
default      0
balance      0
housing      0
loan         0
contact      0
day          0
month        0
duration     0
campaign     0
pdays      0
previous     0
poutcome    0
y            0
dtype: int64

```

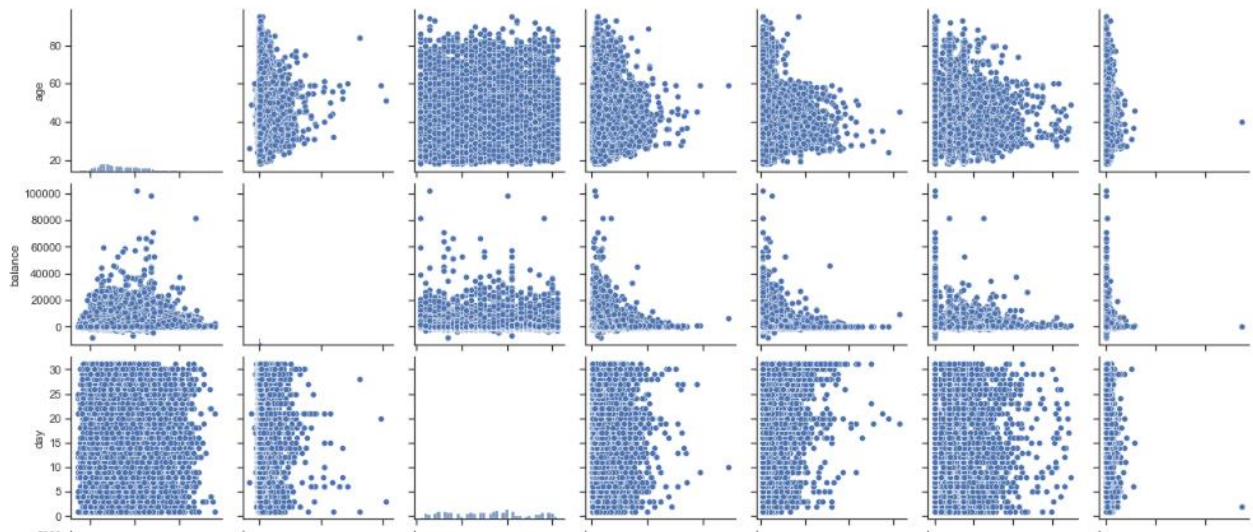
## 6) Построение графиков для понимания структуры данных

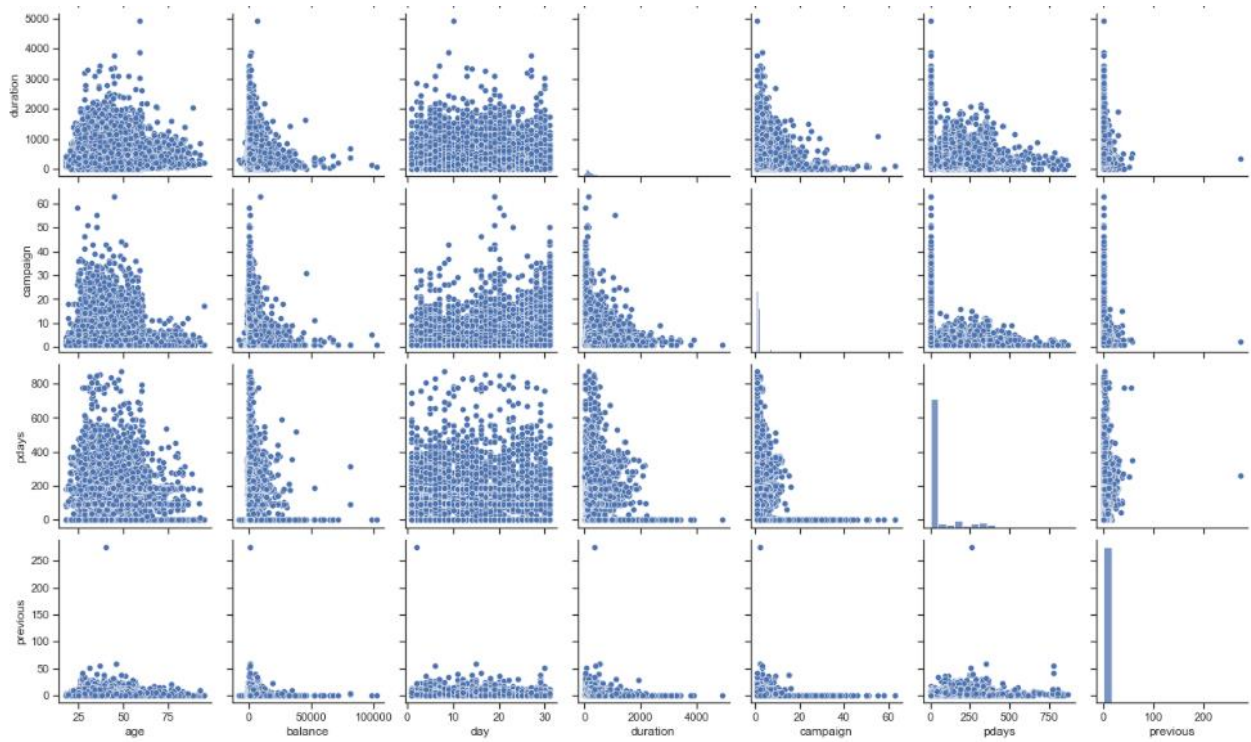
```

1 # Парные диаграммы
2 sns.pairplot(data)

```

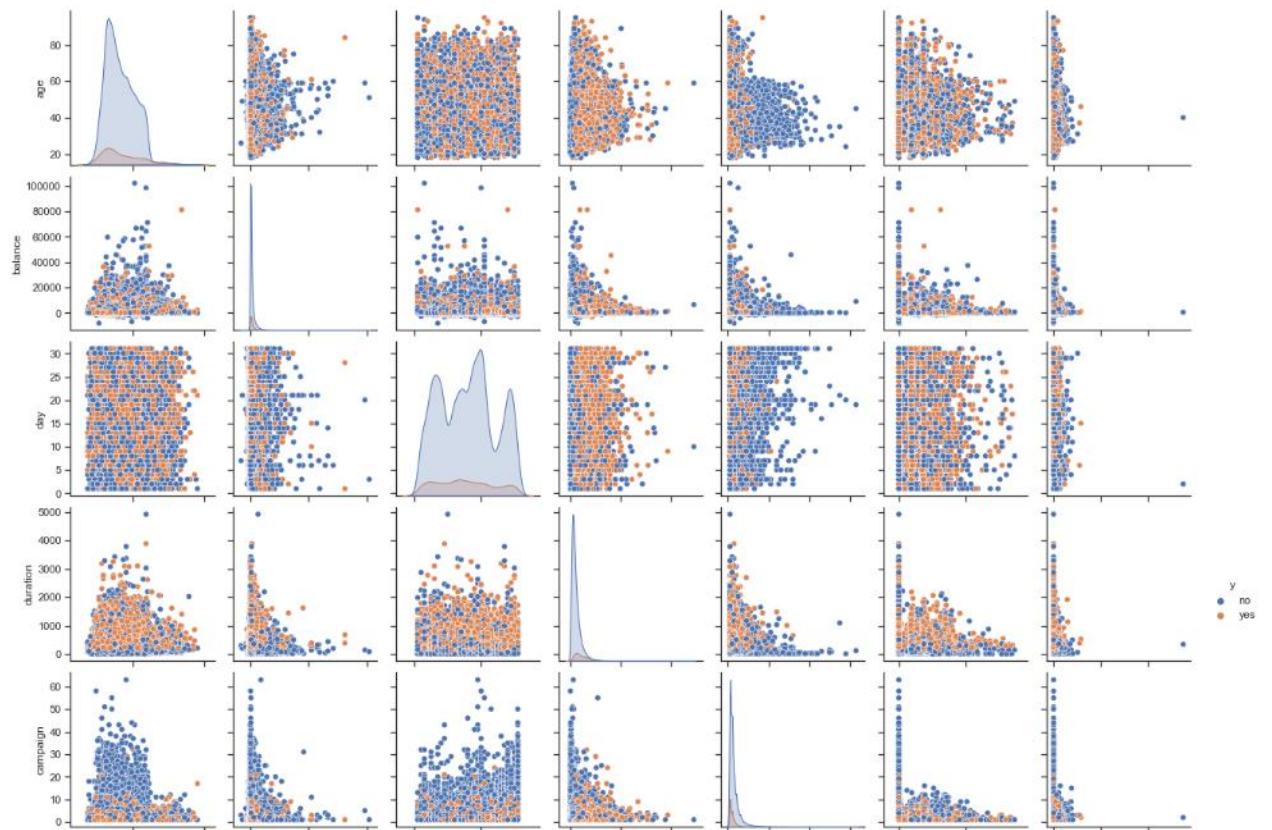
<seaborn.axisgrid.PairGrid at 0x20ba3a47a60>



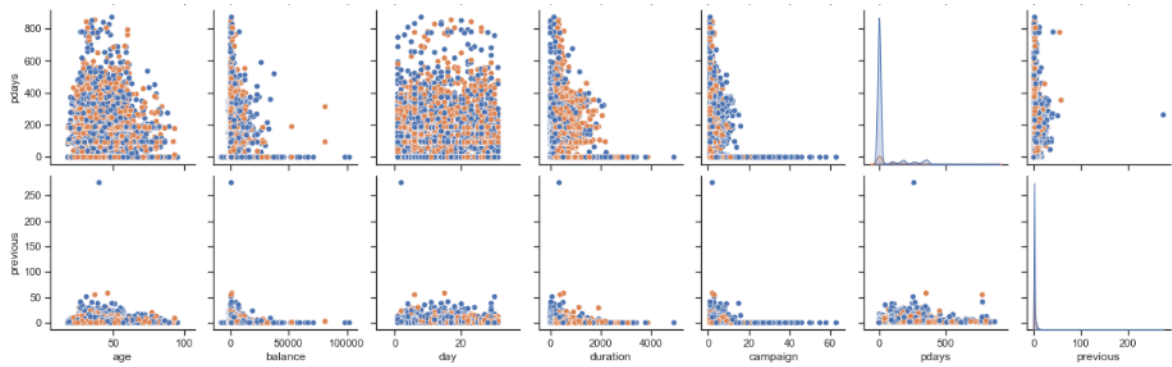


```
1 sns.pairplot(data, hue="y")
```

<seaborn.axisgrid.PairGrid at 0x20bac4a4d90>



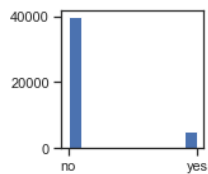




```
1 # Убедимся, что целевой признак
2 # для задачи бинарной классификации содержит только 'no' и 'yes', потом преобразуем к 0 и 1
3 data['y'].unique()
```

```
array(['no', 'yes'], dtype=object)
```

```
1 # Оценим дисбаланс классов для stroke
2 fig, ax = plt.subplots(figsize=(2,2))
3 plt.hist(data['y'])
4 plt.show()
```



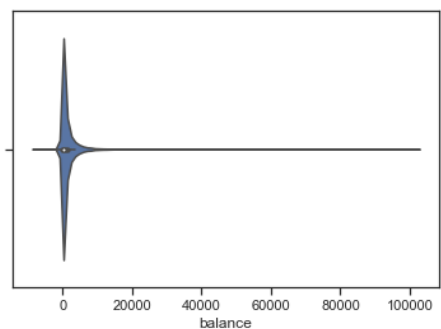
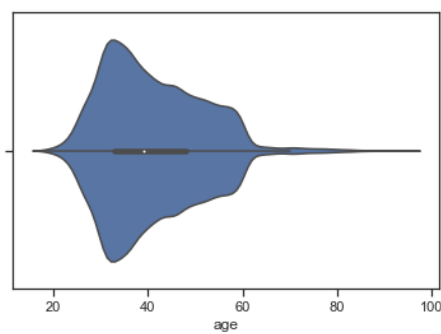
```
1 data['y'].value_counts()
```

```
no      39922
yes      5289
Name: y, dtype: int64
```

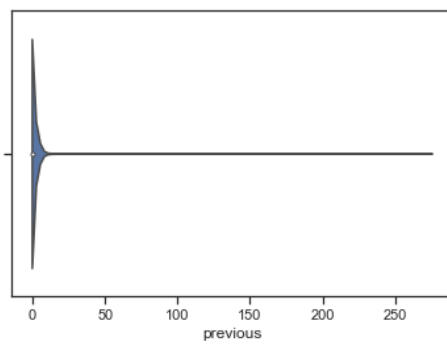
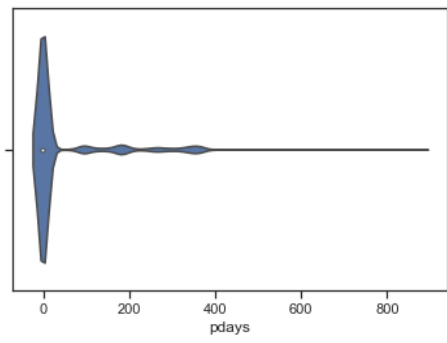
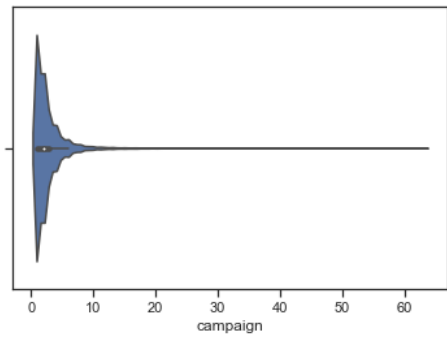
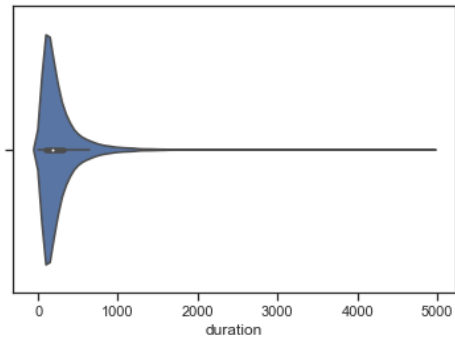
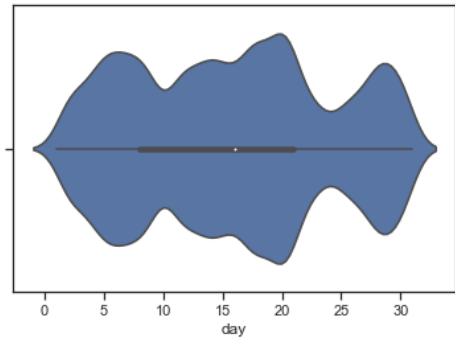
```
1 # посчитаем дисбаланс классов
2 total = data.shape[0]
3 class_0, class_1 = data['y'].value_counts()
4 print('Класс no составляет {}%, а класс yes составляет {}%.'
5       .format(round(class_0 / total, 4)*100, round(class_1 / total, 4)*100))
```

Класс no составляет 88.3%, а класс yes составляет 11.700000000000001%.

```
1 # Скрипичные диаграммы для числовых колонок
2 for col in ['age', 'balance', 'day', 'duration', 'campaign', 'pdays', 'previous']:
3     sns.violinplot(x=data[col])
4     plt.show()
```







## 7) Выбор признаков, подходящих для построения моделей. Кодирование категориальных признаков. Масштабирование данных. Формирование вспомогательных признаков, улучшающих качество моделей

```
1 data.dtypes
```

```
age          int64
job          object
marital      object
education    object
default      object
balance      int64
housing      object
loan         object
contact      object
day          int64
month        object
duration     int64
campaign     int64
pdays       int64
previous     int64
poutcome     object
y            object
dtype: object
```

Категориальные признаки присутствуют, закодируем их.

```
1 data['job'].unique()
```

```
array(['management', 'technician', 'entrepreneur', 'blue-collar',
       'unknown', 'retired', 'admin.', 'services', 'self-employed',
       'unemployed', 'housemaid', 'student'], dtype=object)
```

```
1 job = LabelEncoder()
2 code_job = job.fit_transform(data["job"])
3 data["job"] = code_job
4 data = data.astype({"job": "int64"})
5 np.unique(code_job)
```

```
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
```

```
1 data['marital'].unique()
```

```
array(['married', 'single', 'divorced'], dtype=object)
```

```
1 marital = LabelEncoder()
2 code_marital = marital.fit_transform(data["marital"])
3 data["marital"] = code_marital
4 data = data.astype({"marital": "int64"})
5 np.unique(code_marital)
```

```
array([0, 1, 2])
```

```
1 data['education'].unique()
```

```
array(['tertiary', 'secondary', 'unknown', 'primary'], dtype=object)
```

```
1 education = LabelEncoder()
2 code_education = education.fit_transform(data["education"])
3 data["education"] = code_education
4 data = data.astype({"education": "int64"})
5 np.unique(code_education)
```

```
array([0, 1, 2, 3])
```

```
1 data['default'].unique()
```

```
array(['no', 'yes'], dtype=object)
```

```
1 default = LabelEncoder()
2 code_default = default.fit_transform(data["default"])
3 data["default"] = code_default
4 data = data.astype({"default": "int64"})
5 np.unique(code_default)
```

```
array([0, 1])
```

```
1 data['housing'].unique()
```

```
array(['yes', 'no'], dtype=object)
```

```

1 housing = LabelEncoder()
2 code_housing = housing.fit_transform(data["housing"])
3 data["housing"] = code_housing
4 data = data.astype({"housing": "int64"})
5 np.unique(code_housing)

array([0, 1])

1 data['loan'].unique()

array(['no', 'yes'], dtype=object)

1 loan = LabelEncoder()
2 code_loan = loan.fit_transform(data["loan"])
3 data["loan"] = code_loan
4 data = data.astype({"loan": "int64"})
5 np.unique(code_loan)

array([0, 1])

1 data['contact'].unique()

array(['unknown', 'cellular', 'telephone'], dtype=object)

1 contact = LabelEncoder()
2 code_contact = contact.fit_transform(data["contact"])
3 data["contact"] = code_contact
4 data = data.astype({"contact": "int64"})
5 np.unique(code_contact)

array([0, 1, 2])

1 data['month'].unique()

array(['may', 'jun', 'jul', 'aug', 'oct', 'nov', 'dec', 'jan', 'feb',
      'mar', 'apr', 'sep'], dtype=object)

1 month = LabelEncoder()
2 code_month = month.fit_transform(data["month"])
3 data["month"] = code_month
4 data = data.astype({"month": "int64"})
5 np.unique(code_month)

array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])

1 data['poutcome'].unique()

array(['unknown', 'failure', 'other', 'success'], dtype=object)

1 poutcome = LabelEncoder()
2 code_poutcome = poutcome.fit_transform(data["poutcome"])
3 data["poutcome"] = code_poutcome
4 data = data.astype({"poutcome": "int64"})
5 np.unique(code_poutcome)

array([0, 1, 2, 3])

1 data['y'].unique()

array(['no', 'yes'], dtype=object)

1 y = LabelEncoder()
2 code_y = y.fit_transform(data["y"])
3 data["y"] = code_y
4 data = data.astype({"y": "int64"})
5 np.unique(code_y)

array([0, 1])

1 # Числовые колонки для масштабирования
2 scale_cols = ['age', 'balance', 'day', 'duration', 'campaign', 'pdays', 'previous']

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

```

```

1 # Добавим масштабированные данные в набор данных
2 for i in range(len(scale_cols)):
3     col = scale_cols[i]
4     new_col_name = col + '_scaled'
5     data[new_col_name] = sc1_data[:,i]

```

```
1 data.head()
```

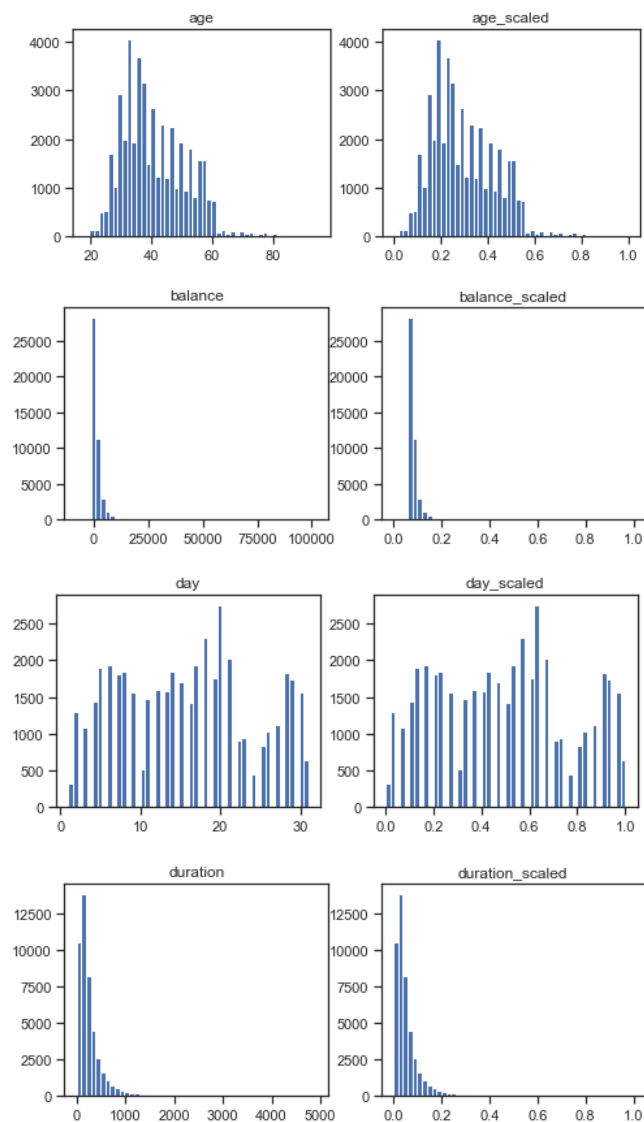
	age	job	marital	education	default	balance	housing	loan	contact	day	...	previous	poutcome	y	age_scaled	balance_scaled	day_scaled	duration_sc
0	58	4	1	2	0	2143	1	0	2	5	...	0	3	0	0.519481	0.092259	0.133333	0.05
1	44	9	2	1	0	29	1	0	2	5	...	0	3	0	0.337662	0.073067	0.133333	0.03
2	33	2	1	1	0	2	1	1	2	5	...	0	3	0	0.194805	0.072822	0.133333	0.01
3	47	1	1	3	0	1506	1	0	2	5	...	0	3	0	0.376623	0.086476	0.133333	0.01
4	33	11	2	3	0	1	0	0	2	5	...	0	3	0	0.194805	0.072812	0.133333	0.04

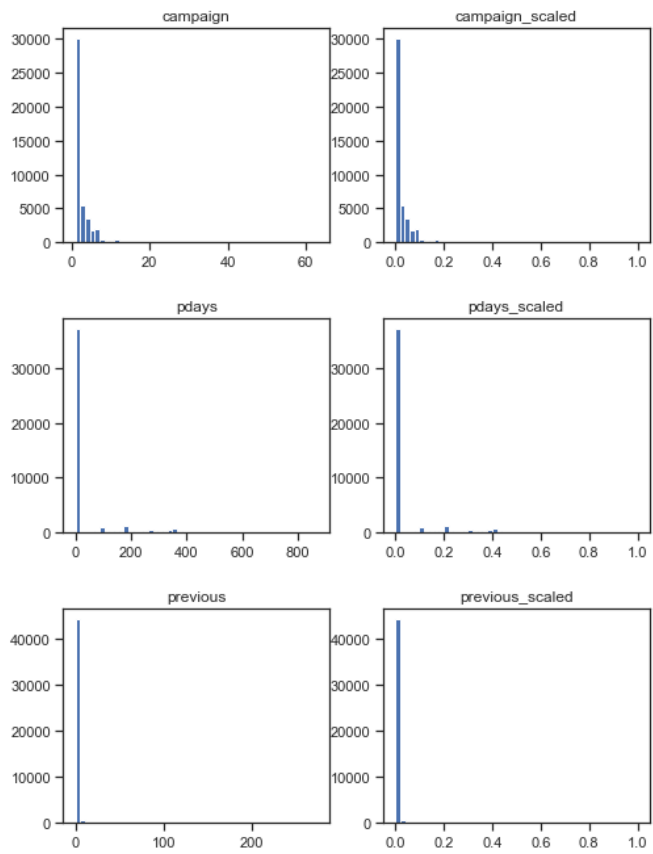
5 rows x 24 columns

```

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[col], 50)
7     ax[1].hist(data[col_scaled], 50)
8     ax[0].title.set_text(col)
9     ax[1].title.set_text(col_scaled)
10    plt.show()

```





## 8) Проведение корреляционного анализа данных. Формирование промежуточных выводов о возможности построения моделей машинного обучения

```
1 # Воспользуемся наличием тестовых выборок,  
2 # включив их в корреляционную матрицу  
3 corr_cols_1 = scale_cols + ['y']  
4 corr_cols_1
```

```
['age', 'balance', 'day', 'duration', 'campaign', 'pdays', 'previous', 'y']
```

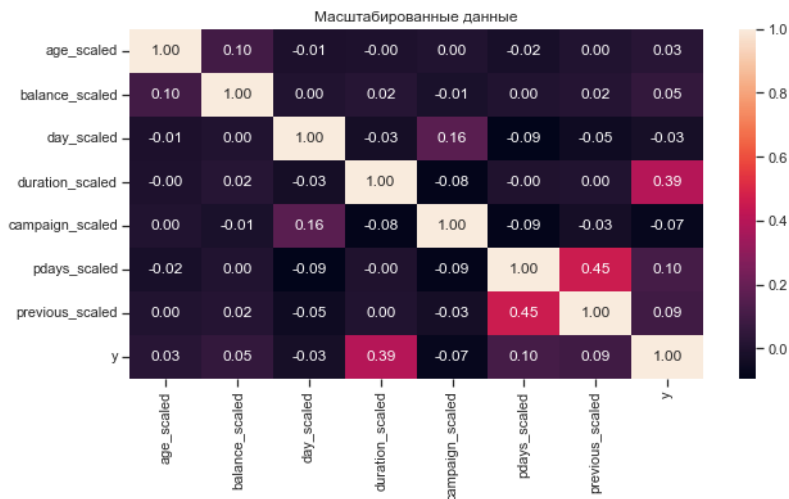
```
1 scale_cols_postfix = [x+'_scaled' for x in scale_cols]  
2 corr_cols_2 = scale_cols_postfix + ['y']  
3 corr_cols_2
```

```
['age_scaled',  
'balance_scaled',  
'day_scaled',  
'duration_scaled',  
'campaign_scaled',  
'pdays_scaled',  
'previous_scaled',  
'y']
```

```
1 fig, ax = plt.subplots(figsize=(10,5))  
2 sns.heatmap(data[corr_cols_1].corr(), annot=True, fmt='.2f')  
3 ax.set_title('Исходные данные (до масштабирования)')  
4 plt.show()
```



```
: 1 fig, ax = plt.subplots(figsize=(10,5))  
2 sns.heatmap(data[corr_cols_2].corr(), annot=True, fmt='.2f')  
3 ax.set_title('Масштабированные данные')  
4 plt.show()
```



На основе корреляционной матрицы можно сделать следующие выводы:

Корреляционные матрицы для исходных и масштабированных данных совпадают.

Целевой признак классификации "y" наиболее сильно коррелирует с duration (0.39) и pdays (0.10). Эти признаки обязательно следует оставить в модели классификации.

## 9) Выбор метрик для последующей оценки качества моделей

В качестве метрик для решения задачи классификации будем использовать:

Метрики, формируемые на основе матрицы ошибок:

Метрика precision: Доля верно предсказанных классификатором положительных объектов, из всех объектов, которые классификатор верно или неверно определил как положительные.

Метрика recall (полнота): Доля верно предсказанных классификатором положительных объектов, из всех действительно положительных объектов.

Метрика F1-мера: Для того, чтобы объединить precision и recall в единую метрику используется Fβ-мера, которая вычисляется как среднее гармоническое от precision и recall:

Метрика ROC AUC:

Идеальная ROC-кривая проходит через точки (0,0)-(0,1)-(1,1), то есть через верхний левый угол графика.

Чем сильнее отклоняется кривая от верхнего левого угла графика, тем хуже качество классификации.

В качестве количественной метрики используется площадь под кривой - ROC AUC (Area Under the Receiver Operating Characteristic Curve). Чем ниже проходит кривая тем меньше ее площадь и тем хуже качество классификатора.

Для получения ROC AUC используется функция roc\_auc\_score.



## 10) Сохранение и визуализация метрик

Разработаем класс, который позволит сохранять метрики качества построенных моделей и реализует визуализацию метрик качества.

```
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, inplace = True)
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()
```

## 11) Выбор наиболее подходящих моделей для решения задачи классификации или регрессии

Для задачи классификации будем использовать следующие модели:

- Логистическая регрессия
- Метод ближайших соседей
- Машина опорных векторов
- Решающее дерево
- Случайный лес
- Градиентный бустинг

## 12) Формирование обучающей и тестовой выборок на основе исходного набора данных

```
1 X_train, X_test, y_train, y_test = train_test_split(data, data.y, random_state=1)
```

```
1 X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

```
((33908, 24), (33908,), (11303, 24), (11303,))
```

## 13) 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

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

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

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

```

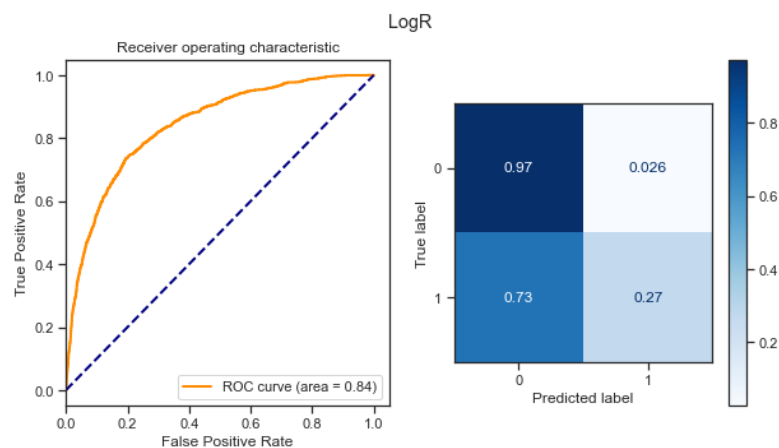
1 def clas_train_model(model_name, model, clasMetricLogger):
2     model.fit(X_train, y_train)
3     # Предсказание значений
4     Y_pred = model.predict(X_test)
5     # Предсказание вероятности класса "1" для roc auc
6     Y_pred_proba_temp = model.predict_proba(X_test)
7     Y_pred_proba = Y_pred_proba_temp[:,1]
8
9     precision = precision_score(y_test.values, Y_pred)
10    recall = recall_score(y_test.values, Y_pred)
11    f1 = f1_score(y_test.values, Y_pred)
12    roc_auc = roc_auc_score(y_test.values, Y_pred_proba)
13
14    clasMetricLogger.add('precision', model_name, precision)
15    clasMetricLogger.add('recall', model_name, recall)
16    clasMetricLogger.add('f1', model_name, f1)
17    clasMetricLogger.add('roc_auc', model_name, roc_auc)
18
19    fig, ax = plt.subplots(ncols=2, figsize=(10,5))
20    draw_roc_curve(y_test.values, Y_pred_proba, ax[0])
21    plot_confusion_matrix(model, X_test, y_test.values, ax=ax[1],
22                          display_labels=['0','1'],
23                          cmap=plt.cm.Blues, normalize='true')
24    fig.suptitle(model_name)
25    plt.show()

```

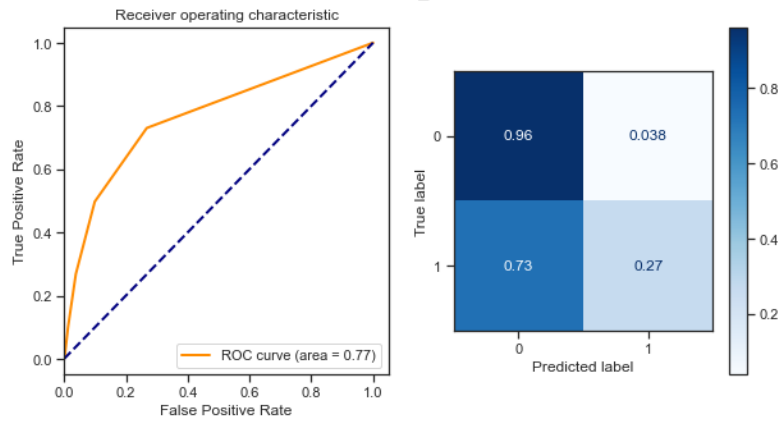
```

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

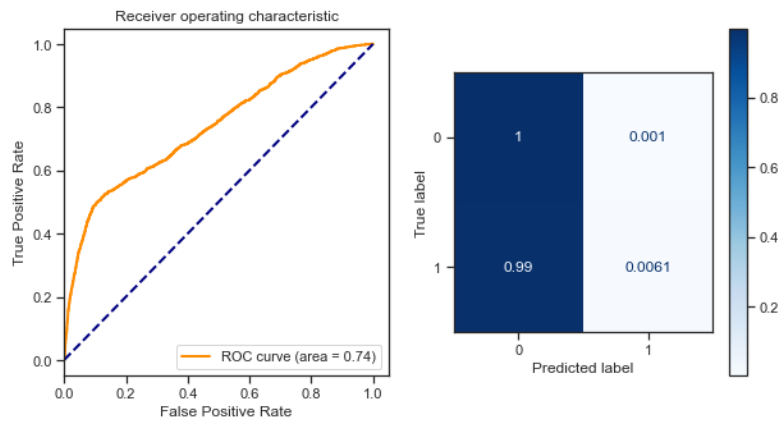
```



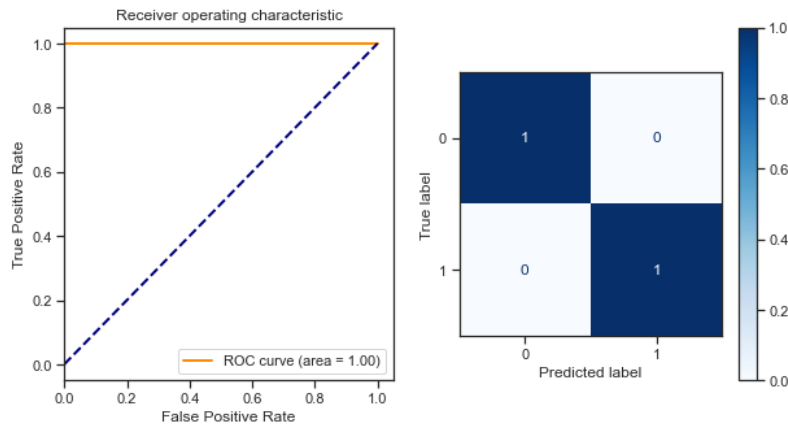
KNN\_5



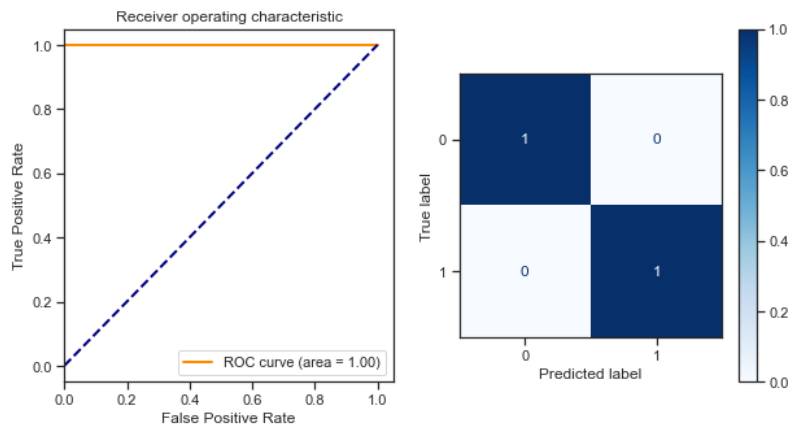
SVC



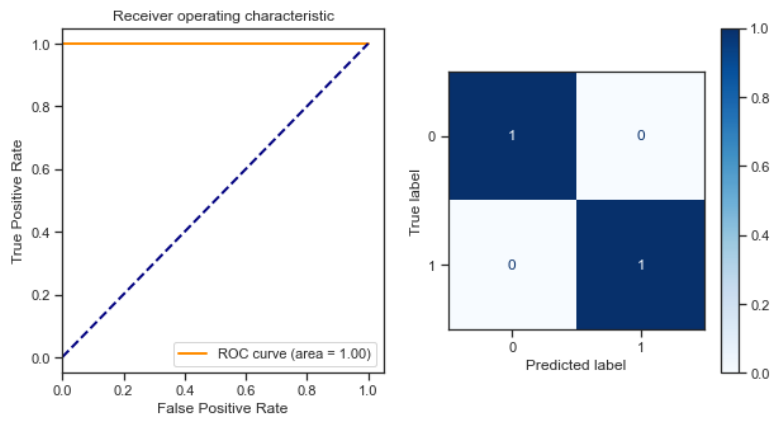
Tree



RF



GB



## 14) Подбор гиперпараметров для выбранных моделей

```
1 X_train.shape
```

```
(33908, 24)
```

```
1 n_range_list = list(range(0,1250,50))
2 n_range_list[0] = 1
```

```
1 n_range = np.array(n_range_list)
2 tuned_parameters = [{'n_neighbors': n_range}]
3 tuned_parameters
```

```
[{'n_neighbors': array([ 1, 50, 100, 150, 200, 250, 300, 350, 400, 450, 500,
550, 600, 650, 700, 750, 800, 850, 900, 950, 1000, 1050,
1100, 1150, 1200])}]
```

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

```
CPU times: total: 9min 21s
```

```
Wall time: 6min 48s
```

```
GridSearchCV(cv=5, estimator=KNeighborsClassifier(),
             param_grid=[{'n_neighbors': array([ 1, 50, 100, 150, 200, 250, 300, 350, 400, 450, 500,
550, 600, 650, 700, 750, 800, 850, 900, 950, 1000, 1050,
1100, 1150, 1200])}],
             scoring='roc_auc')
```

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

```
KNeighborsClassifier(n_neighbors=100)
```

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

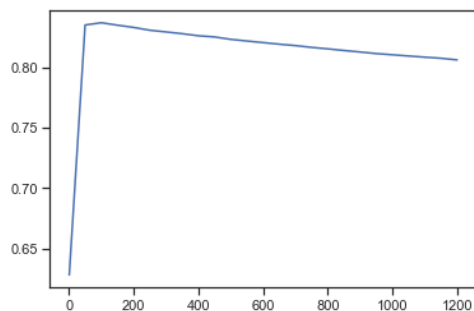
```
{'n_neighbors': 100}
```

```
1 clf_gs_best_params_txt = str(clf_gs.best_params_['n_neighbors'])
2 clf_gs_best_params_txt
```

```
'100'
```

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

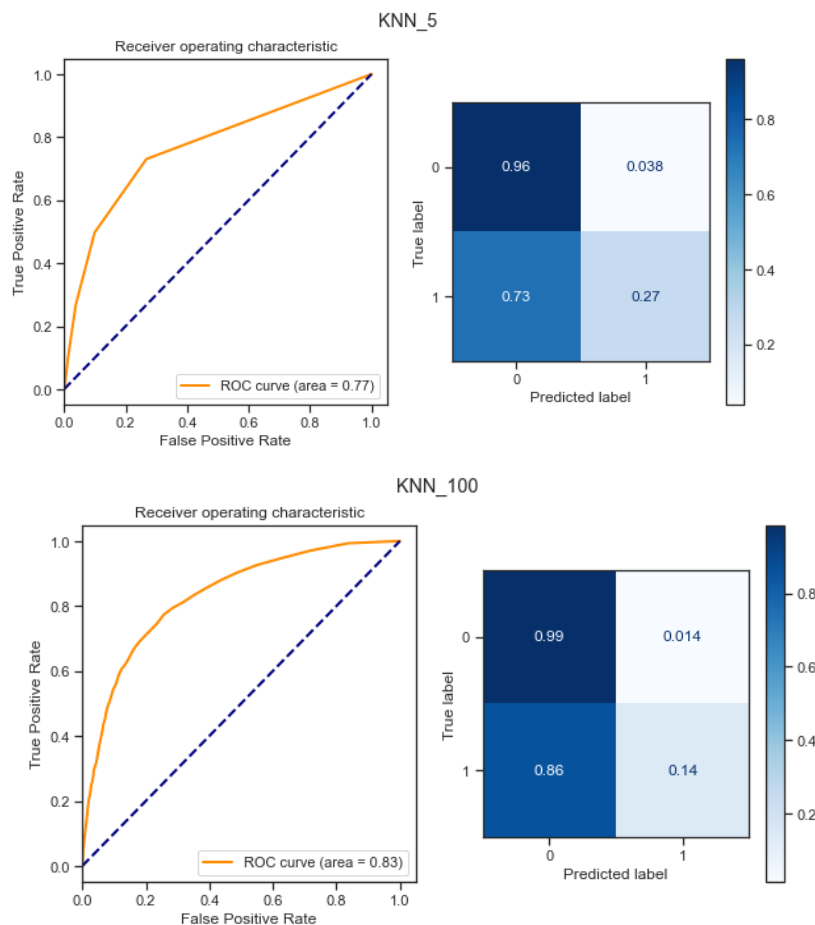
```
[<matplotlib.lines.Line2D at 0x20bb227ab90>]
```



## 15) Сравнение качества полученных моделей с качеством baseline-моделей

```
1 clas_models_grid = {'KNN_5':KNeighborsClassifier(n_neighbors=5),
2                     str('KNN_' + clf_gs_best_params_txt):clf_gs.best_estimator_}
```

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

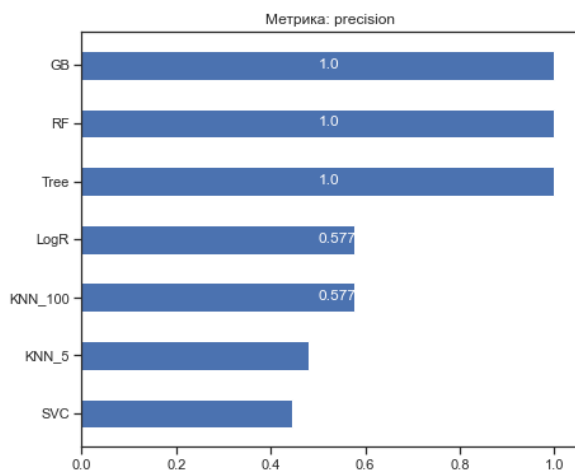


## 16) Формирование выводов о качестве построенных моделей на основе выбранных метрик

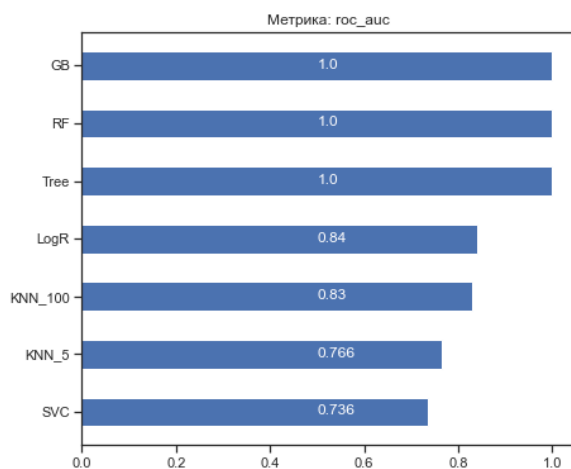
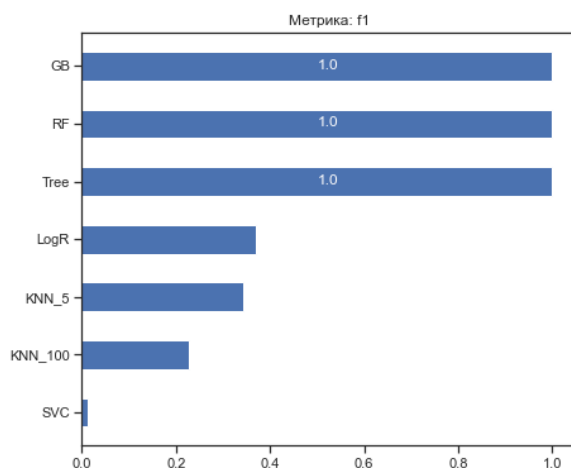
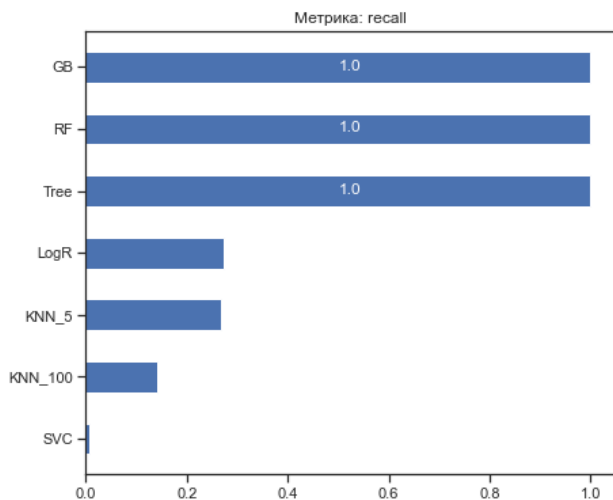
```
1 # Метрики качества модели
2 clas_metrics = clasMetricLogger.df['metric'].unique()
3 clas_metrics
```

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

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







## 17) Заключение

Три модели: градиентный бустинг, дерево и случайный лес показали одинаково высокий результат.

## 18) Список использованных источников информации

1. GitHub репозиторий курса «Технологии машинного обучения» 2022 год.  
URL: [https://github.com/ugapanyuk/ml\\_course\\_2022/wiki/COURSE\\_TMO](https://github.com/ugapanyuk/ml_course_2022/wiki/COURSE_TMO)
2. Matplotlib URL: [https://matplotlib.org/stable/api/as\\_gen/matplotlib.pyplot.html](https://matplotlib.org/stable/api/as_gen/matplotlib.pyplot.html)
3. Scikit-learn URL: <https://scikit-learn.org/stable/index.html>