

Содержание

- 1 Подготовку данных
 - 1.1 Обобщение данных на обучающей выборке и нахожим точность каждой из них на валидационной
 - 1.2 Подготовку данных для машинного обучения
 - 1.3 Создаем учебную, валидационную и тестовые выборки
 - 1.4 Масштабирование
- 2 Исследование задачи
 - 2.1 Обобщение данных на обучающей выборке и нахожим точность каждой из них на валидационной
 - 2.2 Проверка на адекватность моделей
 - 2.2 Борьба с дисбалансом
 - 3.1 Для увеличения точности модели необходимо прийти к балансу классов данных, дающих положительный и отрицательные ответы в целевом признаке
 - 3.2 Обобщение моделей на сбалансированной выборке
 - 4 Тестирование модели

Отток клиентов

Или «Банка» стали уходить клиенты. Каждый месяц. Немного, но заметно. Банковские маркетологи посчитали: сохранять текущих клиентов дешевле, чем привлекать новых.

Нужно спрогнозировать, уйдет клиент из банка в ближайшее время или нет. Вам предоставлены исторические данные о поведении клиентов и условиях договоров с банком.

Постройте модель с предельно большим значением F1-меры. Чтобы сдать проект успешно, нужно довести метрику до 0.59. Проверьте F1-меру на тестовой выборке самостоятельно.

Дополнительно используйте [AUC-ROC](#), сравниваяйте же значение с F1-мерой.

Источники данных: <https://www.kaggle.com/bareillydedicated/bank-customer-churn-modeling>

Подготовка данных

Проверка данных

```
In [1]: # Импортируем необходимые библиотеки для выполнения проекта
import pandas as pd
import numpy as np
import seaborn as sns
from sklearn.cross_validation import train_test_split
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix
from sklearn.metrics import recall_score
from sklearn.metrics import precision_score
from sklearn.metrics import f1_score
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve
from sklearn.metrics import auc
from sklearn.model_selection import GridSearchCV

In [2]: data = pd.read_csv('datasets/Churn.csv')

In [3]: #Общая информация
print('Общая информация:')
display(data.info())
print('Shape: %s')
#Вывод первых пяти строк
print(data.head())
print('Первые пять строк датасета:')
display(data.head())
print('Shape: %s')
print('Название столбцов датасета:')
print('Имя = %s' % data.columns[0])
#Пропуски
print('Информация о пропусках:')
display(data.isnull().sum())
print('Shape: %s')
#Полные дубликаты
print('Полных дубликатов: (data.duplicated().sum())')
print('Shape: %s')

Общая информация:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  --
0   RowNumber              10000 non-null  int64
1   CustomerId             10000 non-null  int64
2   Surname                10000 non-null  object
3   CreditScore             10000 non-null  int64
4   Geography              10000 non-null  object
5   Gender                 10000 non-null  object
6   Age                    10000 non-null  int64
7   Tenure                  9991 non-null   float64
8   Balance                 10000 non-null  float64
9   NumOfProducts          10000 non-null  int64
10  HasCrCard              10000 non-null  int64
11  IsActiveMember         10000 non-null  int64
12  EstimatedSalary         10000 non-null  float64
13  Exited                  10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB

#Полные дубликаты
print('Полных дубликатов: (data.duplicated().sum())')
print('Shape: %s')

Общая информация:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  --
0   RowNumber              10000 non-null  int64
1   CustomerId             10000 non-null  int64
2   Surname                10000 non-null  object
3   CreditScore             10000 non-null  int64
4   Geography              10000 non-null  object
5   Gender                 10000 non-null  object
6   Age                    10000 non-null  int64
7   Tenure                  9991 non-null   float64
8   Balance                 10000 non-null  float64
9   NumOfProducts          10000 non-null  int64
10  HasCrCard              10000 non-null  int64
11  IsActiveMember         10000 non-null  int64
12  EstimatedSalary         10000 non-null  float64
13  Exited                  10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB

#Полные дубликаты
print('Полных дубликатов: (data.duplicated().sum())')
print('Shape: %s')
```

- Приводим название столбцов к нижнему регистру, также преобразуем названия столбцов:

```
In [4]: data = data.rename(columns = {'RowNumber':'row_number', 'CustomerId':'customer_id', 'Surname':'surname', 'CreditScore':'credit_score', 'NumOfProducts':'num_of_products',
                                     'EstimatedSalary':'estimated_salary'})

data.columns = map(str.lower, data.columns)

# Заполняем пропуски нулевыми значениями в значениях Tenure (сколько лет человек является клиентом банка)
data['tenure'] = data['tenure'].fillna(0)

display(data.head())
display(data.info())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  --
0   row_number             10000 non-null  int64
1   customer_id            10000 non-null  int64
2   surname                10000 non-null  object
3   credit_score           10000 non-null  int64
4   geography              10000 non-null  object
5   gender                 10000 non-null  object
6   age                    10000 non-null  int64
7   tenure                  9991 non-null   float64
8   balance                 10000 non-null  float64
9   num_of_products        10000 non-null  int64
10  has_cr_card             10000 non-null  int64
11  is_active_member        10000 non-null  int64
12  estimated_salary        10000 non-null  float64
13  exited                 10000 non-null  int64
dtypes: int64(8), float64(3), object(3)
memory usage: 1.1+ MB

Данные корректны, пропуски и явные дубликаты нет
```

Подготовка данных для машинного обучения

- Устанавливаем целевой признак - 'exited' - факт ухода клиента из Банка
- Удаляем столбцы данных которых не понадобится в исследовании

```
In [5]: data_new = data.drop(['customer_id', 'surname', 'row_number'], axis=1)

display(data_new.head())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  --
0   credit_score           10000 non-null  int64
1   geography              10000 non-null  object
2   gender                 10000 non-null  object
3   age                    10000 non-null  int64
4   tenure                  9991 non-null   float64
5   balance                 10000 non-null  float64
6   num_of_products        10000 non-null  int64
7   has_cr_card             10000 non-null  int64
8   is_active_member        10000 non-null  int64
9   estimated_salary        10000 non-null  float64
10  exited                 10000 non-null  int64
11  Exited                  10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB

Данные корректны, пропуски и явные дубликаты нет
```

- Преобразуем категориальные признаки в числовые для gender, geography и age используя метод ONE:

```
In [6]: data_new = pd.get_dummies(data_new, drop_first=True)

display(data_new.head())
data_new.shape

<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  --
0   credit_score           10000 non-null  int64
1   age                    9991 non-null   float64
2   tenure                  9991 non-null   float64
3   balance                 10000 non-null  float64
4   num_of_products        10000 non-null  int64
5   has_cr_card             10000 non-null  int64
6   is_active_member        10000 non-null  int64
7   estimated_salary        10000 non-null  float64
8   exited                 10000 non-null  int64
9   Exited                  10000 non-null  int64
10  geography_Germany       10000 non-null  int64
11  geography_Spain         10000 non-null  int64
12  geography_Italy          10000 non-null  int64
13  geography_France        10000 non-null  int64
14  geography_UK             10000 non-null  int64
15  gender_Male              10000 non-null  int64
16  gender_Female           10000 non-null  int64
17  gender_Male              10000 non-null  int64
18  gender_Female           10000 non-null  int64
19  gender_Male              10000 non-null  int64
dtypes: float64(3), int64(16), object(2)
memory usage: 1.1+ MB
```

Валидационная выборка (в пропорции 69:40 от общей датасета)

```
In [8]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [9]: features = data_new.drop(['exited'], axis=1)
target = data_new['exited']

display(features.head())
display(target.head())
```

Разделение данных на признаки и целевой признак

```
In [10]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [11]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [12]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [13]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [14]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [15]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [16]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [17]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [18]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [19]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [20]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [21]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [22]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [23]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [24]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [25]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [26]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [27]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [28]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [29]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [30]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [31]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [32]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [33]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [34]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [35]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [36]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [37]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [38]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [39]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [40]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [41]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [42]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [43]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [44]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [45]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [46]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [47]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [48]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [49]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [50]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [51]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [52]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [53]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [54]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [55]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [56]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [57]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [58]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [59]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [60]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [61]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [62]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```

Разделение данных на признаки и целевой признак

```
In [63]: features_train, features_valid, target_train, target_valid_test, = train_test_split(features, target,
                                                train_size=0.66, random_state=12345, stratify=target)
```