

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ «НАЦИОНАЛЬНЫЙ
ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»

ФАКУЛЬТЕТ КОМПЬЮТЕРНЫХ НАУК
ОСНОВНАЯ ОБРАЗОВАТЕЛЬНАЯ ПРОГРАММА
«ПРИКЛАДНАЯ МАТЕМАТИКА И ИНФОРМАТИКА»

КУРСОВАЯ РАБОТА

Программный проект на тему:
«Кредитный скоринг. Сравнение линейных моделей с более сложными
моделями машинного обучения»

Выполнила: студентка группы БПМИ194 3 курса,
Смирнова Анна Романовна

Руководитель КР:
Преподаватель Воробьева Мария Сергеевна

Оглавление

Аннотация.....	3
Ключевые слова	5
Введение	6
Обзор литературы	8
1. Анализ предоставленных данных	9
Разведочный анализ данных.....	11
Обработка данных	12
2. Отбор важных переменных	13
Одномерные методы	13
Многомерные методы	16
3. Построение моделей.....	20
Логистическая регрессия на WOE	21
Логистическая регрессия с кодированием признаков	25
RidgeClassifier	26
Древесные алгоритмы	27
DecisionTreeRegressor.....	28
Градиентный бустинг.....	Ошибка! Закладка не определена.
XGBoost.....	29
CatBoost	31
Выбор лучшей модели	32
Заключение.....	32
Список литературы.....	33

Аннотация

Программный проект на тему: «Нейронные сети на табличных данных на примере задачи кредитного скоринга»

Выполнила: Смирнова А. Р.

Руководитель: Воробьева М. С.

В ходе данной курсовой работы был выполнен анализ существующих методов решения задачи кредитного скоринга. Главной целью являлось реализовать и протестировать существующие методы решения данной задачи, после чего метрически оценить результаты и выбрать лучшую модель. Методы, которые были изучены: Ridge-регрессия, логистическая регрессия, метод опорных векторов, дерево решений, случайный лес.

Abstract

Program course work: «Credit scoring. Linear models vs modern data science models»

Student: Smirnova A. R.

Teacher: Vorobeva M. S.

In the course work, the analysis of the quality of existing methods of solving the problem of credit scoring was carried out. The main objective was to analyze and test existing methods of solving this problem. The conclusion included a detailed description of each method and a comparative analysis of quality. Methods that have been studied: Ridge-regression, logistic regression, support vector machine, decision tree, and random forest.

Ключевые слова

- Кредитный скоринг – это
- Древесные алгоритмы – это
- Логистическая регрессия – это
- Метрика качества – это
- Weight of Evidence (WoE) – это
- Information Value (IV) – это
- Экономический эффект – это

Введение

Задача кредитного скоринга возникает в ситуации, когда банку необходимо принять решение о выдаче или отказе по кредиту, при этом принимая во внимание множество несвязанных факторов. Одним из решений данной ситуации является субъективное заключение кредитного эксперта, однако человеческий фактор не всегда позволяет учесть все входные данные, потому что с течением времени количество факторов, влияющих на принятие решение о выдаче кредита, растет.

Банковская сфера активно расширяется, и каждой компании необходимо внедрять качественные автоматизированные системы, которые помогут эффективно обеспечивать контроль работы бизнес-процессов. Решить этот вопрос активно помогают системы кредитного скоринга, такие как, например, логистическая регрессия, древесные алгоритмы, или другие модели бинарной классификации, с помощью которых можно оценить большой набор признаков и учесть предыдущий опыт работы с кредиторами.

Кредитный скоринг, несмотря на свое название, применяется широко за пределами обычных кредитов. Например, аналогичный анализ проводят для выдачи кредитных карт, ипотеки и других долговых финансовых операциях. От оценки кредитоспособности заемщика так же зависит одобряемая сумма кредита, срок и процентная ставка. Кроме этого, кредитный скоринг используется не только для принятия решения о выдаче кредита, а также в процессе исполнения заемщиком своих обязательств, для того чтобы банк мог изучить поведение клиента до выплаты задолженности, и как можно раньше предвидеть возможный дефолт.

Главной **целью** этой работы является самостоятельно реализовать существующие методы решения задачи кредитного скоринга, после чего составить подробное описание каждого из методов и провести сравнительный анализ качества. **Актуальность** работы подтверждается постоянным совершенствованием банковской системы, и растущей необходимостью работать с большими объемами данных. Численность населения непременно

растет, в связи с чем увеличивается и количество предлагаемых банками услуг, в том числе различного рода кредитов, отличных от привычного понимания. Например, появлением лизинга¹ или популярности микрозаймов², которые так же нуждаются в анализе платежеспособности своих клиентов. Кроме расширения сферы применения, несмотря на популярность задачи, основным методом решения до сих пор является логистическая регрессия с весовыми признаками (WoE).

Основные **задачи**, которые предстоит выполнить в ходе работы:

1. Провести разведочный анализ данных
2. Осуществить обработку данных
3. Реализовать различные методы машинного обучения: в контексте моей работы это различные древесные алгоритмы и вариации логистической регрессии, которые было решено использовать для сравнения.
4. Провести сравнительный анализ полученных результатов используя подходящие метрики качества: F1, ROC AUC, Precision, GINI, Accuracy и другие.
5. Сформулировать предложения о том, можно ли внедрить лучшую модель для решения бизнес-задач и эффективной работы в банковской сфере.

¹ Лизинг — это долгосрочная аренда определенных объектов собственности (оборудование, машины, сооружения) с погашением задолженности в течение нескольких лет.

² Микрозаймы — это кредиты на небольшие денежные суммы, которые обычно даются на короткий период.

Обзор литературы

Сделаю в конце, когда будут все источники в одном месте

1. Анализ предоставленных данных

Предварительная обработка данных необходима перед решением любой задачи машинного обучения. При сборе данных для обучения модели могут быть использованы различные источники, что может послужить наличию выбросов, пропущенных данных, расхождения формата и настоящего значения признака и другого.

Зачастую предобработка данных занимает больше всего времени, однако тщательное и подробное изучение имеющегося набора признаков — залог успешного решения задачи, поскольку более качественные данные помогают получить более совершенную модель. Для того чтобы принять решение, как поступать с той или иной проблемой, существует множество методов. Мы посмотрим на подробное описание каждой из них, и постепенно применим наиболее подходящие к нашему набору данных.

Обзор входных данных

В качестве входных данных использовался набор «Python_Credit_Scoring»³ из соревнования на Kaggle. В нем 26 различных признаков, не включая целевую переменную. Ниже представлены краткие описания каждого из признаков, которые мы рассмотрим подробнее позже:

1. *loan_amnt* — Запрашиваемая у банка сумма для кредита.
2. *funded_amnt* — Общая сумма обязательств по кредиту на данный момент.
3. *funded_amnt_inv* — Общая сумма обязательств инвесторов.
4. *term* — Количество платежей по запрошенной сумме.
5. *int_rate* — Процентная ставка кредита.
6. *installment* — Размер первоначального взноса.
7. *grade* — Оценка кредитного риска.

³ Ссылка на источник с данными: <https://www.kaggle.com/code/aashofteh/python-credit-scoring-ml-elastic-net-regression/data>

8. *emp_title* — Должность, представленная заемщиком при получении кредита.
9. *emp_length* — Трудовой опыт в годах.
10. *home_ownership* — Статус собственности жилья на момент открытия кредита.
11. *annual_inc* — Годовой доход.
12. *verification_status* — Статус верификации.
13. *issue_d* — Месяц, в который получено финансирование.
14. *purpose* — Цель взятия кредита.
15. *addr_state* — Государство, указанное в заявке на получение кредита.
16. *dti* — Процент ежемесячного валового дохода потребителя, который идет на выплату долгов.
17. *delinq_2yrs* — Количество просроченных платежей более чем на 30 дней за последние два года.
18. *earliest_cr_line* — Месяц открытия самой ранней кредитной истории.
19. *inq_last_6mths* — Количество обращений кредитора в бюро кредитных историй за последние 6 месяцев.
20. *open_acc* — Это количество открытых кредитов в данный момент.
21. *revol_bal* — Общий кредитный оборотный остаток.
22. *revol_util* — Доля утилизации кредита
23. *total_acc* — Общее количество активных и закрытых кредитов.
24. *out_prncp* — Оставшаяся непогашенная сумма.
25. *total_pymnt* — Выплаты, полученные на сегодняшний день.
26. *loan_status* — Текущий статус кредита.
27. *risk* — Дефолт (1) или не дефолт (0).

Создание новых признаков

Иногда отдельный признак на первый взгляд может не нести полезной информации, однако, как мы обсудим позже, может быть очень полезен в совокупности с другими, или при создании новых признаков. В нашем наборе данных нет очевидных комбинаций для построения признаков, которые могли бы положительно влиять, но, чтобы попробовать этот метод, я добавила колонку «mean_pauement», которая считает среднемесячный платеж по кредиту, без учета процентов.

Разведочный анализ данных

Exploratory Data Analysis (EDA) – это важнейшая ступень подготовки к построению модели. Этот этап обычно занимает больше всего времени, потому что мы вручную исследуем каждый признак на предмет значимости для решения нашей задачи.

Для рассмотрения категориальных признаков основной задачей является посмотреть на уникальные значения и их интерпретируемость. Например, в изначальном наборе данных у нас есть переменная *loan_status*, которая принимает такие значения: Fully Paid, Charged Off, Default, In Grace Period, Late (16–30 days), Late (31–120 days). Если мы знаем каждое из значений, то можем более точно оценить, какие категориальные переменные подойдут нам лучше. Здесь нам важнее оставить только два из них – Fully Paid, Charged Off, поскольку нас интересует только два сценария.

Для числовых признаков ситуация обстоит иначе: нам часто сложно интерпретировать значения, поэтому важнее будет посмотреть на их математические характеристики. Например, на график распределения, который поможет применить правильные модели отбора переменных в дальнейшем, на «ящик с усами» для проверки наличия выбросов, наличие отрицательных значений, пропусков.

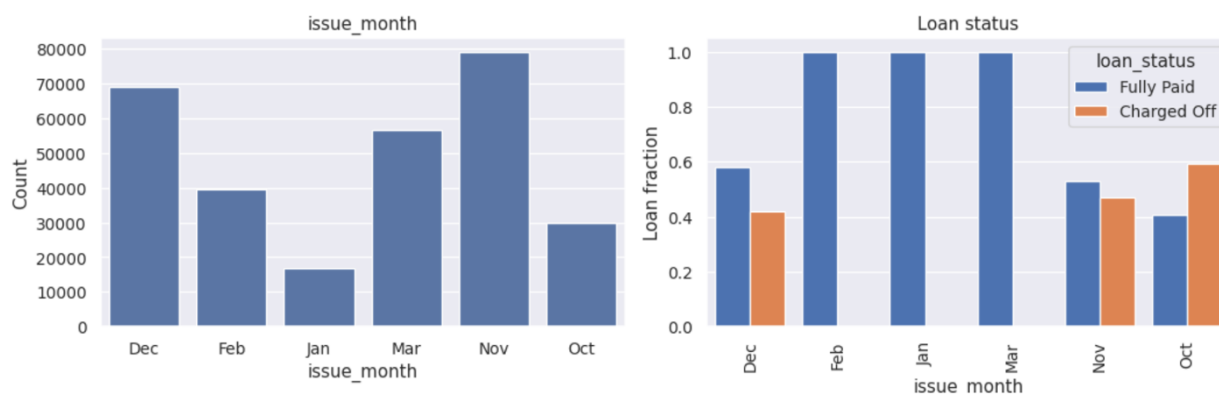


Рисунок 1.1 Необычная зависимость месяца, в который было получено финансирование, и погашения кредита

Обработка данных

Перед тем как переходить к точечным методам, стоит провести общую обработку данных. Она будет включать в себя следующие шаги:

Поскольку в моей задаче достаточное количество таргета обоих значений, я буду удалять любые признаки с более чем 85% пропущенных значений, не боясь потерять ценные данные. Удалению также подлежат признаки, у которых только одно уникальное значение в совокупности с пропущенными данными. Таких колонок в моем дата сете обнаружено не было.

Удалить признаки, значения которых повторяют друг друга, чтобы избежать мультиколлинеарности⁴. В нашем наборе признаков, точно повторяющих другие нет, но есть те, которые включают в себя значения другого признака. Например: *funded_amnt* и *funded_amnt_inv*, распределение которых достаточно близки. В такой ситуации имеет смысл оставить только один признак, для нас это *funded_amnt_inv*, так как он имеет бóльший объем данных. Можно это увидеть так же на рисунке 1.4, где коэффициент Пирсона между этими признаками равен единице.

⁴ Мультиколлинеарность — случай, при котором наблюдается высокая корреляция между признаками. Для проверки на мультиколлинеарность можно посмотреть на коэффициент Пирсона между ними.

2. Отбор важных переменных

После удаления основных пропусков и повторений, стоит рассмотреть каждый признак в отдельности. Данный подход не всегда является оптимальным, поскольку признаков может быть слишком много, но в нашем случае это поможет построить более качественную модель. При рассмотрении каждого признака можно обратиться к нескольким методам.

При рассмотрении каждого признака в отдельности, могут возникнуть сложности при решении, будет нам полезна данная характеристика или нет. В спорных ситуациях важно помнить, что даже если отдельный признак кажется слабо коррелирующим с целевой переменной, в совокупности с другим он может показывать хорошую разделяющую способность. По этой причине, при отборе переменных мы будем рассматривать одномерные и многомерные методы.

Одномерные методы

Одномерными называются методы, отражающие зависимость между конкретным признаком и целевой переменной.

Самый популярный пример данного типа – корреляция. Корреляция – это показатель степени зависимости между двумя переменными. Интервал значения корреляции может быть от -1 до 1 . Положительная корреляция говорит о том, что при увеличении значения одной переменной будет увеличиваться и другая, и наоборот. При работе с данными чаще всего используют корреляции *Спирмена* и *Пирсона*.

$$r_s = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

Рисунок 1.2 Корреляция Спирмена

Чаще всего при построении матрицы корреляции используется корреляция Спирмена. Здесь N – количество элементов в выборке, а значения

коэффициент принимает от $[-1, 1]$, в которых 1 – самая сильная прямая линейная зависимость.

Есть и другие, менее известные коэффициенты. При выборе коэффициента стоит учитывать, к какому типу принадлежит исследуемая переменная. Например, при нормальном распределении переменной стоит смотреть на значение коэффициента Пирсона, поскольку она подходит для непрерывных метрических переменных, а в случае ненормального распределения – Спирмена, как мере, которая работает с признаками, измеренными в ранговой шкале.

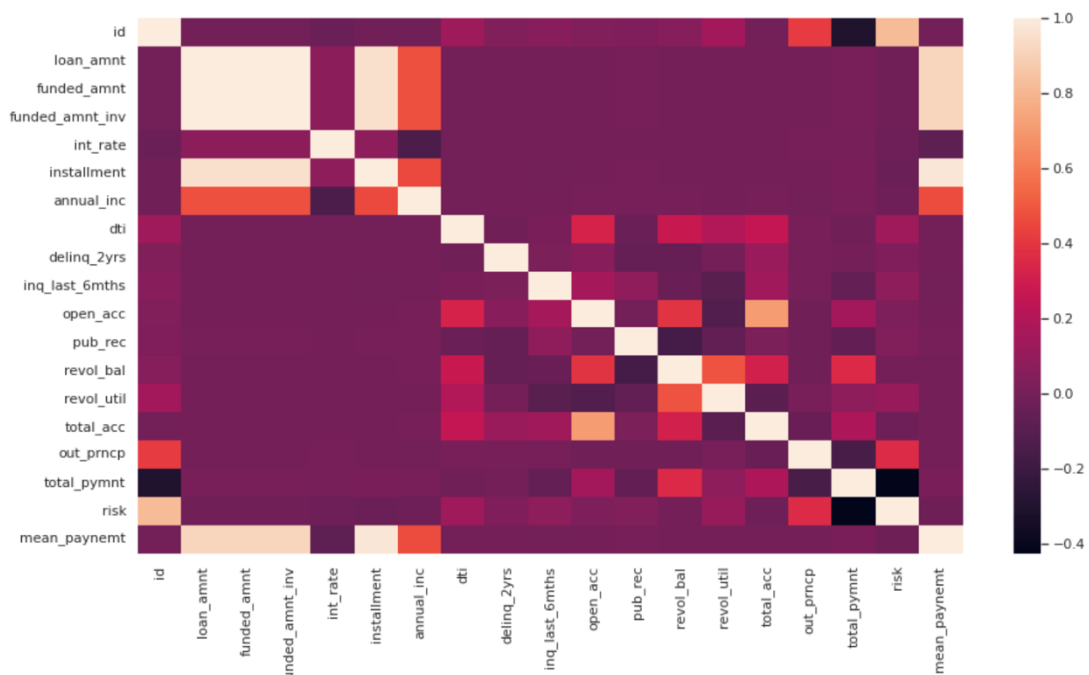


Рисунок 1.3 Матрица корреляции Спирмена.

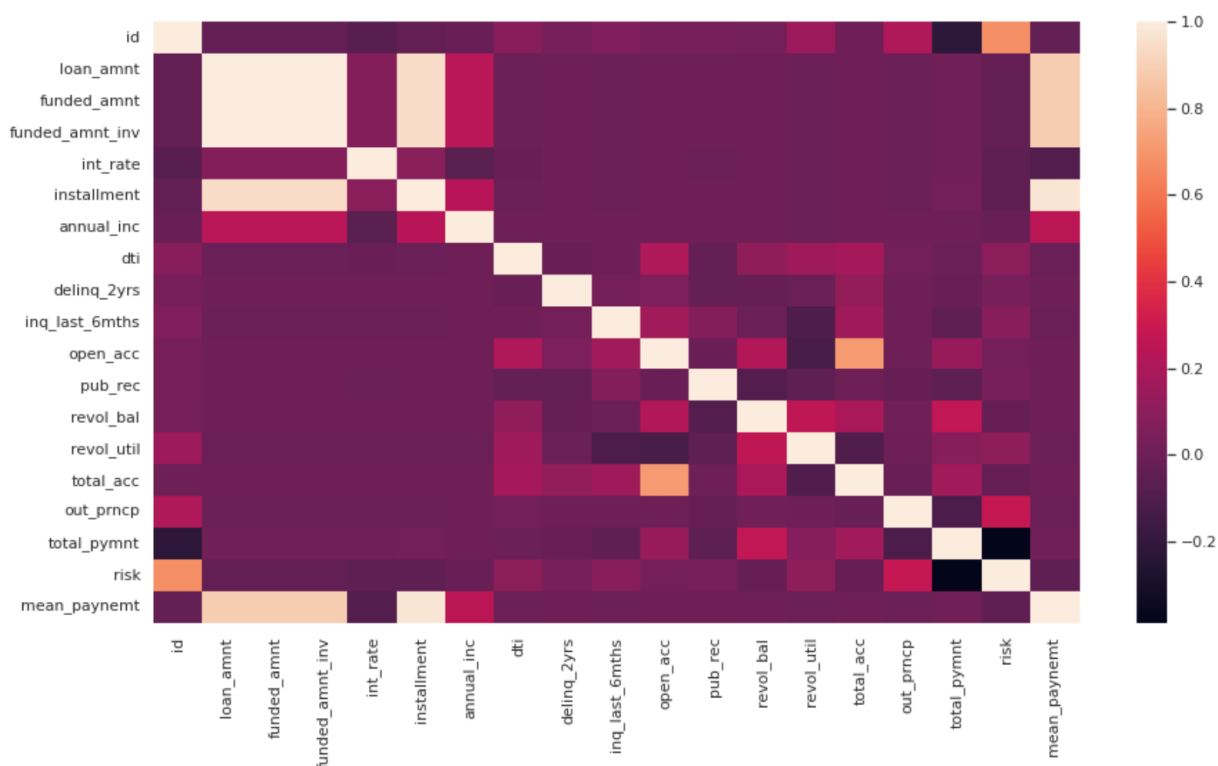


Рисунок 1.4 Матрица корреляции Пирсона.

Другой известной одномерной характеристикой является *Хи-квадрат* – статистический критерий, который используется для анализа связи между двумя категориальными переменными. Данный метод используется с самого зарождения задачи кредитного скоринга. В 1941 г. в «National Bureau of Economic Research»⁵ Дэвид Дюран опубликовал свое исследование, в котором использовал более семи тысяч «хороших» и «плохих» кредитных историй. Тогда хи-квадрат был использован для поиска отличительных черт, которые наиболее остро делили два типа кредитных историй, после чего Дэвидом был разработан индекс эффективности, предназначенный для демонстрации того, насколько эффективна данная характеристика для дифференциации степени риска среди заявителей на кредиты.

$$\chi^2 = \sum \frac{(\text{Observed value} - \text{Expected value})^2}{\text{Expected value}}$$

Рисунок 1.5 Формула критерия Хи-квадрат

⁵ Исследование «Risk Elements in Consumer Instalment Financing»

Для примера можем принять за нулевую гипотезу, что переменные «Оценка кредитного риска» и «Дефолт» не зависят друг от друга. Для этого реализуем тест Chi-Square и получаем Р-значение равным $0.03267 < 0.05$, что говорит о необходимости отвергнуть нулевую гипотезу и принять, что эти признаки друг от друга зависят.

Многомерные методы

Если с одномерными переменными мы имеем четкий и структурированный набор методов, то в случае оценки эффективности нескольких переменных мы чаще пользуемся метриками качества, результаты которых сравниваем при разных наборах переменных.

Создание новых признаков может значительно улучшить качество работы модели, ровно как и наоборот. Ситуаций, в которых мы считаем полученную комбинацию признаков или новый созданный признак удачным, может быть несколько, как например:

- Повышение точности предсказаний. Измеряется с помощью различных метрик, о которых подробнее поговорим позже.
- Более легкая интерпретируемость результатов.

Стоит так же уточнить, что в некоторых ситуациях признаки, не содержащие значительной информации, могут быть полезны для создания новых, более качественных признаков. Это особенно стоит брать во внимание, когда общее количество переменных небольшое и «выбрасывать» из рассмотрения целый столбец может быть неэффективно.

Посмотрим на несколько примеров многомерных методов. Один из самых наглядных – построить *диаграмму рассеяния*. С ее помощью можно наглядно отследить, наблюдается ли корреляция между двумя независимыми признаками, если мы заранее можем предположить сильную связь между ними.

Зависимость между оценкой запрашиваемой суммы и первоначальным взносом



Рис 1.6 Диаграмма рассеяния первоначального взноса и запрашиваемой суммы кредита для разных оценок кредитного риска (А – наивысший уровень кредитоспособности).

Преобразование WOE и Information value

При решении задачи бинарной классификации мы часто имеем дело со множеством категориальных признаков, оценить вклад которых бывает сложно. Для решения этой проблемы разберем две концепции, которые широко применимы при решении задачи кредитного скоринга: «Weight of evidence» и «Information value».

$$WOE = \ln \left(\frac{\text{Event}\%}{\text{Non Event}\%} \right)$$

Рисунок 1.7 Вычисление WOE где Even – дефолт, Non Event – его отсутствие

Вычисление полезности какого-либо признака через WOE преобразование очень близко к вычислению обычной байесовской вероятности, мы смотрим логарифм от соотношения процента «хороших клиентов» к рассматриваемой характеристике к количеству «плохих». Таким

образом, каждое значение признака заменяется размером его «предсказательной силы». После вычисления значения *weight of evidence* у каждого признака, можно попробовать объединить столбцы с близким значением WOE, сильно коррелирующие между собой. Это связано с тем, что эти признаки могут попросту одинаково влиять на целевую переменную, а значит нет необходимости оставлять их обе сепарировано, чтобы не вызвать переобучения.

Некоторые плюсы использования WOE:

- Самостоятельно обрабатывает пропущенные значения
- Работает с выбросами
- Итоговое значение это логарифм распределения, что согласуется со значением логистической регрессии

$$IV = \sum (\text{Event\%} - \text{Non Event\%}) * \ln \left(\frac{\text{Event\%}}{\text{Non Event\%}} \right)$$

Рисунок 1.8 Вычисление IV, где Even – дефолт, Non Event – его отсутствие

Information value – это один из наиболее простых способов оценки значимости признака. После выполнения WOE преобразования можно вычислить «коэффициент полезности», который будет отражать разделительную способность признака, после чего принять решение полезен ли признак в нашей ситуации. Можно заметить, что в формуле IV второй множитель – это как раз вычисленное ранее значение WOE.

Information Value	Variable Predictiveness
Less than 0.02	Not useful for prediction
0.02 to 0.1	Weak predictive Power
0.1 to 0.3	Medium predictive Power
0.3 to 0.5	Strong predictive Power
>0.5	Suspicious Predictive Power

Рисунок 1.8 Согласно «Credit Risk Scorecards»⁶, оценка IV в задаче кредитного скоринга может трактоваться следующими значениями.

Для вычисления обеих характеристик, важно чтобы значение признака было конечным, что не всегда выполняется для численных признаков. В этом используют следующий прием: численный признак разбивается на так называемые «бины», а каждое значение внутри заменяется на соответствующий интервал.

Посмотрим на значения, которые получились после WOE преобразования категориальных и численных признаков на 10 бинах, при этом данные предварительно были вручную обработаны:

categ_feature			IV	feature			IV
0	term	0.003421		0	loan_amnt	0.008603	
1	grade	0.026464		1	funded_amnt	0.008603	
2	emp_length	0.012042		2	int_rate	0.398808	
3	home_ownership	0.000192		3	installment	0.017854	
4	verification_status	0.113406		4	annual_inc	0.006154	
5	purpose	0.020983		5	dti	0.095786	
6	addr_state	0.024436		6	open_acc	0.007700	
7	earliest_cr_line	0.022572		7	revol_bal	0.010455	
8	issue_month	0.469447		8	revol_util	0.078568	
				9	total_acc	0.001475	
				10	mean_payment	0.021894	
				11	earliest_cr_year	0.002613	

Рисунок 1.9 Значения Information values на обработанных данных

⁶ Книга «Credit Risk Scorecards: Developing and Implementing Intelligent Credit Scoring», автор: Naeem Siddiqi

Значения получились совсем низкие, при этом с предсказательной способностью осталось слишком мало признаков. Попробуем посмотреть на совсем «сырых» данных:

			num_feature	IV	
categ_feature	0	term	0.003234		
	1	grade	0.027093		
	2	emp_length	0.011360		
	3	home_ownership	0.000227		
	4	verification_status	0.111073		
	5	issue_d	0.000000		
	6	purpose	0.016257		
	7	addr_state	0.022347		
	8	earliest_cr_line	0.018013		
			0	id	0.000025
			1	loan_amnt	0.015740
			2	funded_amnt	0.015740
			3	funded_amnt_inv	0.015744
			4	int_rate	0.161944
			5	installment	0.013286
			6	annual_inc	0.005194
			7	dti	0.075124
			8	delinq_2yrs	0.004326
			9	inq_last_6mths	0.022072
			10	open_acc	0.003928
			11	pub_rec	0.002155
			12	revol_bal	0.004939
			13	revol_util	0.056900
			14	total_acc	0.002927
			15	out_prncp	0.000000
			16	total_pymnt	1.034497
			17	mean_payment	0.012848

Рисунок 1.10 Значения Information values на необработанных данных

Посмотрим, как оба преобразования данных повлияют на качество моделей.

3. Построение моделей

Логистическая регрессия – это статистическая модель, используемая для решения задачи бинарной классификации, в которой вместо бинарного значения целевой переменной, предсказываются вероятности попадания в положительный класс.

$$\langle w, x_i \rangle = \log\left(\frac{p}{1-p}\right)$$

$$e^{\langle w, x_i \rangle} = \frac{p}{1-p}$$

$$p = \frac{1}{1 + e^{-\langle w, x_i \rangle}}$$

Где соответственно x - вектор со значениями признаков у объекта, а w - вектор весов модели. Обучение модели происходит за счет подбора оптимальных весов w .

Логистическая регрессия на WOE

Наиболее часто применяемым решением в банках является использование логистической регрессии на WOE переменных. Для отбора признаков после преобразования использовались следующие действия:

- Значение IV для признака ≥ 0.015
- Для всех признаков с WOE была посчитана матрица корреляции Спирмина. Для всех пар признаков, сильно коррелирующих между собой (значение корреляции выше 0.8), был оставлен только один с большим значением IV.

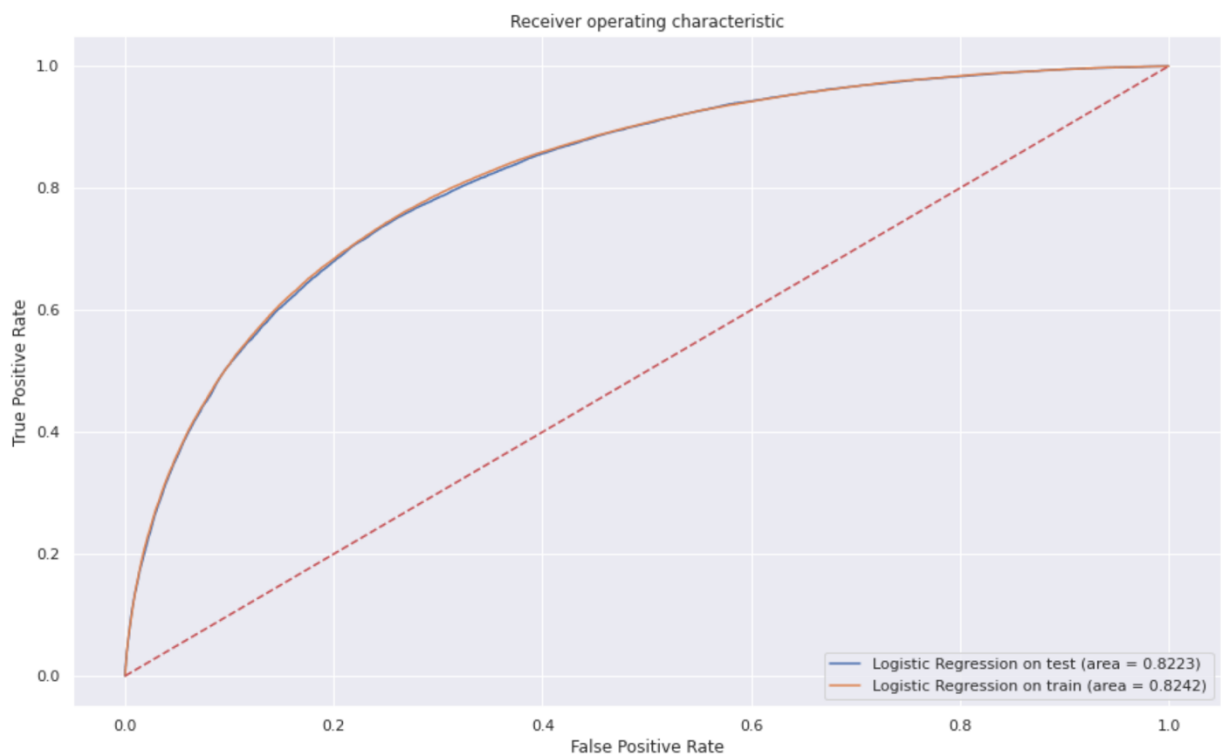


Рисунок 3.1 ROC-AUC кривая для логистической регрессии

После обучения модели можно попробовать улучшить ее качество за счет подбора параметра регуляризации. Для этого я взяла 7 различных значений, и обучив модель с каждым, выбрала наилучшее для обучения на тестовой выборке. Существенного улучшения это не дало: модель стала точнее на 0.01. Все эксперименты были зафиксированы с помощью MLflow и отражены в следующей таблице:

		Metrics		Parameters <			
<input type="checkbox"/>	Run Name	GINI	ROC-AUC	Bin number	Features num	Regularization	Solver
<input type="checkbox"/>	Логрег на WOE с предобработкой	0.74	0.87	10	17	1	liblinear
<input type="checkbox"/>	Логрег на WOE с предобработкой	0.74	0.87	10	17	1	-
<input type="checkbox"/>	Логрег на WOE с предобработкой	0.74	0.87	10	17	-	-
<input type="checkbox"/>	Логрег на WOE без предобработки	0.679	0.84	15	18	0.1	newton-cg
<input type="checkbox"/>	Логрег на WOE без предобработки	0.679	0.84	15	18	0.1	saga
<input type="checkbox"/>	Логрег на WOE без предобработки	0.679	0.84	15	18	0.1	sag
<input type="checkbox"/>	Логрег на WOE без предобработки	0.679	0.84	15	18	0.1	liblinear
<input type="checkbox"/>	Логрег на WOE без предобработки	0.679	0.84	15	18	0.1	-
<input type="checkbox"/>	Логрег на WOE без предобработки	0.675	0.838	20	15	0.1	-
<input type="checkbox"/>	Логрег на WOE без предобработки	0.675	0.838	20	15	-	-
<input type="checkbox"/>	Логрег на WOE без предобработки	0.673	0.837	20	18	-	-

Рисунок 3.2 Промежуточные эксперименты логистической регрессии с WOE преобразованием признаков

Постепенно посмотрим на распространенные подходы к измерению качества моделей, и поговорим подробнее о том, почему были выбраны именно эти метрики для определения качества.

Глобально метрики качества бинарных классификаторов делятся на «пороговые» (single-threshold) и «не зависящие от порога» (threshold-free). Чтобы разделить значения на две категории, необходимо найти «порог», по которому это разделение будет проведено – именно это отличает два этих типа друг от друга. Не зависящие от порога метрики проводят разделение разными способами, которые мы рассмотрим ниже. Считается, что в задаче кредитного скоринга исторически наиболее часто используются именно threshold-free метрики.

Наиболее очевидной метрикой при обучении в задаче классификации является доля правильных ответов – *accuracy*. Она является пороговой, и

сложно интерпретирует полезность полученной модели (особенно на несбалансированных наборах данных), потому что процент положительных ответов не всегда содержит в себе какую-то полезную информацию. Наш набор данных содержит уникальную сбалансированность: практически 50% дефолтных состояний, поэтому ориентироваться на эту метрику не стоит.

Более информативными являются две следующих метрики: *точность* (precision) и *полнота* (recall). Для того чтобы их рассчитать, нам необходимо обратиться к матрице ошибок. Матрица ошибок – это способ разделения всех результатов на четыре категории, в зависимости от предсказания модели и реального ответа: True Positive, False Positive, True Negative, False Negative. Из этих данных точность и полнота высчитываются следующим образом:

$$\text{precision} = \frac{TP}{TP + FP}; \quad \text{recall} = \frac{TP}{TP + FN}.$$

Полнота отражает долю верно выделенных классификатором положительных объектов. Точность отражает долю положительно выделенных объектов, которые действительно являются положительными. Однако, обе метрики стоит особенно осторожно применять на несбалансированных данных. В задаче кредитного скоринга чаще используется комбинация этих метрик, которая создает новый способ оценки ошибки: *F-меру*.

F-мера (или F1) является средним гармоническим точности и полноты. Эта метрика идеально подходит, если мы хотим сравнить точность предсказания используя результаты работы разных моделей. Что отличает эту метрику от трех предыдущих: она полезна даже на несбалансированных наборах данных. Однако, ее сложнее интерпретировать, чем, например, ассигасу. Можно посмотреть на ее значение, но кроме этой метрики использовать и другие.

Чтобы перейти к следующей метрике, нужно ввести еще два понятия: *False Positive Rate* и *True Positive Rate*. Первая отражает долю неверно определенных дефолтных заемщиков, вторая – долю верно определенных. Обе метрики формируются из обозначенной ранее матрицы ошибок:

$$FPR = \frac{FP}{FP + TN}; \quad TPR = \frac{TP}{TP + FN}.$$

Рассмотрим координатную плоскость, где двумя координатами будут FPR и TPR. Значения координат будут соответственно от 0 до 1. Площадь под кривой, полученной из всех комбинаций данных значений, и будет являться данной интегральной метрикой. Таким образом, через количественное значение этой метрики, можно оценить качество классификатора, зная, например что ROC-AUC случайной модели равно 1/2. Эта метрика является одной из самых популярных при решении задачи кредитного скоринга, кроме этого ее можно легко интерпретировать и она подходит даже для несбалансированных данных, поэтому в основном я буду использовать ее.

Коэффициент Джини (Gini Coefficient) — еще одна метрика, которая наиболее часто используется в задачах банковского кредитования. Она используется при оценке качества моделей при решении задач бинарной классификации в условиях сильной несбалансированности классов целевой переменной.

Существует несколько способов посчитать эту метрику, одна из формул напрямую даже линейно связана с предыдущей метрикой ROC-AUC⁷:

$$GINI = 2 * ROC_AUC - 1$$

В задачах кредитного скоринга коэффициент интерпретируется напрямую как качество модели: от 0 до 1, при наилучших результатах. Одним из главных недостатков этой метрики является то, что она не отражает чувствительность модели к разным уровням риска. Однако, даже несмотря на это, коэффициент Джини пока остается одной из самых широко применяемых метрик для оценки модели.

⁷ Доказана зависимость между этими двумя метриками в исследовании «The relationship between Gini terminology and the ROC curve», проведенным Edna Schechtman и Gideon Schechtman

Логистическая регрессия с кодированием признаков

Кроме WOE преобразования, можно попробовать более «точечный» способ предобработки данных. Для того, чтобы вручную не анализировать каждый категориальный признак, я буду использовать различные типы кодирования. Это нужно делать из-за того, что базовые линейные модели часто не могут работать с категориальными переменными. Перед тем как перейти к построению модели, разберем несколько самых известных кодировщиков.

Label Encoder является одним из самых часто применяемых способов кодирования. Принцип очень простой: каждой категории ставится в соответствие уникальное значение. Таким образом, данный кодировщик отлично работает даже с признаками, у которых большое количество разных категорий. Однако у данного способа, есть один существенный недостаток, который стоит учитывать: мы создаем избыточные зависимости в данных. Иногда моделями признаки 1 и 4 могут восприниматься «ближе», чем 1 и 10, хотя реальные значения категорий могут этому не соответствовать. Например, если мы работаем с оценками кредитного риска, и они закодированы как: 1 – «Inadequate», 4 – «Excellent», и 10 – «Adequate».

OneHotEncoding создает новые бинарные признаки в данных. Простыми словами, в колонке нового признака «Grade_ Excellent» будет 1 в случае, если объект принадлежит к этой категории, и 0 иначе. Существенный минус этого подхода в том, что если у нас большое количество категориальных признаков, при этом с большим количеством различных категорий внутри каждого, то полученный набор данных становится просто огромным.

Некой комбинацией двух предыдущих методов служит *Binary Encoder*. Отличает его от OneHot тот факт, что каждая категория кодируется бинарным числом, а не 1 или 0. Таким образом, вместо создания N новых признаков для N категорий, мы получаем $\log(N)$, что существенно меньше при работе с большим количеством категориальных признаков.

			Metrics	Parameters		
<input type="checkbox"/>	↓ Start Time	Duration	F1	Average of F1	Encoding type	Scaling type
<input type="checkbox"/>	✔ 11 seconds ago	8.7s	0.869	weighted	OneHotEncoder	StandardScaler
<input type="checkbox"/>	✔ 36 seconds ago	10.7s	0.589	weighted	OneHotEncoder	MinMaxScaler
<input type="checkbox"/>	✔ 1 minute ago	8.9s	0.596	weighted	OneHotEncoder	StandardScaler
<input type="checkbox"/>	✔ 4 minutes ago	12.0s	0.598	weighted	OneHotEncoder	StandardScaler
<input type="checkbox"/>	✔ 6 minutes ago	10.8s	0.041	binary	OneHotEncoder	StandardScaler
<input type="checkbox"/>	✔ 6 minutes ago	11.8s	0.434	macro	OneHotEncoder	StandardScaler
<input type="checkbox"/>	✔ 7 minutes ago	9.5s	0.434	-	OneHotEncoder	StandardScaler

Рисунок 3.3 Промежуточные эксперименты логистической регрессии с различным наполнением Pipelines

В данном случае кодирование признаков происходило после непосредственной ручной предобработки каждого признака. Наилучшее качество модели по F-мере: 0.869, по ROC_AUC метрике: 0.838, что немного ниже качества, полученного при обучении на WOE переменных.

RidgeClassifier

Кроме классической логистической регрессии, часто используется RidgeClassifier. Это одна из версий алгоритма Ridge Regression, но хорошо работающая с бинарной классификацией. Этот классификатор преобразует целевую переменную в $\{-1, 1\}$, а затем решает ту же задачу, что и классическая Ridge-регрессия. Кроме того, эта модель может быть намного быстрее при работе с большими данными, чем логистическая регрессия, которую мы рассмотрели выше.

Попробуем обучить эту модель аналогично логистической регрессии, а затем подобрать параметр регуляризации и различные типы решений.

<input type="checkbox"/>	Run Name	Metrics		Parameters		
		F1	ROC-AUC	Encoding type	Scaling type	Solver
<input type="checkbox"/>	RidgeClassifier	0.751	0.677	OneHotEncoder	StandardScaler	lbfgs
<input type="checkbox"/>	RidgeClassifier	0.872	0.838	OneHotEncoder	StandardScaler	sag
<input type="checkbox"/>	RidgeClassifier	0.872	0.838	OneHotEncoder	StandardScaler	lsqr
<input type="checkbox"/>	RidgeClassifier	0.71	0.608	OrdinalEncoder	StandardScaler	svd
<input type="checkbox"/>	RidgeClassifier	0.71	0.608	OrdinalEncoder	StandardScaler	-
<input type="checkbox"/>	RidgeClassifier	0.872	0.838	OneHotEncoder	MinMaxScaler	-
<input type="checkbox"/>	RidgeClassifier	0.872	0.838	OneHotEncoder	StandardScaler	-
<input type="checkbox"/>	Логистическая регрессия с кодированием	0.869	0.838	OneHotEncoder	StandardScaler	-
<input type="checkbox"/>	Logreg на WOE с предобработкой	-	0.87	-	-	-

Рисунок 3.4 Промежуточные эксперименты RidgeClassifier

Как видно по результатам, в сравнении с логистической регрессией, получить качество лучше чем у логистической регрессии на WOE переменных не удалось. Качество модели либо сохраняется, либо ведет себя хуже на некоторых видах параметров или типе кодирования.

Древесные алгоритмы

Линейные алгоритмы, которые мы рассмотрели выше, действительно имеют хорошее качество и умеют работать с большими данными. Однако, в задаче скоринга, кроме линейных зависимостей между параметрами (например, между размером кредита и первоначальным взносом, или между оценкой кредитного риска и процентной ставкой) могут быть и другие, более неочевидные зависимости. Для того, чтобы повысить качество предсказаний еще больше, попробуем обучить несколько древесных алгоритмов.

Древесные алгоритмы, или решающие деревья – это новый вид моделей, который может устанавливать нелинейные зависимости между признаками произвольной сложности. Кроме этого, благодаря функциям визуализации, деревья можно будет легко интерпретировать.

DecisionTreeClassifier

Одним из самых базовых алгоритмов решения классификации является DecisionTreeClassifier. В качестве входных данных он использует два набора: массив X с обучающими выборками, и массив Y целочисленных значений целевой переменной. Посмотрим на некоторые параметры:

- *criterion* – это критерий выбора. По умолчанию используется описанный ранее коэффициент GINI.
- *splitter* – это критерий выбора точки разделения. Может быть двух типов: «случайный» или «лучший». «Случайный» выбор оптимальнее использовать, когда у нас очень большой объем данных. В нашем случае подойдет «лучший».
- *max_depth* – это максимальная глубина дерева. Параметр можно оставить по умолчанию, если набор данных небольшой, в противном случае глубину лучше ограничить во избежание переобучения.
- *max_features* – это количество признаков, учитываемых при делении. Если количество признаков небольшое, то можно оставить значение по умолчанию «None», что будет значить, что при разделении учитываются все признаки.
- *class_weight* – этот параметр помогает вручную корректировать веса признаков, в случае если при обучении случается очевидный «перекос». В нашем случае можно попробовать оставить по умолчанию «None» или попробовать «balanced», что автоматически скорректирует веса.

Тут в принципе у меня вообще только два параметра отражены, можно будет позже добавить более подробные эксперименты с разными

			Metrics		Parameters <			
<input type="checkbox"/>	Duration	Run Name	F1	ROC-AUC	Criterion	Encoding type	Max_depth	max_features
<input type="checkbox"/>	11.6s	DecisionTreeClassifier	0.919	0.914	-	OneHotEncoder	10	18
<input type="checkbox"/>	6.3s	DecisionTreeClassifier	0.919	0.914	-	LabelEncoder	10	18
<input type="checkbox"/>	5.9s	DecisionTreeClassifier	0.919	0.914	-	LabelEncoder	10	18
<input type="checkbox"/>	6.0s	DecisionTreeClassifier	0.919	0.914	gini	LabelEncoder	10	18
<input type="checkbox"/>	5.0s	DecisionTreeClassifier	0.815	0.798	gini	LabelEncoder	5	-
<input type="checkbox"/>	6.1s	DecisionTreeClassifier	0.999	0.999	gini	LabelEncoder	20	-
<input type="checkbox"/>	5.8s	DecisionTreeClassifier	0.919	0.914	-	LabelEncoder	10	-
<input type="checkbox"/>	4.7s	DecisionTreeClassifier	0.776	0.769	-	LabelEncoder	-	-

Рисунок 3.5 Промежуточные эксперименты DecisionTreeClassifier

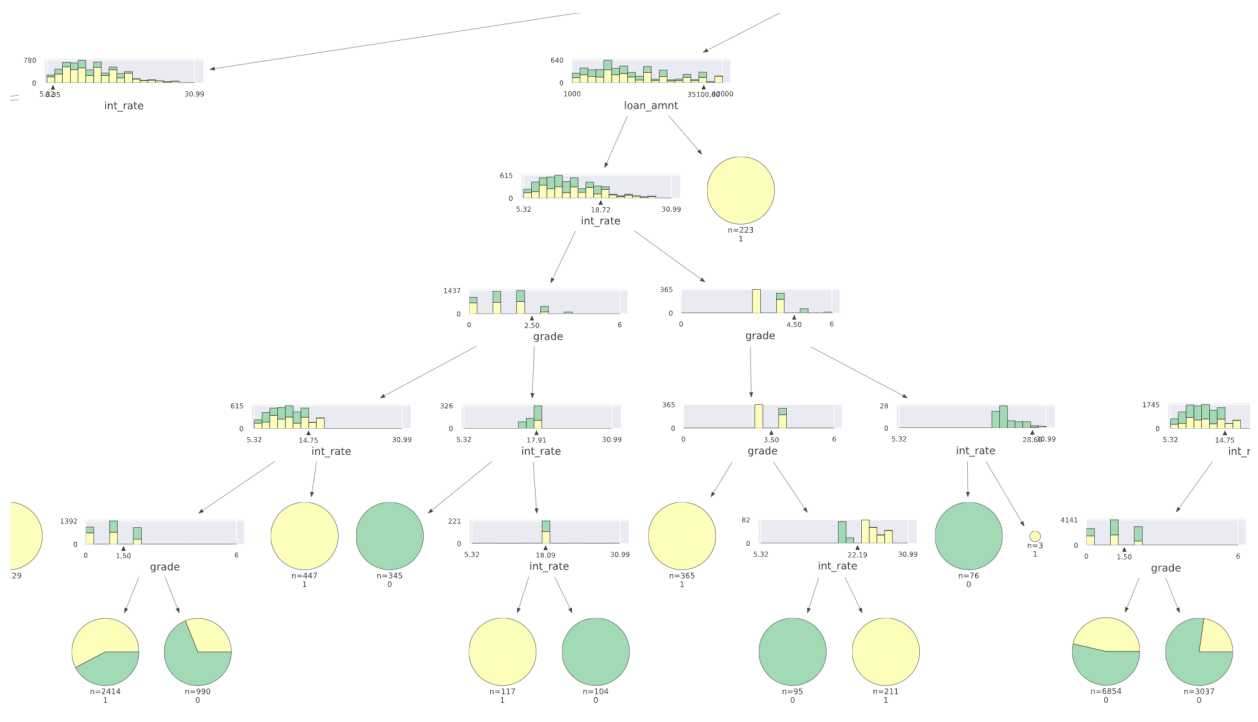


Рисунок 3.6 Часть построенного DecisionTreeClassifier дерева

Нам удалось значительно превзойти качество линейной модели, кроме этого, благодаря визуализации дерева можно наглядно интерпретировать результаты.

XGBoost

Посмотрим на одну из самых успешных реализаций алгоритма градиентного бустинга. В отличие от классического дерева решений, XGBoost использует алгоритмы, строящие деревья связно и постоянно минимизируя

функцию потерь, в то время как алгоритмы как `DecisionTreeClassifier` строят деревья независимо друг от друга. Таким образом, градиентный бустинг позволяет строить более точные модели.

Посмотрим аналогично на основные параметры модели:

- *booster* – это тип нашей модели. Она принимает два значения: «*gbtree*», для использования алгоритмов, основанных на деревьях, и «*gblinear*» для линейных моделей.
- *nthread* – это количество используемых ядер при параллельном обучении. Полезно при большом наборе данных, так как увеличивает скорость модели, но в нашем случае не необходимо.
- *eta* – это скорость обучения, делает модель более устойчивой, с каждым шагом уменьшая вес.
- *max_depth* – это максимальная глубина дерева, с аналогичными свойствами, что и в `DecisionTreeClassifier`.
- *objective* – выбор функции потерь для минимизации. Значение, которое подойдет нам лучше всего: «*binary:logistic*» – логистическая регрессия для бинарной классификации, возвращает предсказанную вероятность.
- *eval_metric* – метрика для проверки данных. Зависит от предыдущего параметра. Возможные значения: *rmse* (среднеквадратическая ошибка), *mae* (средняя абсолютная ошибка), *auc* (площадь под ROC-кривой), *error* (частота ошибки для бинарной классификации).

К сожалению, один из самых главных минусов использования градиентного бустинга является невозможность напрямую интерпретировать результаты модели. В главе с выбором лучшей модели мы посмотрим, как можно применить различные методы интерпретации и посмотреть на значимость отдельных признаков.

Наилучшее качество, которое получилось у этой модели по метрике качества `ROC_AUC`: 0.9453, что значительно превосходит предыдущие алгоритмы.

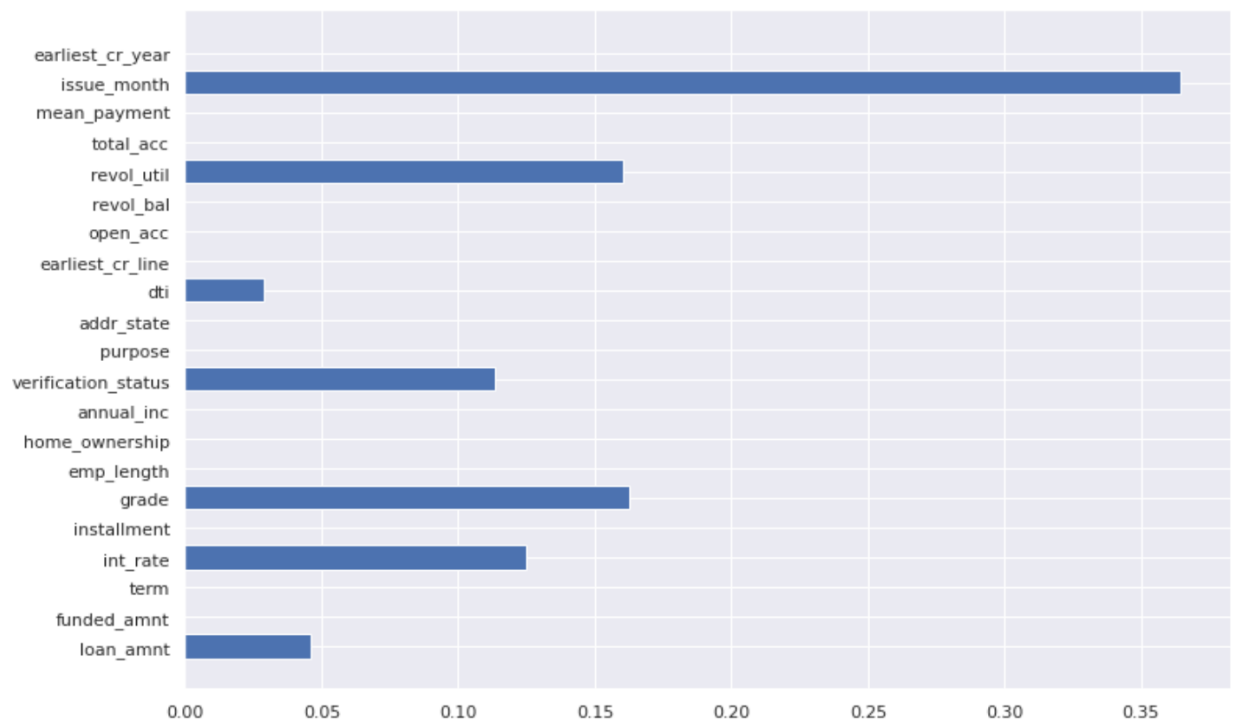


Рисунок 3.7 Индекс важности каждого признака после обучения XGBoost

CatBoost

Выбор лучшей модели

Заключение

Список литературы

- Ahamed, S. (б.д.). *Machine learning algorithms can help us to estimate the risk of a financial decision*. Получено из Towards Data Science: <https://towardsdatascience.com/financial-data-analysis-80ba39149126>
- Complex System Modelling in Engineering Under Industry*. (22 10 2021 г.). Получено из Hindawi: <https://www.hindawi.com/journals/complexity/2021/92222617/>
- Essche, P. O. (2018). *Solvay Brussels School*. Получено из Big Data for Credit Scoring: towards the End of Discrimination on the Credit Market? Evidence from Lending Club.
- Python_Credit Scoring (ML)_Elastic net regression*. (б.д.). Получено из kaggle: <https://www.kaggle.com/code/aashofteh/python-credit-scoring-ml-elastic-net-regression/data>
- Risk Elements in Consumer Instalment Financing*. (б.д.). Получено из National Bureau of Economic Research: <https://www.nber.org/books-and-chapters/risk-elements-consumer-instalment-financing-technical-edition>
- Отбор признаков в задачах машинного обучения*. (б.д.). Получено из Хабр: <https://habr.com/ru/post/550978/>