

Московский государственный университет имени М.В. Ломоносова

Направление 38.03.01 Экономика

Программа «Национальная экономика»

**Прогнозирование рисков банковского кредитования с
использованием технологий искусственного
интеллекта**

Выпускная квалификационная работа

Студент

Касенова Асем Ардаковна

Научный руководитель:

к.э.н., к.ф.-м.н., к.ю.н., доцент

Сидоренко Владимир Николаевич

Астана

2026

ОГЛАВЛЕНИЕ

Глава 1. Тестирование алгоритмов машинного и глубокого обучения на данных по кредитному скорингу	3
1.1 Подбор гиперпараметров для алгоритмов машинного обучения при оценке возврата кредита	3
1.1.1 Деревья решений	4
1.1.2 Случайный лес	4
1.1.3 Градиентный бустинг	4
1.1.4 Нейронная сеть	4
1.2 Ограничения и возможности моделей	4
Приложение 1. Предобработка данных	6

ГЛАВА 1. ТЕСТИРОВАНИЕ АЛГОРИТМОВ МАШИННОГО И ГЛУБОКОГО ОБУЧЕНИЯ НА ДАННЫХ ПО КРЕДИТНОМУ СКОРИНГУ

В предыдущей главе было выяснено, что алгоритмы градиентного бустинга и нейронных сетей не имеют эффектов переобучения и недообучения. В случае градиентного бустинга используются гиперпараметры по умолчанию, результаты метрик которых не ниже метрик других алгоритмов. Это означает, что не требуется поиск гиперпараметров для данного алгоритма и можно остановиться на гиперпараметрах по умолчанию. Идентичная ситуация возникает у нейронной сети с отличием в том, что у нее метрики точности немного хуже, чем у алгоритма градиентного бустинга. Поэтому в этой главе будут подобраны гиперпараметры для алгоритмов деревьев решений и случайного леса, поскольку они обладают эффектом переобучения. Основным инструментом для поиска оптимальных гиперпараметров используются инструменты *GridSearchCV*¹ и *RandomizedSearchCV*². Также в конце главы будут проанализированы прибыль и убыток банка на основании предсказаний четырех рассмотренных алгоритмов.

1.1 Подбор гиперпараметров для алгоритмов машинного обучения при оценке возврата кредита

Гиперпараметры – это такие параметры модели, которые задаются до начала работы модели. Правильно подобранные параметры позволяют улучшить качество предсказаний моделей. Каждому алгоритму МО присущие собственные наборы гиперпараметров, которые будут рассматриваться ниже. Существуют два метода подбора гиперпараметров:

1. *GridSearchCV* – метод подбора оптимальных гиперпараметров для модели с помощью перебора всех возможных комбинаций из заданного набора, что является ресурсозатратным способом.
2. *RandomizedSearchCV* – метод, позволяющий выбирать количество случайных комбинаций из заданного набора гиперпараметров. Является наиболее эффективным методом по времени при работе с большими данными.

¹GridSearchCV [Электронный ресурс] / Python-библиотека для машинного обучения. URL: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html (дата обращения: 25.11.2025).

²RandomizedSearchC (дата обращения: 25.11.2025).

1.1.1 Деревья решений

Для алгоритма классификации решающих деревьев (*Decision Tree Classifier*) основными гиперпараметрами являются:

1. *criterion* – критерий качества разбиения, определяющая показатель, по которому алгоритм лучше классифицирует клиентов.
 - a) *gini* – индекс Джини³, измеряющий неоднородность двух классов (0 и 1) внутри узла. Если индекс равен 0, значит в узле содержатся объекты одного класса. Чем индекс ближе к 0, тем лучше алгоритм определяет класс клиента. Индекс Джини в алгоритме *Decision Tree* используется по умолчанию, поскольку быстрее вычисляется;
 - b) *entropy* – энтропия Шеннона⁴, измеряющий степень смешанности классов. Высокая энтропия свидетельствует о высокой неоднородности классов, чем ниже энтропия, тем узел более однородный. Алгоритм выбирает разбиения с наиболее низкой энтропией, т.е. содержатся большинство объектов одного класса;
 - c) *log_loss* – логарифмическая функция потерь. В данном случае алгоритм не просто выбирает класс, а оценивает их вероятности. За каждый неверный прогноз вводит штраф.

1.1.2 Случайный лес

1.1.3 Градиентный бустинг

1.1.4 Нейронная сеть

1.2 Ограничения и возможности моделей

³Джини К. (1884-1965 гг.) – итальянский статистик, экономист, социолог и демограф.

⁴Шеннон К.Э. (1916-2001 гг.) – американский инженер, криptoаналитик и математик, заложил основы теории информации.

СПИСОК ЛИТЕРАТУРЫ

1. GridSearchCV [Электронный ресурс] / Python-библиотека для машинного обучения. — URL: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html. — Загл. с экрана.

Приложение 1. Предобработка данных

Листинг 1 – Формирование единого датасета

```
1 import pandas as pd
2 import os
3 import pyarrow.parquet as pq
4
5 path = "train_data"
6 for i,file in enumerate(os.listdir(path)):
7     file = 'train_data_' + str(i) + '.pq'
8     dataset = pq.ParquetDataset(os.path.join(path,file))
9     df = dataset.read(use_threads=True).to_pandas()
10    df_gr = df.groupby('id').agg('mean')
11    file_csv = file.replace('.pq', '.csv')
12    df_gr.to_csv(os.path.join('train_data_csv_all',file_csv))
13 os.listdir(path)
14
15 path = 'train_data_csv_all'
16 frames = []
17 for file_csv in os.listdir(path):
18     if file_csv.endswith('.csv'):
19         df = pd.read_csv(os.path.join(path, file_csv))
20         frames.append(df)
21
22 result = pd.concat(frames)
23 file_csv_all = '1_data_csv_all.csv'
24 result.to_csv(file_csv_all)
25 df_all = pd.read_csv(file_csv_all)
```

Листинг 2 – Датасет, содержащий 41 признак

```
1 from sklearn.decomposition import PCA
2
3 file_csv = 'data_csv_small.csv'
4 df = pd.read_csv(file_csv)
5
6 columns = ['pre_pterm', 'pre_fterm',
7             'pre_loans_next_pay_summ', 'pre_loans_outstanding',
8             'pre_loans_total_overdue', 'pre_loans_max_overdue_sum',
9             'pre_loans_credit_cost_rate', 'is_zero_loans5',
10            'is_zero_loans530', 'is_zero_loans3060',
11            'is_zero_loans6090', 'is_zero_loans90',
12            'pre_util', 'pre_maxover2limit',
13            'is_zero_util', 'is_zero_over2limit',
14            'enc_paym_0', 'enc_paym_1', 'enc_paym_2',
15            'enc_paym_3', 'enc_paym_4', 'enc_paym_5',
16            'enc_paym_6', 'enc_paym_7', 'enc_paym_8',
17            'enc_paym_9', 'enc_paym_10', 'enc_paym_11',
18            'enc_paym_12', 'enc_paym_13', 'enc_paym_14',
19            'enc_paym_15', 'enc_paym_16', 'enc_paym_17',
20            'enc_paym_18', 'enc_paym_19', 'enc_paym_20',
21            'enc_paym_21', 'enc_paym_22',
22            'enc_paym_23', 'enc_paym_24']
23
24 X = df.loc[:, columns].copy()
25 pca = PCA()
26 pca.fit(X)
27
28 100*pca.explained_variance_ratio_.round(3)
```

Листинг 3 – Датасет, содержащий 24 признака

```
1 file_csv = 'data_csv_small.csv'
2 df = pd.read_csv(file_csv)
3
4 columns_pay = ['enc_paym_0', 'enc_paym_1', 'enc_paym_2',
5                 'enc_paym_3', 'enc_paym_4', 'enc_paym_5',
6                 'enc_paym_6', 'enc_paym_7', 'enc_paym_8',
7                 'enc_paym_9', 'enc_paym_10', 'enc_paym_11',
8                 'enc_paym_12', 'enc_paym_13', 'enc_paym_14',
9                 'enc_paym_15', 'enc_paym_16', 'enc_paym_17',
10                'enc_paym_18', 'enc_paym_19', 'enc_paym_20',
11                'enc_paym_21', 'enc_paym_22', 'enc_paym_23',
12                'enc_paym_24']
13
14 X_pay = df.loc[:, columns_pay].copy()
15 pca = PCA()
16 pca.fit(X_pay)
17
18 100*pca.explained_variance_ratio_.round(3)
```

Листинг 4 – Алгоритм выявления аномальных клиентов

```
1 df_y = pd.read_csv('train_target.csv')
2 y = df_y.flag.values[:1000000].copy()
3 y_ = df_y.flag.values[:1000000].copy()
4
5 counter_norm, counter_anom = 0, 0
6 ind_norm, ind_anom = [], []
7 for i,x_pay in enumerate(np.array(X_pay)):
8     if y_[i]==0 and min(x_pay)>0:
9         counter_anom +=1
10        ind_anom.append(i)
11        y_[i]=1
12    else:
13        counter_norm +=1
14        ind_norm.append(i)
```

Листинг 5 – Критерий Колмогорова-Смирнова

```
1 col = 'enc_paym_15'
2 T = X_pay.enc_paym_15
3 x00, x01 = T[(y==0) & (y_== 0)], T[ind_anom]
4 x11, x01 = T[(y==1)], T[ind_anom]
5 ind00 = (x00!=0) & (x00!=1) & (x00!=2) & (x00!=3) & (x00!=4) & (x00
   !=5)
6 ind01 = (x01!=0) & (x01!=1) & (x01!=2) & (x01!=3) & (x01!=4) & (x01
   !=5)
7 ind11 = (x11!=0) & (x11!=1) & (x11!=2) & (x11!=3) & (x11!=4) & (x11
   !=5)
8 x00, x11, x01 = x00[ind00], x11[ind11], x01[ind01]
9
10 def my_cdf(x1, x2, step):
11     tm = max(max(x1), max(x2))
12     F1, F2 = [], []
13     n1, n2 = x1[x1>0].shape[0], x2[x2>0].shape[0]
14     for t in np.arange(0, tm+0.1, step):
15         F1.append(x1[(x1>0)&(x1<t)].shape[0])
16         F2.append(x2[(x2>0)&(x2<t)].shape[0])
17     F1, F2 = np.array(F1)/n1, np.array(F2)/n2
18     K = max(abs(F1-F2))*np.sqrt( n1*n2/(n1+n2))
19     return F1, F2, K, tm
20 step = 0.01
21 F01, F11, K, tm = my_cdf(x01, x11, step)
22 F01_, F00, K_, tm = my_cdf(x01, x00, step)
23 K, K_
```

Листинг 6 – Тест Фишера на равенство коэффициентов корреляции выборочной и генеральной совокупностей

```
1 def transform_f(r, R, n):
2     if n > 3:
3         return 0.5*np.log(((1+r)/(1-r))*((1-R)/(1+R)))*np.sqrt(n-3)
4     else:
5         return 0
6
7 mth = 17
8 n_all, n_nod, n_def, n_ano = X_pay.shape[0], len(ind_nod), len(
9     ind_def), len(ind_ano)
10 R_all, r_nod, r_def, r_ano = all_c[mth], nod_c[mth], def_c[mth],
    ano_c[mth]
transform_f(r_nod, R_all, n_nod), transform_f(r_def, R_all, n_def),
transform_f(r_ano, R_all, n_ano)
```

Приложение 2