

## РЕФЕРАТ

Выпускная квалификационная работа содержит 55 страниц, 6 рисунков, 5 таблиц, список использованных источников содержит 15 позиций, 2 приложения.

### МЕТОД МАКСИМАЛЬНОГО ПРАВДОПОДОБИЯ, ПУАССОНОВСКАЯ РЕГРЕССИЯ, ГЕОМЕТРИЧЕСКАЯ РЕГРЕССИЯ, ГЛОБАЛЬНАЯ ОПТИМИЗАЦИЯ, ИНФОРМАЦИОННЫЕ КРИТЕРИИ

В выпускной квалификационной работе показано решение задачи регрессии. Имеется набор данных о случаях сходов с рельсов и крушений грузовых поездов по причине излома боковой рамы. Признак количество подвижных единиц в сходе является целевым. Ключевой особенностью набора данных является малая мощность выборки, а также его разреженность.

Для построения предсказательных моделей использовался метод максимального правдоподобия. Делается предположение о распределении целевого признака (геометрическое, либо Пуассоновское распределение). Целевой признак имеет слабую корреляцию с остальными признаками, по этой причине были введены новые признаки на основе имеющихся.

В работе были рассмотрены модели с различными признаковыми пространствами (8 пространств) и различными функциями связи (7 для Пуассоновской регрессии и 5 для геометрической). Было построено 56 моделей Пуассоновской регрессии и 40 моделей геометрической регрессии.

Скорректированный критерий Акаике  $AIC_c$  – показатель качества. Проведя численные эксперименты, были получены следующие результаты:

- диапазон значений  $AIC_c$  для Пуассоновской регрессии: [356.87, 567.74];
  - диапазон значений  $AIC_c$  для геометрической регрессии: [153.65, 918.43].
- Многие модели геометрической регрессии для заданного набора данных оказалась существенно лучше моделей пуассоновской регрессии.

В приложении к работе был разработан программный комплекс (веб-сервис), реализующий работу с методом максимального правдоподобия. Используются следующие языки и технологии: Python, Java, Spring (Boot, AOP, Web, WebSocket, Data, Security), Log4J, JUnit, Mockito, MyBatis, FlyWay, PostgreSQL, Angular, HTML, CSS, Bootstrap, TypeScript.

## ОГЛАВЛЕНИЕ

	Стр.
<b>ВВЕДЕНИЕ .....</b>	<b>4</b>
<b>ОСНОВНАЯ ЧАСТЬ .....</b>	<b>6</b>
<b>1 ТЕОРЕТИЧЕСКАЯ ЧАСТЬ .....</b>	<b>7</b>
1.1 Постановка задачи машинного обучения .....	7
1.2 Методология решения задач машинного обучения .....	9
1.3 Методы оптимизации .....	12
1.3.1 Градиентный спуск .....	12
1.3.2 Инерционный градиентный спуск .....	14
1.3.3 Стохастический градиентный спуск .....	14
1.3.4 Нормальное уравнение (normal equation) .....	15
1.3.5 Метод имитации отжига .....	16
1.4 Регрессионный анализ .....	16
1.4.1 Проблема переобучения и меры качества .....	17
1.4.2 Информационные критерии .....	19
1.4.3 Линейные модели .....	20
1.4.4 Нелинейные модели .....	24
<b>2 ПРАКТИЧЕСКАЯ ЧАСТЬ .....</b>	<b>28</b>
2.1 Предварительный анализ данных .....	28
2.1.1 Обзор признаков .....	28
2.1.2 Описательные статистики .....	29
2.1.3 Корреляция признаков .....	30
2.1.4 Пропуски в данных .....	32
2.1.5 Экстремальные значения .....	32
2.1.6 Оценка функции вероятности .....	33
2.2 Методы решения задачи .....	34
2.2.1 Признаковые пространства .....	34
2.2.2 Пуассоновская регрессия .....	36
2.2.3 Геометрическая регрессия .....	38
<b>ЗАКЛЮЧЕНИЕ .....</b>	<b>41</b>
<b>СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ .....</b>	<b>43</b>
<b>ПРИЛОЖЕНИЕ А .....</b>	<b>45</b>
<b>ПРИЛОЖЕНИЕ Б .....</b>	<b>50</b>

## ВВЕДЕНИЕ

В странах с большой железнодорожной сетью и большим потоком перемещения поездов таких, как РФ, США, Китай, Индия существует проблема схода составов с рельсов, которые могут быть обусловлены различными факторами, их можно классифицировать на:

- внешние: кривизна пути, профиль пути, состояние транспортного пути, проблемы со стрелочным переводом, погодные условия (при экстремальных температурах рельсы могут сильно расширяться или сжиматься);
- внутренние: количество вагонов в составе, загруженность, скорость, невнимательность машиниста, состояние состава.

Некоторые пути могут проходить через национальные парки, национальные заповедники и другие типы особо охраняемых объектов. По этой причине аварии, произошедшие на таких участках, могут привести к экологической катастрофе, особенно велика опасность, если поезд был грузовым и перевозил легко воспламеняемые грузы (нефть, газ, метан, уголь, древесина) или высокотоксичные грузы. Следует отметить, что помимо экологической проблемы могут возникнуть и другие проблемы, например такие как:

- логистическая - если состав сошел с рельсов, следующим поездам приходится идти в обход, в некоторых случаях обхода может не быть;
- экономическая - связана с издержками транспортной компании по решению экологических проблем, потери части вагонов, локомотива, утрата части груза, временные издержки;
- инфраструктурная - повреждение строения железнодорожного пути, стыков, моста, обрушение тоннеля и др.

В данной работе рассматривается проблема схода состава с рельсов, поскольку проблема является одной из самых опасных. В зависимости от масштаба происшествия сходы классифицируют на аварии и крушения. Согласно [1] за период с 2013 г. по 2016 г. в Российской Федерации имеется 262 протокола сходов с рельсов вагонов как в грузовых поездах, так и в пассажирских, без учета протоколов транспортных происшествий, классифицированных как крушения. Соответственно, в среднем раз в пять с половиной

дней происходит сход с рельсов или крушение грузового поезда, поэтому проблема представляет интерес для железнодорожных компаний.

В данной работе будут построены предсказательные модели числа сошедших вагонов. Для достижения поставленных задач будут использованы методы теории вероятностей и математической статистики, а также методы оптимизации.

## **ОСНОВНАЯ ЧАСТЬ**

## ТЕОРЕТИЧЕСКАЯ ЧАСТЬ

### 1.1 Постановка задачи машинного обучения

Одним из первых, кто дал определение предмету машинного обучения стал американский ученый Артур Самуэль. В 1959 году в своей работе [2], посвященной созданию искусственного интеллекта по игре в шашки с помощью алгоритма *minimax*, Артур Самуэль дал определение тому, что есть машинное обучение – процесс обучения, в результате которого компьютеры способны показывать поведение, которое в них явно не было заложено.

Более современное и точное определение дал Том Митчелл в 1998 году. Корректно поставленная задача обучения определяется следующим образом: говорят, что компьютерная программа обучается на основе опыта  $E$  (experience) по отношению к некоторому классу задач  $T$  (task) и меры качества  $P$  (performance), если качество у задачи из  $T$ , измеренное на основе  $P$ , улучшается с приобретением опыта  $E$ .

Большинство алгоритмов машинного обучения условно можно разбить на несколько классов: обучение с учителем, обучение без учителя, обучение с подкреплением, рекомендательные системы, а также частичное обучение, существуют и другие менее используемые классы алгоритмов.

В алгоритмах обучения с учителем подразумевается обучение на размеченных данных, то есть когда дана матрица, описывающая объекты с помощью признаков (матрица "объекты-признаки") и вектор ответов для каждого объекта. Таким образом методы обучения с учителем можно представлять как функциональную зависимость: на каждый набор признаков  $x \in X$  есть ответы из  $Y$  такие, что  $y : X \rightarrow Y$ , где  $y$  – искомая зависимость.

Рассмотрим такой подход более подробно [8]. Пусть  $X$  – множество объектов,  $Y$  – множество ответов,  $y : X \rightarrow Y$  – неизвестная зависимость, ставящая в соответствие объекты и ответы.

Пусть нам известны:

- $\{x_1, \dots, x_N\} \subset X$ ;

- $y_i, \forall i = \overline{1, N}$  – известное множество результатов.

Ставится задача найти  $a : X \rightarrow Y$  – искомый алгоритм, дающий минимальное расхождение с целевым признаком  $y$ .

Замечание: как правило, множество объектов описывается с помощью признаков. Пусть есть  $p$  признаков, под которыми понимаются следующие отображения:  $f_j : X \rightarrow D_j, \forall j = \overline{1, p}$ . Признаки могут быть: количественными  $D_j = \mathbb{R}$ , бинарными  $D_j = \{A, B\}$ , номинальными  $|D_j| = k < \infty$ , упорядочено номинальными. Один объект может иметь набор признаков разных типов.

Тогда любой объект  $x \in X$  имеет вектор признаков  $(f_1(x), \dots, f_p(x))$ . Следовательно, все объекты можно описать с помощью матрицы ”объекты-признаки”:

$$F = \begin{pmatrix} f_1(x_1) & f_2(x_1) & \dots & f_p(x_1) \\ f_1(x_2) & f_2(x_2) & \dots & f_p(x_2) \\ \dots & \dots & \dots & \dots \\ f_1(x_N) & f_2(x_N) & \dots & f_p(x_N) \end{pmatrix}$$

Каждой строчке соответствуют правильные (ожидаемые) ответы, полученные в результате эксперимента, наблюдения, опроса и т.д. По типу представления множества ответов можно разбить на 3 класса:

- классификация:  $Y = \{A, B\}$ ,  $Y = \{A_1, \dots, A_k\}$ ,  $Y = \{A, B\}^k$ ;
- регрессия:  $Y = \mathbb{R}$  (одномерный случай),  $Y = \mathbb{R}^k$  ( $k$ -мерный случай);
- ранжирование: в отличие от классификации на множестве  $Y$  задано отношение частичного порядка.

Суть прогнозирования заключается в подборе функции  $g : X \times \Theta \rightarrow Y$ , где  $g(x, \theta)$  - функция зависящая от свойств объекта и настраиваемых параметров  $\theta$ .

В обучении без учителя данные об ответах не известны, поэтому можно говорить лишь о том, как данные расположены друг относительно друга. Данные методы машинного обучения происходят без участия экспериментатора и применяются для обнаружения внутренних взаимосвязей. Обычно так решаются задачи кластеризации, понижения размерности, визуализации данных.

Также используются методы обучения с подкреплением, когда набор

данных дается из некоторого потока, а также рекомендательные системы.

Существуют задачи частичного обучения, которые занимают промежуточное положение между задачами регрессии и задачами кластеризации. Задача возникает когда только на части обучающей выборки даны ответы, а другая часть не размечена, обычно такие задачи возникают при больших объемах данных, когда разметить всю выборку либо невозможно, либо очень дорого. При этом задачи частичного обучения не сводятся ни к классификации, ни к кластеризации.

## 1.2 Методология решения задач машинного обучения

Вне зависимости от метода машинного обучения задача состоит из 2-х этапов: обучение и применение. На первой стадии происходит построение оптимального алгоритма  $f$  – функция, дерево, набор инструкций и др. На второй стадии алгоритм выдает ответы для новых объектов.

Оптимальным алгоритмом будем называть такой алгоритм, который на большинстве объектов обучающей выборки дает правильные ответы или достаточно близкие ответы к ожидаемым. Для того чтобы это сделать нужно определять точность или расстояние между объектами, другими словами нужно задать метрику в пространстве объектов. Для этого вводится понятие функции потерь  $\mathcal{L}$  – величина ошибки алгоритма  $f \in A$  на объекте  $x \in X$ :

- $\mathcal{L}(f, x) = \begin{cases} 1, & \text{если } f(x) \neq y(x), \\ 0, & \text{иначе.} \end{cases}$  – индикатор ошибки для случая классификации;
- $\mathcal{L}(f, x) = |f(x) - y(x)|^q$  – для случая регрессии.

На практике для случая регрессии обычно берут  $q = 2$  т.к. при  $q = 1$  возникает проблема с дифференцированием функции потерь. Однако, если в выборке присутствуют большие по модулю выбросы, то можно брать  $q = 1$ , поскольку тогда модель будет меньше реагировать на большие отклонения.

Чтобы в целом оценить алгоритм берут сумму функций потерь по всем объектам из обучающей выборки, деленную на ее мощность, получившуюся величину называют эмпирическим риском (функционал качества алго-



ритма  $f$  на объектах  $X^l$ ):  $Q(f, X^N) = \frac{1}{N} \sum_{i=1}^N \mathfrak{L}(f, x_i)$ .

Таким образом, задача обучения сводится к задаче оптимизации (минимизация функционала эмпирического риска на обучающей выборке):  $\mu(X^N) = \arg \min_{f \in A} Q(f, X^N)$ .

Замечание: иногда сумму делят не на  $N$ , а на  $2 \cdot N$ , с той целью, чтобы при дифференцировании функции эмпирического риска сократились некоторые коэффициенты (для  $q = 2$ ).

Для решения задачи минимизации применяют различные численные методы. Например, метод наименьших квадратов (МНК).

## **CRISP-DM**

Таким образом, любая задача машинного обучения: классификация, регрессия, кластеризация сводится к оптимизационной задаче, возможно с ограничениями. По этой причине были предприняты попытки по созданию единого алгоритма решения задач. Чтобы облегчить процесс решения задач, был разработан и предложен CRISP-DM (CRoss Industry Standard Process for Data Mining) – межотраслевой стандарт решения задач интеллектуального анализа данных. CRISP-DM – модель жизненного цикла исследования данных.

Первая версия данного стандарта была принята в 1999 году. Стандарт призван формализовать схему решения задач анализа данных. Предлагается алгоритм решения произвольной задачи анализа данных в 5 шагов, причем они могут замыкаться в цикле.

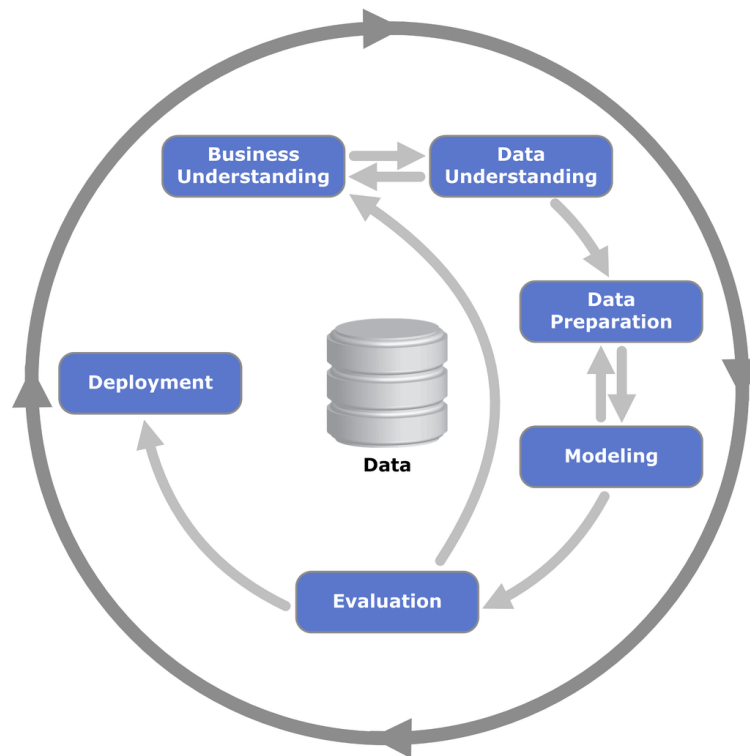


Рисунок 1.1 Жизненный цикл исследования данных [9]

1. Вначале исследователь должен понять предметную область или сферу бизнеса;
2. Далее необходимо понять как собирались данные;
3. На следующем шаге нужно определить есть ли в данных шумы, пропуски, выбросы, все ли признаки несут полезную информацию, можно ли вычислить полезные признаки по уже имеющимся, данный этап можно назвать подготовкой данных;
4. После происходит моделирование или построение предсказательной модели;
5. Полученная модель оценивается с помощью выбранных метрик;
6. Если качество полученной модели удовлетворяет исследователя, модель внедряется в производственные процессы и эксплуатируется.

Помимо CRISP-DM существуют менее известные стандарты: My own, SEMMA и другие. На сайте [6] публикуются результаты опросов по популярности методологий анализа данных.

### 1.3 Методы оптимизации

Во многих задачах науки, экономики и бизнеса возникают проблемы нахождения максимальных/минимальных значений целевой функции. При этом на множество допустимых решений могут быть наложены ограничения. Такие задачи называются задачами оптимизации. Если есть ограничения, то говорят о задаче условной оптимизации, иначе безусловной оптимизации. Рассмотрим некоторые методы решения задач оптимизации.

#### 1.3.1 Градиентный спуск

Метод градиентного спуска является наиболее часто используемым методом ввиду быстроты работы для большинства задач и многообразия модификаций метода (метод наискорейшего спуска, метод сопряженных градиентов, метод Нестерова, метод стохастического спуска и многие другие). Для отыскания экстремальной точки в градиентном спуске используют итеративную формулу:

$$\theta := \theta - \alpha \nabla g(\theta)$$

Где  $\nabla g(\theta)$  – градиент функции  $g(\theta)$ , в векторной форме градиент можно записать как  $\nabla g(\theta) = \frac{\partial g(\theta)}{\partial \theta_1} \vec{i}_1 + \frac{\partial g(\theta)}{\partial \theta_2} \vec{i}_2 + \dots + \frac{\partial g(\theta)}{\partial \theta_n} \vec{i}_n$ , причем  $\vec{i}_1, \vec{i}_2, \dots, \vec{i}_n$  – единичные векторы. Направление градиента указывает на то направление из точки  $\theta$ , которое имеет наибольшую скорость роста функции из данной точки. Соответственно, антиградиент  $-\nabla g(\theta)$  показывает направление наискорейшего убывания;

$\alpha$  – шаг градиента (в машинном обучении называют скоростью обучения). Данный параметр выбирается вручную, при этом, если  $\alpha$  мал, то методу потребуется много итераций для сходимости, если же  $\alpha$  большой, то метод градиентного спуска может начать расходиться, то есть значение  $\theta$  будет отдаляться от точки экстремума. Существует условие Липшица [11], которое дает достаточное условие сходимости градиентного спуска: если  $\exists L : \forall \theta_1, \theta_2 \hookrightarrow \|\nabla g(\theta_1) - \nabla g(\theta_2)\| \leq L \|\theta_1 - \theta_2\|$ , то для  $\forall \alpha < \frac{2}{L}$  гарантируется убывание функции  $g(\theta)$ . На практике обычно выбирают  $\alpha = 10^{-2}$ ,

либо еще меньше в зависимости от того, какая точность необходима.

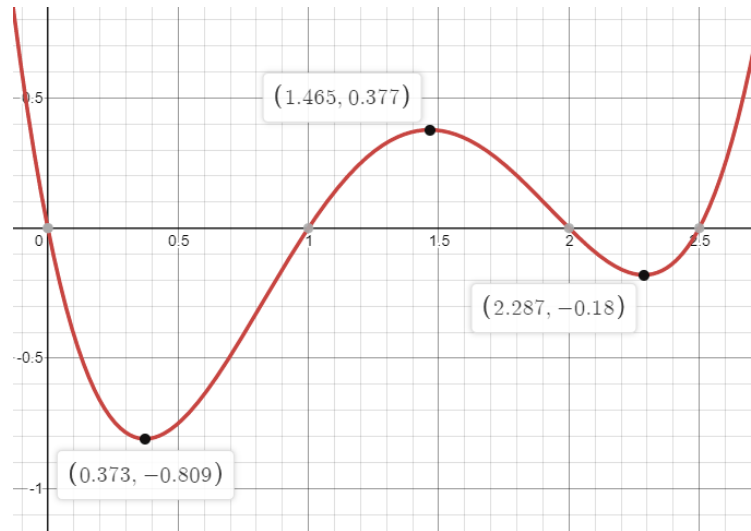


Рисунок 1.2 Проблема много экстремальности функции

Важной особенностью градиентных методов является выбор начального приближения  $\theta_0$ , так, если производится поиск минимума и функция имеет несколько минимумов, как на рисунке 1.2, то выбор  $\theta_0$  будет определять найденный минимум. Например, если  $\theta_0 = 2$ , метод найдет минимум  $g(\theta^*) = -0.18$ , если  $\theta_0 = 1$ , то  $g(\theta^*) = -0.809$ , если же  $\theta_0 = 1.5$ , то метод не сможет найти ни один минимум функции, поскольку в точке  $\theta = 1.5$  производная функции  $g(\theta)$  равна 0, поэтому  $\nabla g(\theta) = 0$  и формула градиентного спуска вырождается в  $\theta := \theta$ . Для того чтобы метод искал глобальный минимум можно запустить данный метод несколько раз из разных начальных точек. Выбрать начальные приближения можно как самому, так и выбрать их случайным образом из равномерного распределения на отрезке  $[a, b]$ . Для проверки того, что был найден глобальный минимум можно запустить метод из точек  $\theta^* - \theta_n$ ,  $\theta^* + \theta_n$ , где  $\theta^*$  – найденный минимум,  $\theta_n$  – точка с большими по модулю координатами. Если метод не сойдется к тому же решению, то необходимо сравнить 2 минимума и выбрать наименьший.

При поиске экстремума в овражных функциях градиентный спуск сходится медленно, причем если число аргументов функции велико, то довольно часто можно встретить овражные области. Для решения данной проблемы используют модификации градиентного спуска, так называемые, овражные методы.

Для ускорения сходимости градиентного спуска применяют различные техники. Можно определять  $\alpha$  на каждом шаге по следующей формуле:  $\alpha_k = \arg \min_{\alpha \in [0, \infty)} g(\theta^k - \alpha g'(\theta_k))$  метод с таким выбором градиентного шага носит название метода наискорейшего спуска, также можно уменьшать параметр  $\alpha$  на каждом шаге. Идея заключается в том, что после каждого шага расстояние до экстремума уменьшается, следовательно, нужно уменьшить градиентный шаг.

### 1.3.2 Инерционный градиентный спуск

Метод известен с середины 20 века и изначально носил название метод тяжелого шарика. Идея данного метода – добавить в формулу градиентного спуска свойство инерционности, то есть чтобы на каждом следующем шаге аргумент  $\theta_{k+1}$  зависел не только от значения антиградиента, но и также от значения предыдущего шага. Для этого предлагается добавить в формулу градиентного спуска следующее слагаемое  $\beta(\theta_k - \theta_{k-1})$ , где  $\beta$  – коэффициент инерции, также как и  $\alpha$  является гиперпараметром, то есть задается вручную. Как правило коэффициент инерции берется немного меньше единицы.

Все семейство инерционных градиентных методов можно описать следующей формулой:

$$\theta_{k+1} = \theta_k - \alpha_k \nabla g(\theta_k) + \beta_k (\theta_k - \theta_{k-1})$$

Добавление нового слагаемого в метод градиентного спуска не меняет асимптотической сложности алгоритма, при этом, если выбрать оптимальные значения в паре  $(\alpha_k, \beta_k)$ , можно добиться ускорения метода на порядок.

### 1.3.3 Стохастический градиентный спуск

Идея метода заключается в использовании вместо градиента функции  $\nabla g(\theta)$  другую функцию (случайный процесс)  $u(\theta, \xi)$  такую, что математическое ожидание  $E[u(\theta, \xi)] = \nabla g(\theta)$ , где  $\xi$  – случайная величина. Метод

стохастического градиента можно описать следующей формулой:

$$\theta_{k+1} = \theta_k - \alpha_k u(\theta_k, \xi_k)$$

Если  $\theta$  – вектор с большим количеством компонент, то вычисление градиента может происходить продолжительное время. Поэтому если использовать метод стохастического градиента и в качестве  $\xi$  взять случайный индекс у  $\theta$ , то есть  $\theta_i$ , где  $i$  – случайная величина, то вычисление градиента сведется к вычислению частной производной по аргументу со случайным индексом. За счет этого происходит существенное ускорение сходимости, это называется процедурой Роббинса-Монро и Кифера–Вулфовица, представленные в 1951 и 1952 годах соответственно [12]. В настоящее время данный метод активно применяется в алгоритмах машинного обучения и в нейронных сетях, где количество признаков у объектов может быть велико. Следует сказать, что данный метод реализован в крупных библиотеках машинного обучения: TensorFlow, PyTorch. Метод используется для обучения моделей с большим и сверх большим объемом данных из-за того, что даже на небольшой подвыборке объектов модель может хорошо обучиться.

#### 1.3.4 Нормальное уравнение (normal equation)

Нормальное уравнение позволяет найти экстремум функционала аналитически, без необходимости применять итеративный подход. Пусть функционал задан следующим образом:  $\mathfrak{J}_\theta(g, y) = (g(\theta) - y)^2$  поставим задачу минимизации этого функционала по всем параметрам  $\theta$ , при условии, что известны  $N$  пар  $(g^{(i)}(\theta), y^{(i)})$ . Составим из  $g^{(i)}(\theta)$  матрицу:  $X = (g^{(1)}(\theta), g^{(2)}(\theta), \dots, g^{(N)}(\theta))$ , тогда  $\theta = (X^T X)^{-1} X^T y$ .

В отличие от градиентных методов в данном методе не нужно проводить итерационную процедуру, а также в методе отсутствуют гиперпараметры. Однако для вычисления ответа нужно вычислить обратную матрицу порядка  $p \times p$ , алгоритм вычисления обратной матрицы работает за  $O(n^3)$  по времени, поэтому для объектов с большим числом признаков, например  $n = 10^6$  потребуется  $C \cdot 10^{18}$  операций, где  $C$  – некоторая константа из алгоритма обратной матрицы. Современные компьютеры могут обрабатывать

порядка  $10^9$  операций в секунду, учитывая нужное количество операций алгоритм не успеет за разумное время вычислить оптимум функции. Поэтому область применимости метода – объекты с небольшим числом параметров ( $< 10^4$ ). На объектах с большим числом признаков хорошо работает стохастический градиентный спуск.

### **1.3.5 Метод имитации отжига**

В основе данного метода глобальной оптимизации лежит физический процесс, который происходит при рекристаллизации вещества. В металлургии под отжигом понимают такую термическую обработку металлов, при которой вначале его нагревают до определенной температуры, выдерживают какое-то количество времени, а затем медленно охлаждают до комнатной температуры. При этом процессе атомы образуют кристаллическую решетку, однако отдельные атомы все еще могут переходить в соседнюю ячейку с некоторой вероятностью. При снижении температуры снижается и средняя кинетическая энергия вещества, а значит снижается и средняя квадратичная скорость движения атомов в веществе, следовательно вероятность перехода атома в соседнюю ячейку кристаллической решетки с уменьшением температуры снижается. Минимальная энергия атомов соответствует устойчивой кристаллической решетке. Поэтому, при отсутствии внешнего воздействия атом переходит в состояние с меньшей энергией, либо остается на месте.

## **1.4 Регрессионный анализ**

Задача регрессионного анализа состоит в определении наиболее оптимальной функции  $f(x)$ , которая бы наилучшим образом восстанавливала закономерность во входных парах данных  $(x_i, y_i)$ . Тогда можно записать следующее соотношение:  $y = f(x, \theta) + \varepsilon$ , где  $f(x)$  – функция, предсказывающая поведение  $y$ ,  $\varepsilon$  – случайная величина с нулевым математическим ожиданием. В большинстве задач полагается, что  $\varepsilon$  имеет нормальное распределение.

Задача исследователя состоит в определении вида функции  $f(x)$ . После происходит параметризация данной функции, выбирается функция штрафа за отклонение от целевого значения, например, сумма квадратов отклонений. Нахождение оптимальных параметров  $\theta$  означает минимизацию функции потерь. Таким образом задача регрессии сводится к задаче оптимизации. Рассмотрим конкретные виды аппроксимирующих моделей.

#### 1.4.1 Проблема переобучения и меры качества

Проблема переобучения или эффект переподгонки – одна из самых частых проблем машинного обучения. Данная проблема возникает, когда модель очень сильно подстроилась под объекты на обучающей выборке и из-за этого потеряла обобщающую способность. Эффект можно обнаружить, если сопоставить функции эмпирического риска, вычисленные для обучающей и тестовой выборок, тогда на значение функционала на обучающей выборке будет близко к нулю, а значение на тестовой выборке будет существенно больше.

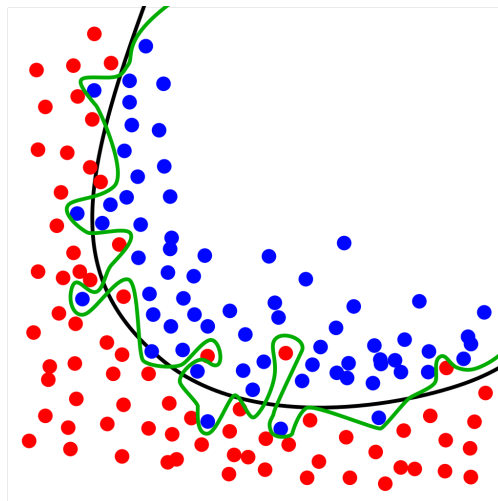


Рисунок 1.3 Пример переобучения модели [14]

На рисунке 1.3 изображены две кривые, соответствующие двум модели, решающие задачу классификации. Видно, что черная разделительная кривая классифицирует данные более общо, чем зеленая кривая. Можно говорить, что модель с зеленой разделительной кривой переобучена. Эффект связан



с тем, что в выборке существуют случайные отклонения, которые могут быть обусловлены ошибкой измерения или погрешностью, при этом, если модель обладает высокой выразительной способностью, она настраивается на этих данных, что приводит к ухудшению качества на тестовой выборке.

Для того чтобы измерять эффект переобучения считают функционалы, измеряющие качество построенной модели на данных на которых модель не обучалась.

Рассмотрим 3 таких метода:

- Отложенная выборка (hold-out) – деление всей выборки на 2 части (обучающая и валидационная). Недостатком данного функционала качества является то, что можно неудачно сделать разбиение так, что оценка будет смещенной, субъективной. Важным параметром данного метода является отношение размеров разбиений, если обучающая выборка мала, оценка качества будет пессимистической, если же валидационная выборка мала, оценка будет неточной. В большинстве случаев обучающую выборку берут в размере 70-80% от общего размера;
- Скользящий контроль (leave-one-out) – по очереди будем выбирать один объект из выборки, на оставшейся части будем обучаться, а на выбранном объекте тестировать. После результаты тестирований усредняются. Таким образом невозможно получить смещенную выборку как в hold-out, однако обучение нужно проводить столько же раз, сколько объектов в обучающей выборке;
- Кросс-валидация (cross-validation) – выберем число групп  $N$ , примерно одинакового размера, равномерно распределим объекты выборки по группам. После будем проводить обучение на  $N - 1$  группе, а тестирование на оставшейся. Результаты усредняются. При выборе малого числа групп оценки будут пессимистичными, но при этом точными. При большом числе групп оценки будут несмещенными, но с большой дисперсией. Обычно выбирают число групп от 5 до 10.

Данные методы основываются на предположении об одинаковости распределения в группах данных и их независимости. Если данные зависимы или их можно разбить на группы так, чтобы в разных группах распределение были разными, возникает риск переобучения.

Также потеря качества может возникнуть при неправильном проведении отбора признаков или понижения размерности. Корректно в начале разбить выборку на обучающую и валидационную, а отбор признаков или понижение размерностей необходимо проводить не над всей выборкой, а только над обучающей. После проводится проверка качества на валидационной выборке. Было бы ошибочно проводить отбор признаков или понижение размерности на всей выборке до разбиения.

Также, если модель регрессии линейная, то в роли меры качества можно считать след оценки ковариационной матрицы  $tr(\hat{K}) = s^2(X^T X)^{-1}$ , где  $s^2 = \frac{e^T e}{N - p - 1}$  – выборочная дисперсия,  $N$  – мощность выборки,  $p$  – количество параметров,  $e$  – вектор невязки,  $X = (x_1, x_2, \dots, x_N)$ . При наличии проблемы переобучения данная мера будет весьма большой.

#### 1.4.2 Информационные критерии

Информационный критерий – еще один способ измерения качества моделей, который учитывает степень переобучения в соответствии с корректировкой на штраф за количество параметров в модели. Таким образом, получается компромисс между сложностью модели и ее точностью. Различные информационные критерии по разному определяют грань этого компромисса.

- Критерий Акаике:  $AIC = 2k - 2l$ , где  $k$  – количество параметров,  $l$  – значение функции логарифмического правдоподобия. При сравнении моделей, имеющие разные мощности выборок  $AIC$  нормируется на размер выборки.
- Для малых выборок используется скорректированный критерий Акаике:  

$$AIC_c = AIC + \frac{2k(k+1)}{N-k-1}.$$
- Критерий Шварца:  $BIC = k \ln(N) - 2l$ .

### 1.4.3 Линейные модели

Иногда, искомую зависимость можно хорошо аппроксимировать линейными моделями, например это могут быть: прогноз стоимости дома или квартиры, классификация типа опухоли, предсказание оттока клиентов, прогноз оклада по описанию вакансии, некоторые физические законы (Гука, Ома, Паскаля).

#### Линейная регрессия

Модель линейной регрессии называется так из-за того, что функция прогноза выглядит как линейная комбинация компонент объекта и компонент вектора параметров. Пусть  $f(x, \theta) = \sum_{i=1}^N x_i \theta_i = \langle \theta, x \rangle$ . Чтобы записать последнее равенство в вектор  $x$  добавляют фиктивную компоненту, равную единице для того, чтобы размерности вектора параметров и вектора объекта совпадали. В качестве функции штрафа можно взять, например  $SSE = \sum_{i=1}^N (y_i - f(x_i, \theta))^2 = \sum_{i=1}^N (y_i - \langle \theta, x_i \rangle)^2$ . Тогда, определив, с помощью методов теории оптимизации, вектор  $\theta$ , будет найдена искомая функция.

В общем случае можно считать что модель  $f(x, \theta)$  задает гиперплоскость в  $n$ -мерном пространстве и чем больше расстояние объекта до гиперплоскости тем больше штраф.

Если среди признаков есть линейно зависимые, то модель теряет свою точность на тестовой выборке, хотя на обучающей выборке можно достичь максимальной точности. Данная проблема называется мультиколлинеарностью линейной модели. При этом возможны 2 случая: функциональная зависимость, когда набор признаков однозначно и точно определяет другой признак и частичную мультиколлинеарность, когда линейная комбинация части признаков сильно коррелирует с другим признаком. При функциональной зависимости вектор параметров определяется неоднозначно, следовательно появляется степень свободы при выборе параметров и можно подобрать их так, чтобы на тестовой выборке они давали завышенные метрики качества. Частичная мультиколлинеарность приводит к неустойчивости оценок.

Существует несколько способов борьбы с мультиколлинеарностью:

- Метод главных компонент – позволяет уменьшить размерность пространства признаков и следовательно избавиться от их коррелированности;
- гребневая регрессия – добавим в функционал эмпирического риска штрафное слагаемое за большие по модулю значения вектора параметров:  $Q = \frac{1}{N} \mathcal{L}(\theta, x_i, y_i, f) + \frac{\tau}{2} \|\theta\|^2$ , где  $\tau$  – коэффициент регуляризации. При минимизации  $Q$  будет минимизироваться как сумма функций потерь, так и квадрат нормы вектора параметров. Подбор параметра  $\tau$  можно сделать вручную, посмотрев на качество метрик на текстовой выборке или определить по критерию скользящего контроля;
- метод LASSO – (Least Absolute Shrinkage and Selection Operator) в отличие от гребневой регрессии, в методе LASSO берется сумма модулей компонент вектора параметров:  $Q = \frac{1}{N} \mathcal{L}(\theta, x_i, y_i, f) + \tau \sum_{i=1}^N |\theta_i|$ , также некоторые значения  $\theta_j$  могут стать в точности нулем, следовательно соответствующий признак больше не будет учитываться. Таким образом в данном методе происходит отбор признаков.

Гребневая регрессия и метод LASSO позволяют ограничить вектор параметров и тем самым избежать проблемы мультиколлинеарности, при этом в методе LASSO можно провести селекцию признаков то есть убрать часть из них.

Говорят, что регуляризация приводит к сокращению размерности пространства, хотя само пространство остается той же размерности, сокращается эффективная размерность, поскольку мы накладываем на вектор параметров ограничение.

## Полиномиальная регрессия

Линейная регрессия требует чтобы между целевой переменной и переменными признаков была линейная зависимость. На практике возникают случаи, когда зависимость между данными невозможно описать линейным образом, для этого рассмотрим модель регрессии с полиномом. Определим модель следующего вида:  $f(x, \theta) = \sum_{i=0}^k \theta_i x^i = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_k x^k$ . Тогда  $f$  задает полиномиальную регрессию. При этом  $x$  может быть как скалярной величиной, так и векторной. Такая модель все еще остается линейной

так как степени при весах остаются первого порядка. Следует отметить, что полиномиальная регрессия может привести к эффекту переобучения, если наибольшая степень  $k$  будет высокой, поскольку выразительная способность (число степеней свободы) модели будет велико, то модель сможет подстроиться под имеющиеся данные и не будет обладать обобщающей способностью.

## Пуассоновская регрессия

Пуассоновская регрессия применяется, когда целевая переменная имеет пуассоновское распределение или в основе которой лежат события, счетчиком которых она является. При этом частота возникновения событий не обязательно стационарна, а может меняться со временем. Пуассоновское распределение имеют, например количество звонков в колл-центр за период времени или число вакцинированных людей за период времени. Для описания частотности событий введем обозначение:  $\lambda$ , которое возможно зависит от времени.

Пуассоновская регрессия основывается на Пуассоновском распределении, которое имеет следующее распределение вероятности:  $p(k) = \frac{e^{-\lambda t} (\lambda t)^k}{k!}$  – вероятность получения  $k$  событий за время  $t$  с интенсивностью потока  $\lambda$ . Математическое ожидание распределения Пуассона численно равно интенсивности потока  $\lambda$ , следовательно для интервала времени  $[0, t]$  можно ожидать  $\lambda t$  событий.

Если целевая переменная определяется распределением Пуассона и нет оснований говорить о нестационарности потока событий, в качестве оценки можно использовать значение интенсивности потока  $\lambda$ .

Если же  $\lambda$  может изменяться от одного события к другому, то предполагается, что  $\lambda$  зависит от объясняющих признаков по некоторому закону. Задача регрессии в таком случае состоит в приближении значений целевой переменной  $y$  к  $\lambda(x)$ . Для этого функция гипотезы  $\lambda(x)$  параметризуется. Как правило [1]  $\lambda(x, \theta) = e^{(x, \theta)}$ , где  $\theta$  – вектор обучаемых параметров,  $x$  – вектор признаков, описывающих объект.

Таким образом вероятность обнаружить  $y_i$  событий для объекта  $x_i$  равна

условной вероятности  $p(y_i; \theta | x_i) = \frac{e^{-\lambda(x_i, \theta)} \lambda^{y_i}(x_i, \theta)}{y_i!}$ .

Для нахождения  $\theta$  воспользуемся методом максимального правдоподобия (ММП). Функция правдоподобия имеет следующий вид:

$$L = L(x_1, \dots, x_N, y_1, \dots, y_N, \theta) = \prod_{i=1}^N p(y_i, x_i | \theta) = \prod_{i=1}^N \frac{e^{-\lambda(x_i, \theta)} \lambda^{y_i}(x_i, \theta)}{y_i!}.$$

Максимизируем полученное выражение по  $\theta$  (переменная  $\lambda$  зависит от  $\theta$ ), для удобства будем максимизировать не  $L$ , а  $l = \ln L$  на значение максимума это не повлияет, поскольку логарифмическая функция монотонно возрастающая. Данное преобразование удобно тем, что при вычислении производной удобнее считать производную суммы, чем производную произведения. По свойству логарифма произведения:

$$l = \sum_{i=1}^N \ln \left( \frac{e^{-\lambda(x_i, \theta)} \lambda^{y_i}(x_i, \theta)}{y_i!} \right) = \sum_{i=1}^N (-\lambda(x_i, \theta) + y_i \ln(\lambda(x_i, \theta)) - \ln(y_i!)).$$

Функция  $l$  называется логарифмической функцией правдоподобия. Найдя максимум этой функции по  $\theta$  одним из численных методов, получим такие значения  $\theta$ , которые при подстановке в исходную функцию правдоподобия максимизируют ее.

В данном методе метрикой качества является значение функции максимального правдоподобия. Та модель лучше, у которой большее значение функции, учитывая, что эффекта переобучения нет, либо его влиянием можно пренебречь.

Когда модель обучилась, то есть подобран вектор параметров для получения прогноза количества событий для объекта  $x_{predict}$  вычисляется значение  $\lambda(x_{predict}, \theta)$ , далее искомая величина  $y_{predict}$  вычисляется по следующей формуле:  $y_{predict} = \arg \max_{k \in \mathbb{Z}^+} p(k | x_{predict}, \theta)$ .

Помимо Пуассоновской регрессии существуют и другие модели, основанные на подсчетах: модель Пуассона с нулевым завышением, отрицательная биномиальная регрессия, обобщенная Пуассоновская регрессия [5].

## Геометрическая регрессия

Геометрическая регрессия применяется, когда целевая функция имеет геометрическое распределение. При помощи такого распределения моделируется число неудач в серии до первого успеха. Функция вероятности имеет вид:  $P(X = n) = (1 - p)^n p$ . В задаче моделирования значения

СВ с натуральными значениями  $n$  функция  $p$  выбирается исследователем, причем на данных из набора данная функция должна принимать значения от нуля до единицы, из-за вероятностных ограничений. Для нахождения  $\theta$  воспользуемся ММП. Функция правдоподобия имеет следующий вид:

$$L = \prod_{i=1}^N (1 - p(x_i, \theta))^{y_i} p(x_i, \theta).$$

Используя численные методы вычислим максимум полученной функции, однако для удобства будем искать максимум не  $L(\theta, X, y)$ , а  $l(\theta|X, y) = \ln L(\theta|X, y)$ . Тогда функция логарифмического правдоподобия будет иметь следующий вид:  $l(\theta, X, y) = \sum_{i=1}^N (y_i \ln(1 - p(x_i, \theta)) + \ln(p(x_i, \theta)))$ .

Как и в модели пуассоновской регрессии метрикой качества является значение функции максимального правдоподобия. Та модель лучше, у которой большее значение функции, учитывая, что эффекта переобучения нет, либо его влиянием можно пренебречь.

Как и в пуассоновской регрессии для формирования прогноза используется следующая формула:  $y_{predict} = \arg \max_{k \in \mathbb{Z}^+} p(k|x_{predict}, \theta)$ .

#### 1.4.4 Нелинейные модели

Нелинейные модели регрессии – такие методы регрессионного анализа, в котором в моделирующей функции параметры модели, ошибки входят нелинейным образом.

#### Решающие деревья

Данный метод исходно разрабатывался для решения задач классификации, однако в настоящее время существует алгоритм CART, расширяющий класс решаемых задач методом решающих деревьев на случай регрессии, поэтому данный метод помещен в раздел "Регрессионный анализ".

Решающие деревья – попытка формализовать человеческое мышление при принятии решений. Решающее дерево можно приближенно проиллюстрировать на примере работы врача: врач задает пациенту уточняющие вопросы, исходя из его ответов происходит спуск по дереву на уровень ниже. Вершина, в которую нужно перейти определяется ответом на вопрос.

Если у вершины нет дочерних вершин, она листовая в данном примере листовые вершины – диагноз больного или совет. Таким образом за конечное количество вопросов можно дойти от корневой вершины до листовой.

Все вершины дерева можно разделить на листовые и внутренние. Любая внутренняя вершина  $v$  содержит предикат  $\beta_v : X \rightarrow A$ , для случая бинарного решающего дерева  $A = \{0, 1\}$ . Любая листовая вершина  $v$  содержит метку класса  $c_v$ .

Для построения решающего дерева используется алгоритм Induction of Decision Tree (ID3). На вход алгоритм может принимать часть выборки  $U$ .

```
def LearnID3(U):
    if все объекты из U лежат в одном классе:
        return (новый лист v, c_v = c)
    найти предикат с максимальной информативностью beta = argmax I(beta, U)
    разбиваем выборку на U0 и U1 по предикату beta
    if len(U0) == 0 or len(U1) == 0: # не смогли найти информативный предикат
        return (новый лист v, c_v = мажоритарный_класс(U))
    создать новую вершину v
    построить левое дерево: L_v = LearnID3(U0)
    построить правое дерево: R_v = LearnID3(U1)
    return v
```

Рассмотрим случай, когда при разбиении входной выборки  $U$  по найденному предикату  $\beta$  один из классов  $U_0$  или  $U_1$  оказывается пустым. В этом случае критерий информативности  $I$  не смог разделить выборку на 2 класса, хотя в выборке были представители как одного, так и другого классов. В этом случае образуем вершину по данному предикату и отнесем ее к тому классу, которых больше в выборке. При этом в этой вершине будут ошибки.

Для поиска предиката с максимальной информативностью используются различные критерии ветвления:

- критерий Джини:  $I(\beta, X) = |\{(x_i, x_j) : y_i \neq y_j \ \& \ \beta(x_i) = \beta(x_j)\}|$  Из определения следует, что критерий Джини позволяет определить число вершин из одного класса, которые были классифицированы предикатом одинаково. При нормировании критерия на число вершин критерий будет показывать частоту объединения вершин из одинаковых классов.



- D-критерий Донского:  $I(\beta, X) = |\{(x_i, x_j) : y_i \neq y_j \ \& \ \beta(x_i) \neq \beta(x_j)\}|$  позволяет определить число вершин из разных классов, которые были классифицированы предикатом в разные ветви дерева. При нормировании критерия на число вершин критерий будет показывать частоту разделения вершин из разных классов.

Для оптимальной обработки пропусков в данных на стадии обучения для каждой вершины считают частоты прохождения объекта в левую и правую ветви. На этапе классификации, если для объекта невозможно вычислить предикат, он отправляется в оба поддерева с определенными весами, после ответ усредняется по ним и выбирается наиболее вероятный класс.

Для избежания эффекта переобучения и улучшения качества алгоритма можно ограничить максимальную глубину дерева, процедура называется *pruning* и реализована в алгоритме C4.5, предложенным Джоном Квинланом.

Рассмотрим обобщение алгоритма ID3 для случая регрессии, для этого рассмотрим метод CART (Classification And Regression Trees).

Будем считать, что  $Y = \mathbb{R}$ ,  $c_v = \mathbb{R}$ . Пусть  $U_v$  – множество объектов, дошедших до вершины  $v$ . Тогда терминальные вершины  $c_v$  определим как среднее значение по всем достигшим вершины  $v$  объектам:  $c_v = \frac{1}{|U_v|} \sum_{x_i \in U_v} y_i$ .

В качестве критерия ветвления возьмем среднеквадратичную ошибку:  $I(\beta, U_v) = \sum_{x_i \in U_v} (\hat{y}_i(\beta) - y_i)^2$ , где  $\hat{y}_i(\beta) = \beta(x_i)\hat{y}_i(U_{v1}) + (1 - \beta(x_i))\hat{y}_i(U_{v0})$  – прогноз  $\beta$  и разбиения  $U = U_{v0} \sqcup U_{v1}$ .

При обнаружении эффекта переобучения можно воспользоваться методом MCMP (Minimal Cost-Complexity Pruning). Идея метода заключается в регуляризации количества терминальных вершин, для этого составляется критерий из двух частей: среднеквадратичная ошибка и  $\alpha|V_{terminal}|$ , где  $|V_{terminal}|$  – количество терминальных вершин,  $\alpha$  – коэффициент регуляризации. При увеличении  $\alpha$  дерево решений упрощается.

## Градиентный бустинг

Если качество базовых алгоритмов машинного обучения для решения задачи не удовлетворяет исследователя, можно воспользоваться композицией

данных алгоритмов, качество которой будет выше отдельно взятого алгоритма данной композиции. Градиентный бустинг – один из таких методов, объединяющий в себе несколько алгоритмов. Алгоритмы в градиентном бустинге обучаются независимо друг от друга. Каждому алгоритму  $a_i(x)$  сопоставлен коэффициент  $w_i$ . Пусть всего алгоритмов  $T$ , тогда линейной композицией будем называть взвешенную сумму:  $a(x) = C \left( \sum_{i=1}^T w_i a_i(x) \right)$ , где  $C : \mathbb{R} \rightarrow Y$  – решающее правило.

Решающее правило необходимо для обобщения градиентного бустинга на задачи классификации. Так, если решается задача регрессии, то решающее правило не учитывается  $C(x) = x$ , для задачи классификации на 2 класса  $C(x) = \text{sign}(x)$ .

Выберем произвольную функцию потерь  $\mathcal{L}(a(x), y)$  и пусть  $T - 1$  алгоритм уже обучен и соответствующие им коэффициенты  $w_i$  определены, тогда запишем функционал эмпирического риска как функцию от нового алгоритма  $a$  и веса  $w$ :  $Q(w, a) = \sum_{i=1}^N \mathcal{L} \left( \sum_{t=1}^{T-1} w_t a_t(x_i) + w a(x_i), y_i \right) \rightarrow \min_{w, a}$ .

Для решения такой задачи оптимизации вначале находится базовый алгоритм  $a(x)$ , после решается задача одномерной оптимизации поиска  $w$ . После решения задачи функционал эмпирического риска обновляется с учетом новой пары  $(a(x), w)$ . Это и есть алгоритм градиентного бустинга.

## ПРАКТИЧЕСКАЯ ЧАСТЬ

### 2.1 Предварительный анализ данных

Анализ данных будет проведен при помощи языка программирования Python3. Данный язык был выбран по нескольким причинам:

- большое количество модулей для анализа данных;
- подробная документация как языка, так и его модулей;
- удобство и простота работы с данными в форматах csv, xlsx;
- множество встроенных функций и выразительность языка.

#### 2.1.1 Обзор признаков

В имеющемся наборе данных о случаях сходов с рельсов и крушений грузовых поездов наличествуют следующие свойства у объектов:

- дата – дата происшествия;
- кол-во вагонов – кол-во вагонов в составе (без локомотивов);
- максимальное число вагонов в сходе – кол-во вагонов с конца до сошедшего поезда;
- общее кол-во вагонов – кол-во вагонов в составе (с локомотивами);
- кол-во сошедших вагонов – целевая переменная; кол-во вагонов, сошедших с рельс;
- скорость – средняя скорость состава;
- вес – масса состава, включая груз, выраженная в тоннах;
- загрузка – отношение текущего веса к максимальному весу;
- стрелочный перевод – наличие стрелочного перевода в месте схода;
- кривизна – обратная величина к радиусу кривизны;
- профиль пути – знак величины определяет направление (положительный, если подъем и отрицательный, если спуск);
- Режим движения – дискретная величина, обозначающая тягу, выбег, торможение;

Вывод первых пяти записей из набора, таблица разделена на две части:

Таблица 2.1 Первые 5 записей в наборе данных

№	Дата	Кол-во вагонов	Макс. число вагонов в сходе	Общее кол-во вагонов	Кол-во подвиж. ед. в сходе
1	2013-01-08	56.0	19.0	58.0	1
2	2013-01-09	60.0	25.0	62.0	1
3	2013-01-10	60.0	4.0	64.0	1
4	2013-01-12	66.0	63.0	68.0	21
5	2013-01-19	67.0	34.0	69.0	1

Продолжение таблицы 2.1

№	Скорость	Вес	Загрузка	Стрелочный перевод	Кривизна	Профиль пути	Режим движения
1	57.0	3402.0	0.547101	0	0.000000	0.0007	NaN
2	72.0	4082.0	0.652657	0	0.000000	0.0009	NaN
3	15.0	4420.0	0.734300	0	0.001639	NaN	3.0
4	67.0	5699.0	0.918094	0	0.002326	0.0060	NaN
5	69.0	5854.0	0.932944	0	0.000000	0.0006	2.0

## 2.1.2 Описательные статистики

Мощность выборки равна 56, при этом в записях содержится существенное количество пропусков.

Вывод описательных статистик для набора данных (примечание: признак 'Дата' исключен):

Таблица 2.2 Описательные статистики

	Кол-во вагонов	Макс. число вагонов в сходе	Общее кол-во вагонов	Кол-во сошедших вагонов	Скорость
count	54	51	54	56	53
mean	63.87	37.13	66.41	3.87	49.15
std	9.79	21.54	10.05	6.08	18.45
min	24	2	26	1	9
25%	60	17.5	62.5	1	35
50%	66	43	68	1	51
75%	68	56.5	71.75	2.25	64
max	96	72	100	26	78

## Продолжение таблицы 2.2

	Вес	Загрузка	Стрелочный перевод	Кривизна	Профиль пути	Режим движения
count	54	54	56	46	44	33
mean	5126.63	0.81	0.1	0.0008	-0.0003	1.67
std	1438.74	0.24	0.31	0.001	0.006	0.78
min	998	0.18	0	0	-0.01	1
25%	4155.5	0.69	0	0	-0.005	1
50%	5722	0.93	0	0	0	1
75%	6010.25	0.99	0	0.001	0.002	2
max	8806	1.07	1	0.01	0.01	3

Примечание к таблице:

- count – количество не пустых значений у свойства;
- mean – среднее значение свойства по всем объектам;
- std – стандартное отклонение свойства по всем объектам;
- min (max) – минимальное (максимальное) значение свойства по всем объектам;
- 25%, 50%, 75% – квантили соответствующих уровней.

Следует отметить, что такие признаки как: 'Режим движения', 'Кривизна', 'Профиль пути' имеют наибольшее количество пропусков в данных: 41%, 17%, 21% пропусков соответственно.

### 2.1.3 Корреляция признаков

Построим матрицу корреляции признаков:

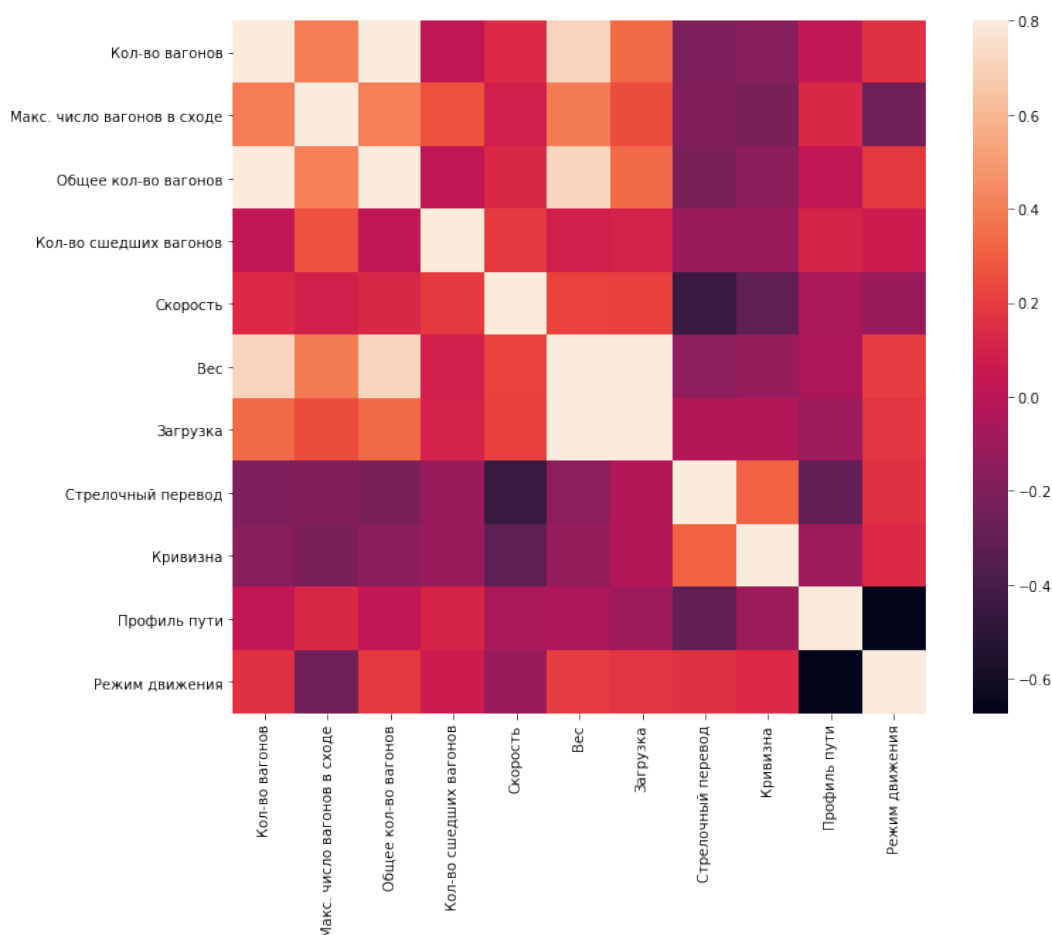


Рисунок 2.1 Корреляция признаков

Из рисунка 2.1 видно, что признаки 'Количество вагонов' и 'Общее число вагонов' имеют сильную корреляцию, что можно объяснить тем, что общее число вагонов равно количеству вагонов плюс количество локомотивов, при этом количество локомотивов в составах как правило равняется двум, следовательно получается функциональная зависимость. Также можно заметить, что признаки 'Вес' и 'Загрузка' сильно коррелируют. Менее сильная корреляция наблюдается у признаков 'Вес' и 'Общее число вагонов'. Заметим, что у 'Профиль пути' и 'Режим движения' наблюдается сильная обратная корреляция. Многие зависимости можно нетрудно объяснить: чем больше вагонов в составе, тем больше вес, чем больший вес, тем, как правило, большая загруженность. Таким образом, можно прийти к выводу, что в данные в наборе избыточны, поскольку несколько признаков несут одинаковое количество информации. Поэтому эти зависимости приводят к проблеме мультиколлинеарности, что приведет к эффекту переобучения в

линейных моделях. Для решения данной проблемы нужно исключить избыточные признаки.

#### **2.1.4 Пропуски в данных**

Из таблицы 2.2 по строке 'count' видно, что в последних трёх признаках присутствуют пропуски в данных. Для решения проблемы пропусков в данных существует ряд методов:

- удалить все записи из используемого в модели подмножества признаков, в которых есть хотя бы одно пустое поле. При использовании этого метода для данного набора данных существует риск того, что оставшегося множества записей не хватит для получения приемлемого качества построенной модели;
- заменить пропуски на средние значения по признаку;
- заменить пропуски на медианные значения по признаку. В отличие от среднего значения замена на медианное позволяет избежать сильного влияния выбросов на итоговое значение;
- заменить пропуски на крайне большие значения, если в качестве модели так или иначе используется структура данных дерево. Таким образом, все записи, имеющие пропуски будут выделены в отдельную ветку дерева;
- заменить пропуски нулевыми значениями, таким образом для логистической регрессии пропуски в данных не будут влиять на предсказанные значения.

При решении данной задачи будет использован метод удаления записей с пропусками. Замена значений на медианное, среднее для данного набора данных может оказаться не оптимальным выбором, поскольку мощность выборки не велика.

#### **2.1.5 Экстремальные значения**

Для поиска выбросов построим графики, изображающие отношения между парами признаков, и гистограмм распределений.

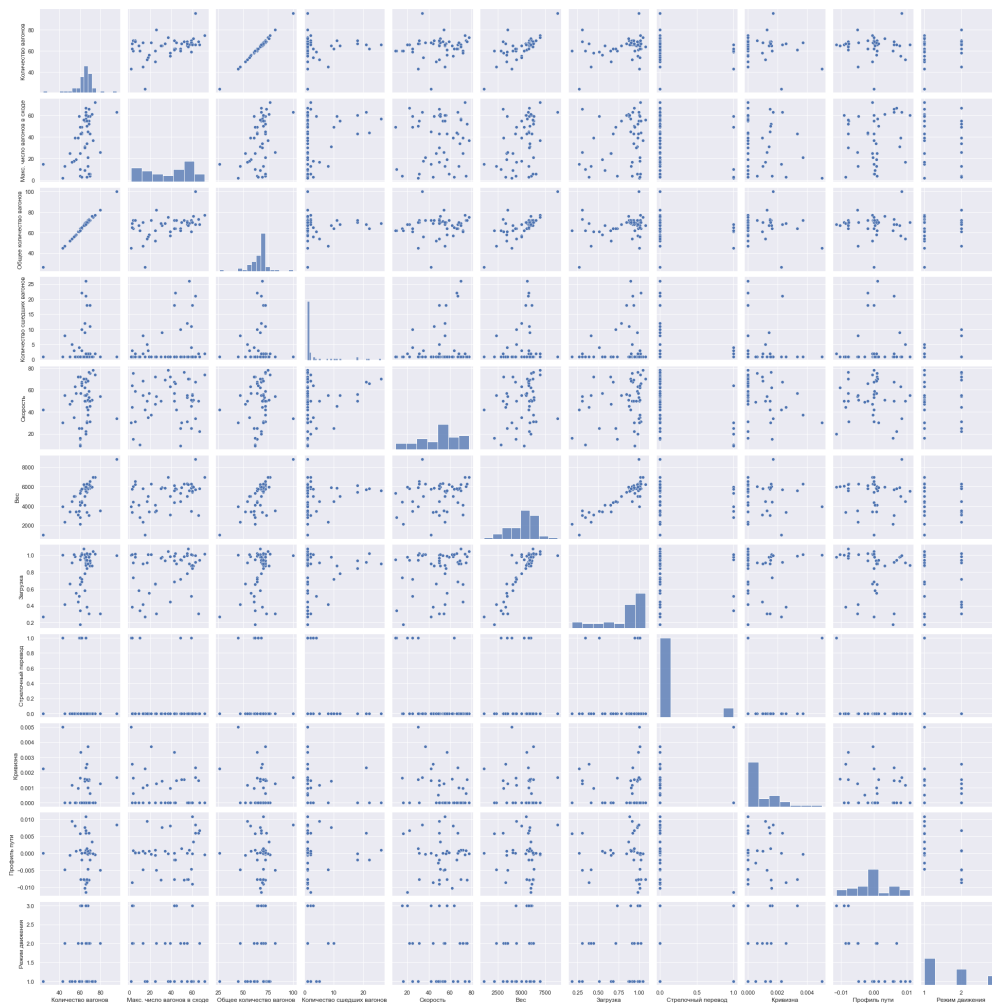


Рисунок 2.2 Пары признаков

Изучив таблицу с описательными статистиками 2.2, а также при детальном рассмотрении графиков пар признаков и распределений на рисунке 2.2 сильных выбросов в данных обнаружить не удалось.

### 2.1.6 Оценка функции вероятности

Построим частотную оценку функции вероятности  $P(\xi = n)$  признака количества сошедших с рельсов подвижных единиц грузового поезда в виде гистограммы, поскольку количество вагонов – дискретная величина.



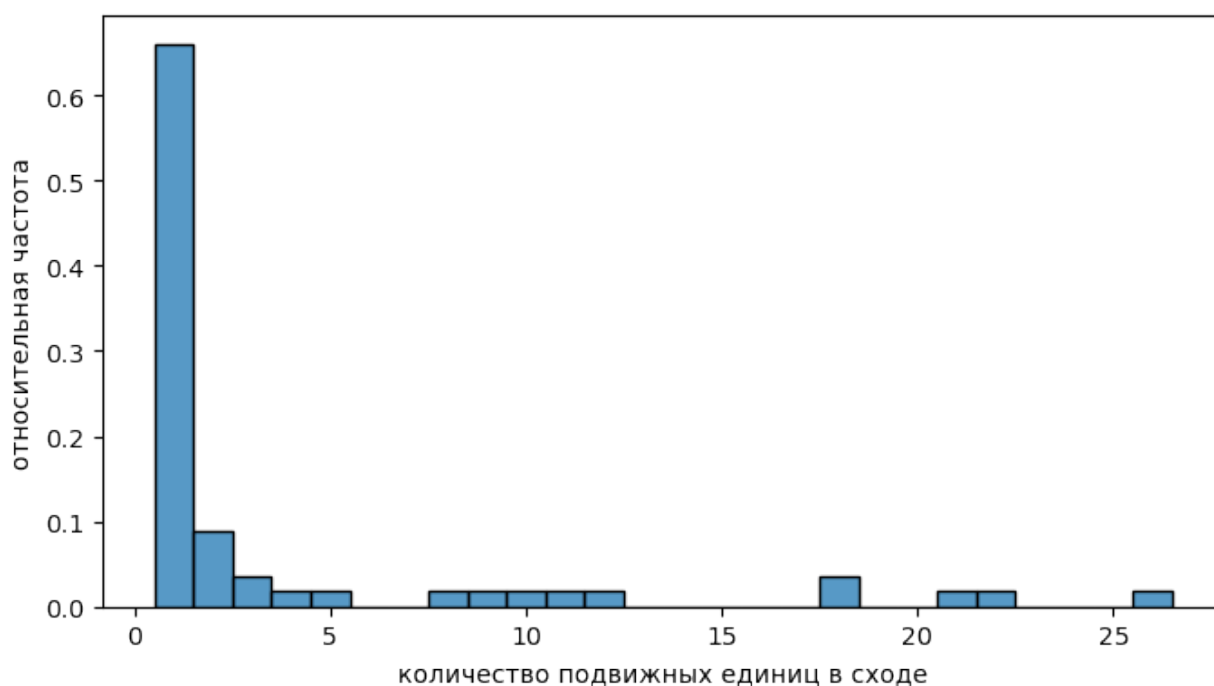


Рисунок 2.3 Диаграмма относительных частот

Из графика видно, что намного чаще сходит небольшое количество подвижных единиц.

Пусть  $C$  – случайная величина, характеризующая количество подвижных единиц в сходе. Введем обозначение  $\tilde{C} = C - 1$ . Будем предполагать, что  $\tilde{C}|x$  имеет распределение Пуассона, либо геометрическое распределение.

## 2.2 Методы решения задачи

В данной работе была написана программная реализация метода максимального правдоподобия, которая описана в приложении А. Также был написан программный комплекс, представляющий собой веб-сервис по работе с методом максимального правдоподобия, описан в приложении Б.

### 2.2.1 Признаковые пространства

Сформируем признаковые пространства для моделей, основанных как на Пуассоновской регрессии, так и на геометрической регрессии:

1.  $features_1$  : (кривизна);
2.  $features_2$  : (кривизна, профиль пути);
3.  $features_3$  : (кривизна, профиль пути · макс. число вагонов в сходе);
4.  $features_4$  : (кривизна,  $1 - \frac{\text{макс. число вагонов в сходе}}{\text{общее кол-во вагонов}}$ );
5.  $features_5$  : (кривизна, профиль пути, скорость · загрузка);
6.  $features_6$  : (кривизна, профиль пути, скорость · загрузка,  $1 - \frac{\text{макс. число вагонов в сходе}}{\text{общее кол-во вагонов}}$ );
7.  $features_7$  : (кривизна, скорость · загрузка,  $1 - \frac{\text{макс. число вагонов в сходе}}{\text{общее кол-во вагонов}}$ );
8.  $features_8$  : (скорость · загрузка,  $1 - \frac{\text{макс. число вагонов в сходе}}{\text{общее кол-во вагонов}}$ ).

Также в каждый набор признаков добавим новый признак *Intercept*, реализация которого всегда равна единице. Данный признак необходим для появления свободного члена  $\theta_0$  в результате скалярного произведения:  $\langle \theta, x \rangle = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ , где  $n + 1$  – размерность соответствующего признакового пространства.

Были добавлены новые признаки, такие как: профиль пути · макс. число вагонов в сходе,  $1 - \frac{\text{макс. число вагонов в сходе}}{\text{общее кол-во вагонов}}$ , скорость · загрузка. Вычислим коэффициенты корреляции новых признаков с целевой целевым признаком количество сошедших вагонов, а также вычислим коэффициенты корреляции между самими признаками (для удобства чтения таблицы сделаем соответствующие переименование признаков: target,  $f_1$ ,  $f_2$ ,  $f_3$ ):

Таблица 2.3 Корреляция целевого признака с введенными

	target	$f_1$	$f_2$	$f_3$
target	1.0	0.101375	-0.286535	0.198847
$f_1$	0.101375	1.0	-0.086693	-0.228508
$f_2$	-0.286535	-0.086693	1.0	-0.124420
$f_3$	0.198847	-0.228508	-0.124420	1.0

Из таблицы видно, что признаки  $f_1$  и  $f_3$  имеют значительную величину корреляции между собой. При этом признак  $f_2$  слабо коррелирует как с  $f_1$ , так и с  $f_3$ . Из этого следует, что в один набор признаков нежелательно включать  $f_1$  и  $f_3$  вместе.

Также можно заметить, что признак  $f_2$  сильно коррелирует с целевым признаком.

## 2.2.2 Пуассоновская регрессия

Предположим, что количество сошедших вагонов имеет Пуассоновское распределение. Тогда функция вероятности  $p(k) = \frac{\lambda^k}{k!} e^{-\lambda}$ . Пусть плотность потока событий  $\lambda = \lambda(\theta, x)$ , где  $\theta$  – вектор параметров,  $x$  – вектор, описывающий объект. Выберем набор функций  $\lambda$ , параметризованных по  $\theta$ :

1.  $\lambda_1(\theta, x) = e^{\langle \theta, x \rangle}$ ;
2.  $\lambda_2(\theta, x) = e^{-[\langle \theta, x \rangle]^2}$ ;
3.  $\lambda_3(\theta, x) = \sqrt{|5^2 - (\langle \theta, x \rangle - 5)^2|} + 1$ ;
4.  $\lambda_4(\theta, x) = (\langle \theta, x \rangle - 1)^2$ ;
5.  $\lambda_5(\theta, x) = \frac{1}{1 + (\langle \theta, x \rangle)^2}$ ;
6.  $\lambda_6(\theta, x) = \langle \theta, x \rangle \left( \frac{\pi}{2} + \arctan(\langle \theta, x \rangle) \right) + 1$ ;
7.  $\lambda_7(\theta, x) = \log(1 + (\langle \theta, x \rangle)^2) + 1$ .

В качестве оптимизационного метода был выбран топологический метод глобальной оптимизации SHGO (Simplicial Homology Global Optimisation), реализованный в модуле `scipy.optimize`. В качестве граничного множества для искомых параметров был взят гиперкуб со стороной 2000 и центром в начале координат (т.е.  $-1000 \leq \theta_i \leq 1000 \quad i = \overline{0, n}$ , где  $n$  – размерность соответствующего признакового пространства).

Ранее была получена функция логарифмического правдоподобия:  $l(\theta, x) = \sum_{i=1}^N (-\lambda(\theta, x) + y_i \ln(\lambda(\theta, x)) - \ln(y_i!))$ . Программно реализовав данную функцию и запустив метод `fit_all` у объекта класса MLM, были получены следующие результаты:

Таблица 2.4 Модели для Пуассоновской регрессии с признаковыми пространствами  $\{features_i\}_{i=1}^8$  и  $\{\lambda_i\}_{i=1}^7$

$n = 46$	модели с признаковым пространством $features_1$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	466.96	567.74	466.62	466.7	567.74	486.46	467.74
$\hat{\theta}$	[1.18, -271.63]	[0.00, -0.01]	[0.46, -287.10]	[2.80, -208.16]	[0.00, -0.00]	[999.96, -7.02]	[2.66, -549.74]

## Продолжение таблицы 2.4

$n = 41$	модели с признаковым пространством $features_2$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	444.09	548.71	425.81	438.43	551.73	470.21	429.43
$\hat{\theta}$	[1.17, -160.63, 39.27]	[0.08, -378.08, 99.63]	[0.38, -484.09, -46.55]	[2.83, -212.71, 56.93]	[0.08, -348.32, 93.04]	[999.85, 727.28, 609.78]	[-3.28, 486.69, -276.38]
$n = 37$	модели с признаковым пространством $features_3$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	418.35	534.39	408.81	411.31	538.51	450.29	401.29
$\hat{\theta}$	[1.29, -201.68, 0.82]	[-0.09, 588.72, -3.29]	[10.77, -204.22, 1.40]	[2.94, -280.76, 1.44]	[0.09, -583.40, 3.24]	[617.90, 984.90, -945.55]	[-3.91, 1000.00, -7.26]
$n = 42$	модели с признаковым пространством $features_4$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	400.59	561.05	419.03	394.93	561.41	465.11	411.03
$\hat{\theta}$	[2.11, -180.45, -2.36]	[0.19, -105.94, -0.63]	[-1.46, -12.36, 1.59]	[3.70, -207.40, -1.96]	[-0.16, 86.41, 0.54]	[437.60, -209.12, -457.37]	[5.94, -74.48, -6.39]
$n = 42$	модели с признаковым пространством $features_5$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	412.99	537.99	421.99	402.35	538.95	458.69	408.95
$\hat{\theta}$	[-0.09, -237.32, 0.03]	[0.82, 120.24, -0.01]	[-0.11, -53.49, 0.01]	[1.56, -216.06, 0.03]	[-0.73, -97.15, 0.01]	[-377.36, -871.14, 131.07]	[0.39, 432.09, -0.07]
$n = 37$	модели с признаковым пространством $features_6$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	392.23	520.73	429.17	384.45	523.69	442.69	393.29
$\hat{\theta}$	[-0.15, -148.03, 53.18, 0.03]	[0.40, -373.62, 111.98, -0.01]	[-4.41, -29.28, -174.92, 0.18]	[1.71, -244.78, 52.47, 0.02]	[0.45, -362.51, 112.78, -0.01]	[-96.52, 103.35, -840.07, 35.09]	[-0.31, -633.85, 7.06, 0.08]
$n = 35$	модели с признаковым пространством $features_7$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	356.87	522.63	378.31	401.53	523.29	437.51	371.07
$\hat{\theta}$	[0.92, -105.33, 34.34, 0.02, -2.24]	[-0.26, 356.46, -105.69, 0.00, -0.03]	[9.37, -141.15, 76.38, 0.03, -0.75]	[2.84, -1000.00, -169.76, 0.00, -0.56]	[0.31, 425.76, -162.84, -0.00, -1.26]	[800.97, -788.99, 172.08, 7.95, -975.81]	[-4.47, 445.37, -652.01, -0.08, 10.05]
$n = 37$	модели с признаковым пространством $features_8$						
$\lambda_i$	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$
$AIC_c$	367.37	531.39	451.59	376.29	531.61	446.69	371.21
$\hat{\theta}$	[0.55, 39.28, 0.03, -2.14]	[-0.94, 5.26, 0.01, 0.81]	[17.77, 88.80, -0.06, -16.75]	[-0.49, -15.96, -0.02, 1.44]	[-0.81, 3.67, 0.01, 0.68]	[-614.64, 544.88, 37.68, 613.23]	[2.73, -541.54, 0.09, -8.99]

Общий диапазон значений критерия Акаике для всех моделей:  $AIC_c \in [356.87, 567.74]$ .

Наилучшее значение по показателям качества получилось для модели  $(features_7, \lambda_1)$ . Таким образом, в моделях Пуассоновской регрессии показатель качества по критерию Акаике для наилучшей модели равняется 356.87.

Будем считать модель тривиальной, если ее признаковое пространство имеет размерность, равную единице. При этом для всех наблюдений данное свойство имеет постоянное значение. Вычислим значение критерия Акаике для тривиальной модели. В качестве исходных наблюдений возьмем такой же набор, как в модели, соответствующей  $(features_7, \lambda_1)$ . Проведя вычисления, было получено, что  $AIC_{c,trivial} = 551.27$ . Найдем отношение  $AIC_c$ , соответствующей наилучшей модели, к  $AIC_{c,trivial}$ . По данному отношению можно считать качество модели, которое может изменяться от нуля до еди-

ницы. При этом чем ближе значение к нулю, тем больше качество модели. Рассчитав значение данного отношения, получим, что  $\frac{AIC_c}{AIC_{c,trivial}} = 0.54$ .

### 2.2.3 Геометрическая регрессия

Предположим, что количество сошедших вагонов имеет геометрическое распределение. Тогда функция вероятности  $p(n) = (1 - p)^n p$ . Пусть вероятность успеха в серии испытаний Бернулли  $p = p(\theta, x)$ , где  $\theta$  – вектор параметров,  $x$  – вектор, описывающий объект. Выберем набор функций  $p$ , параметризованных по  $\theta$ :

1.  $p_1(\theta, x) = e^{\langle \theta, x \rangle}$ ;
2.  $p_2(\theta, x) = e^{-[\langle \theta, x \rangle]^2}$ ;
3.  $p_3(\theta, x) = \frac{1}{1 + e^{-\langle \theta, x \rangle}}$ ;
4.  $p_4(\theta, x) = \frac{1}{1 + [\langle \theta, x \rangle]^2}$ ;
5.  $p_5(\theta, x) = \langle \theta, x \rangle \left( \frac{\pi}{2} + \arctan(\langle \theta, x \rangle) \right) + 1$ .

Метод SHGO, использованный при оптимизации в Пуассоновской регрессии, оказался не успешным для геометрической регрессии и данного набора данных. Ни для одной модели геометрической регрессии не удалось найти оптимальную точку. Вероятно, это связано со сложным видом оптимизационной функции в признаковых пространствах. По этой причине были использованы другие методы глобальной оптимизации (Dual Annealing, Differential Evolution, Basin-hopping), также реализованные в модуле `scipy.optimize`.

Наилучшие результаты были получены с помощью метода Dual Annealing (алгоритм имитации обжига). В качестве граничного множества для искомых параметров был взят гиперкуб со стороной 2000 и центром в начале координат (т.е.  $-1000 \leq \theta_i \leq 1000 \quad i = \overline{0, n}$ , где  $n$  – размерность соответствующего признакового пространства).

Ранее была получена функция логарифмического правдоподобия:  $l = \sum_{i=1}^N (y_i \ln(1 - p(\theta, x_i)) + \ln(p(\theta, x_i)))$ . Программно реализовав данную функцию и запустив метод `fit_all` у объекта класса MLM, были получены следующие результаты:

Таблица 2.5 Модели геометрической регрессии с признаковыми пространствами  $\{features_i\}_{i=1}^8$  и  $\{p_i\}_{i=1}^5$

$n = 46$	модели с признаковым пространством $features_1$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	200.42	200.66	200.7	200.36	200.3
$\hat{\theta}$	[−1.46, 215.57]	[−1.20, 90.38]	[−1.18, 278.31]	[−1.82, 233.20]	[−0.95, 150.96]
$n = 41$	модели с признаковым пространством $features_2$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	181.75	186.19	185.53	178.59	180.45
$\hat{\theta}$	[−1.53, 269.04, −81.41]	[−1.14, 853.95, 106.70]	[−1.27, 388.36, −103.97]	[1.97, −392.25, 123.78]	[−0.99, 171.43, −53.30]
$n = 37$	модели с признаковым пространством $features_3$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	170.23	173.19	172.89	165.23	168.99
$\hat{\theta}$	[−1.54, 267.47, −1.94]	[1.26, −171.58, 0.90]	[−1.40, 560.16, −3.07]	[−2.05, 523.11, −3.38]	[−1.00, 147.49, −1.52]
$n = 42$	модели с признаковым пространством $features_4$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	176.65	176.33	176.67	174.73	346.83
$\hat{\theta}$	[−2.06, 137.66, 1.43]	[−1.57, 111.19, 0.86]	[−2.35, 342.68, 2.72]	[−2.75, 255.75, 1.99]	[−10.78, 23.77, 11.16]
$n = 42$	модели с признаковым пространством $features_5$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	171.15	173.51	173.91	168.97	170.15
$\hat{\theta}$	[0.04, 131.28, −0.03]	[0.58, −111.49, 0.01]	[0.80, 332.34, −0.04]	[−0.39, 279.20, −0.03]	[0.01, 89.40, −0.02]
$n = 37$	модели с признаковым пространством $features_6$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	167.23	164.09	164.29	160.83	267.79
$\hat{\theta}$	[−1.14, 163.99, −95.46, −0.01]	[−0.72, 145.51, −28.23, −0.01]	[0.39, 439.12, −89.71, −0.03]	[0.91, −413.70, 76.62, 0.02]	[−6.40, 358.05, −222.08, 0.05]
$n = 35$	модели с признаковым пространством $features_7$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	595.67	170.01	158.03	153.65	363.71
$\hat{\theta}$	[7.99, −357.44, −652.34, −0.25, −10.33]	[1.12, −818.27, −100.66, 0.00, −0.22]	[−1.46, 456.63, −52.58, −0.01, 2.33]	[−2.01, 29.36, 226.20, −0.03, 3.85]	[−25.49, −793.87, −999.99, 0.19, 14.44]
$n = 37$	модели с признаковым пространством $features_8$				
$p_i$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$AIC_c$	918.43	177.75	164.23	159.93	430.49
$\hat{\theta}$	[−32.01, −696.21, 0.31, 12.04]	[−0.18, −49.54, −0.03, 1.64]	[−0.59, −15.12, −0.02, 1.88]	[1.26, −99.36, 0.04, −5.14]	[−34.82, −883.84, 0.32, 14.54]

Общий диапазон значений критерия Акаике для всех моделей:  $AIC_c \in [153.65, 918.43]$ .

Наилучшее значение по показателям качества получилось для модели  $(features_7, p_4)$ . Таким образом, в моделях геометрической регрессии показатель качества по критерию Акаике для наилучшей модели равняется 153.65.

Так же как и для Пуассоновской регрессии вычислим отношение значения критерия Акаике лучшей модели к соответствующей ей тривиальной модели:  $\frac{AIC_c}{AIC_{c,trivial}} = \frac{153.65}{249.89} = 0.61$ .

Следует отметить, что большая часть полученных значений критерия Акаике в геометрической регрессии меньше, чем полученные значения критерия Акаике в Пуассоновской регрессии, что говорит о том, что модели геометрической регрессии для данной задачи являются более предпочтительными по сравнению с моделями Пуассоновской регрессии.

## ЗАКЛЮЧЕНИЕ

В данной работе была рассмотрена проблема схода подвижных единиц с рельсов по причине излома боковой рамы. В ходе предварительного анализа данных были вычислены описательные статистики, была построена матрица корреляции признаков, также был проведен анализ пропусков в данных и построены парные графики. При визуальном анализе парных графиков, а также при анализе описательных статистик было установлено, что выбросов в данных нет. При рассмотрении матрицы корреляции был сделан вывод, что данные в наборе имеют свойство избыточности. Так, например, признаки вес и загрузка имеют сильную корреляцию. Построив оценку функции вероятности количества сошедших вагонов были сделаны предположения, что она имеет пуассоновское распределение, либо геометрическое распределение.

Поскольку с целевым признаком остальные признаки слабо коррелируют, были предложены новые признаки, построенные на основе имеющихся. Удалось построить признак  $f_2$ , имеющий сильную по сравнению с остальными признаками корреляцию с целевым признаком, при этом слабую корреляцию с остальными признаками. Другие сконструированные признаки  $f_1$  и  $f_3$  имеют меньшую корреляцию с целевым признаком и значимую корреляцию между собой.

В данной работе была написана программная реализация метода максимального правдоподобия в общем виде. В конструктор данного класса передаются: ссылка на логарифмическую функцию правдоподобия, ссылка на метод оптимизации, граничные условия, ссылка на функцию предсказания, список параметрических функций, список признаков пространств, название целевой переменной, набор данных.

После были построены регрессионные модели. Для пуассоновской регрессии был выбран набор параметрических функций, обычно в качестве параметрических функций исследователи берут экспоненциальный вид функции, однако были рассмотрены и другие варианты. Аналогично был создан набор параметрических функций для геометрической регрессии. Для обеих регрессий были сформированы признаковые пространства, включающие



сконструированные признаки.

Мерой качества был выбран скорректированный критерий Акаике  $AIC_c$ . Проведя численные эксперименты, были получены следующие результаты:

- диапазон значений  $AIC_c$  для Пуассоновской регрессии:  $[356.87, 567.74]$ ;
- диапазон значений  $AIC_c$  для геометрической регрессии:  $[153.65, 918.43]$ .

Геометрические модели регрессии для заданного набора данных оказались существенно лучше моделей пуассоновской регрессии. Большинство моделей Пуассоновской регрессии оказались хуже моделей геометрической регрессии по критерию Акаике.

## СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Замышляев А.М., Игнатов А.Н., Кибзун А.И., Новожилов Е.О. Функциональная зависимость между количеством вагонов в сходе из-за неисправностей вагонов или пути и факторами движения // Надежность. – 2018. – С.1-15.
2. Samuel Arthur. Some Studies in Machine Learning Using the Game of Checkers // IBM Journal. – 1959. – №3. – 20 с. [Electronic resource]. URL: <http://www.cs.virginia.edu/evans/greatworks/samuel1959.pdf> (date of treatment: 25.11.2021).
3. Andrew Ng. Machine Learning // Stanford Online [Electronic resource]. URL: <https://www.coursera.org/learn/machine-learning> (date of treatment: 03.12.2021).
4. Xiao Zhang. Poisson Regression // Microsoft Documentation. – 2019. [Electronic resource]. URL: <https://docs.microsoft.com/en-us/previous-versions/azure/machine-learning/studio-module-reference/poisson-regression> (date of treatment: 02.12.2021).
5. Sachin Date. An Illustrated Guide to the Poisson Regression Model // Towards Data Science. – 2019. [Electronic resource]. URL: <https://towardsdatascience.com/an-illustrated-guide-to-the-poisson-regression-model-50cccba15958> (date of treatment: 10.12.2021).
6. Gregory Piatetsky. CRISP-DM, still the top methodology for analytics, data mining, or data science projects // KDnuggets. – 2014. [Electronic resource]. URL: <https://www.kdnuggets.com/2014/10/crisp-dm-top-methodology-analytics-data-mining-data-science-projects.html> (date of treatment: 15.12.2021).
7. Alexey Grigorev. Normal Equation [Electronic resource]. URL: [http://mlwiki.org/index.php/Normal\\_Equation](http://mlwiki.org/index.php/Normal_Equation) (date of treatment: 15.12.2021).
8. Воронцов К.В. Введение в машинное обучение // НИУ ВШЭ, Yandex School of Data Analysis [Электронный ресурс]. URL:

- <https://www.coursera.org/learn/vvedenie-mashinnoe-obuchenie> (дата обращения: 04.12.2021).
9. Константин Коточигов. CRISP-DM. // ГК ЛАНИТ. – 2019. [Электронный ресурс] URL: <https://habr.com/ru/company/lanit/blog/328858/> (дата обращения: 16.12.2021).
  10. Анатолий Карпов. Основы статистики. // Институт биоинформатики. – 2020. [Электронный ресурс]. URL: <https://stepik.org/course/76/info> (дата обращения: 04.12.2021).
  11. Николай Мальковский. Обзор градиентных методов в задачах математической оптимизации // 2018. [Электронный ресурс]. URL: <https://habr.com/ru/post/413853/> (дата обращения: 16.12.2021).
  12. Юрий Кашницкий. Метод стохастической аппроксимации // Open Data Science. – 2017. [Электронный ресурс]. URL: <https://habr.com/ru/company/ods/blog/326418/> (дата обращения: 21.12.2021).
  13. Регрессионный анализ [Электронный ресурс]. URL: [http://www.machinelearning.ru/wiki/index.php?title=Регрессионный\\_анализ](http://www.machinelearning.ru/wiki/index.php?title=Регрессионный_анализ) // 2016. (дата обращения: 03.12.2021).
  14. Воронцов К.В. Недообучение и переобучение в машинном интеллекте // ПостНаука. – 2020. [Электронный ресурс]. URL: <https://postnauka.ru/video/154955> (дата обращения: 03.12.2021).
  15. Илья Полосухин. Классификация и регрессия с помощью деревьев принятия решений // 2011. [Электронный ресурс]. URL: <https://habr.com/ru/post/116385/> (дата обращения: 03.12.2021).

## ПРИЛОЖЕНИЕ А

### Программная реализация метода максимального правдоподобия.

Для построения предсказательной модели использованы Пуассоновская и геометрическая регрессии. Оба метода в своей основе используют метод максимального правдоподобия. Именно поэтому был реализован общий класс для метода максимального правдоподобия. Сигнатура конструктора данного класса выглядит следующим образом:

```
def MLM(log_likelihood_fun,  
        optimization_method,  
        borders, predict_fun,  
        count_of_param_fun,  
        features, target, df)
```

где:

- `log_likelihood_fun` – функция логарифмического правдоподобия;
- `optimization_method` – метод глобальной оптимизации;
- `borders` – границы поиска оптимальной точки;
- `predict_fun` – функция предсказания целевой переменной по заданному объекту;
- `count_of_param_fun` – количество параметризованных функций;
- `features` – вектор наборов названий признаков из признаковых пространств;
- `target` – название целевого признака;
- `df` – набор данных;

Таким образом, для построения регрессионной модели с заданным распределением целевой переменной достаточно определить функцию, вычисляющую логарифмическую функцию правдоподобия. Также можно переопределить функцию предсказания результата `predict`, либо воспользоваться готовой реализацией. Оставшиеся параметры в конструкторе класса не требуют от пользователя реализации функций.

## Листинг программы

---

```
1 class MLM():
2
3     log_likelihood_fun = 0
4
5     predict_fun = None
6
7     optimization_method = shgo
8     borders = 10 * [(-1000, 1000)]
9
10    df = None
11    spaces = None
12    target = None
13
14    count_of_param_fun = 0
15
16    thetas_for_tex = []
17    AICs_c_for_tex = []
18
19
20    def __init__(self, df, spaces, target,
21                 log_likelihood_fun,
22                 optimization_method, borders, predict_fun,
23                 count_of_param_fun):
24        self.log_likelihood_fun = log_likelihood_fun
25        self.optimization_method = optimization_method
26        self.borders = borders
27        self.predict_fun = predict_fun
28        self.count_of_param_fun = count_of_param_fun
29        self.spaces = spaces
30        self.target = target
31        self.df = df
32
33
34    def neg_log_likelihood_fun(self, theta, X, y, lambda_index):
35        return -self.log_likelihood_fun(theta, X, y, lambda_index)
36
37
38    def calc_aic_c(self, logL, space) -> float:
39        k = len(space)
40        n = len(self.y)
41        eps = 1e-3
42        return round(2 * (k - logL) + 2 * k * (k + 1) / (n - k - 1 + eps), 2)
43
```

```

44
45 def round_and_format_params(self, a, is_tex_output=False):
46     if is_tex_output:
47         result_str = ''
48         for i in a:
49             result_str += ('%.2f$, ' % i)
50             result_str = '[' + result_str[:-2] + ']'
51         return result_str
52     else:
53         result = []
54         for i in a:
55             result.append(round(i, 2))
56         return result
57
58
59 def print_tex(self):
60     print('\\newcommand{\\lambdasTab}{\\lambda_i$}')
61     print('\\newcommand{\\criteriaTab}{AIC_c$}')
62     print('\\newcommand{\\paramsTab}{\\hat{\\theta}$}')
63     print('\\newcommand{\\lenFirstColumnTab}{1.5cm}')
64     print('\\begin{center}\\n')
65     space_ind = 0
66     for space in self.spaces:
67         space_ind += 1
68         data = self.df[space + self.target].dropna()
69         self.y = np.array(data[self.target])
70
71         lambdas_str = ''
72         for i in range(self.count_of_param_fun):
73             lambdas_str += ' & $\\lambda_' + str(i + 1) + '$'
74
75         print('\\resizebox{\\textwidth}{!}{')
76         print('\\t\\begin{tabular}{|p{\\lenFirstColumnTab}||' +
77             '\\p{4cm}|' * (self.count_of_param_fun) + '}')
78
79         print('\\t\\t\\hline')
80         print('\\t\\t$n = ' + str(len(self.y)) + '$ & \\multicolumn{' +
81             str(self.count_of_param_fun) + '}{c|}{модели с признаковым
82             пространством $features_' + str(space_ind) + '$} \\\\
83             \\hline\\hline')
84         print('\\t\\t\\lambdasTab' + lambdas_str + ' \\\\ \\hline')
85         print('\\t\\t\\criteriaTab & ' +
86             str(self.AICs_c_for_tex[space_ind-1]) + ' \\\\ \\hline')
87         print('\\t\\t\\paramsTab & ' + str(self.thetas_for_tex[space_ind-1])
88             + ' \\\\ \\hline')

```

```

84         print('\t\\end{tabular}')
85         print('}\\n')
86     print('\\\\end{center}')
87
88
89     def fit_one(self, space, link_fun_index):
90         count_of_attempt = 0
91         result = self.optimization_method(self.neg_log_likelihood_fun,
92         self.borders[:len(space)],
93         args=(self.X, self.y, link_fun_index))#, maxiter=1500)
94
95         while count_of_attempt < 20 and math.isnan(result['fun']):
96             result = self.optimization_method(self.neg_log_likelihood_fun,
97             self.borders[:len(space)],
98             args=(self.X, self.y, link_fun_index))#, maxiter=1500)
99             count_of_attempt += 1
100             if count_of_attempt == 20:
101                 print('count_of_attempt = 20')
102
103         logL = round(-result['fun'], 2)
104         AIC_c = self.calc_aic_c(logL, space)
105
106         return {
107             "status": result['success'],
108             "metrics": {
109                 "logL": logL,
110                 "AIC_c": AIC_c
111             },
112             "parameters": result['x']
113         }
114
115
116     def fit_all(self, get_tex_code=False):
117         space_ind = 0
118         results = []
119         self.thetas_for_tex = []
120         self.AICs_c_for_tex = []
121         for space in self.spaces:
122             space_ind += 1
123
124             data = self.df[space + self.target].dropna()
125             self.y = data[self.target]
126             self.y = np.array(self.y)
127             self.X = data.drop(self.target, axis=1)
128             self.X = np.matrix(self.X)
129

```

```

130         thetas_for_tex_line = ''
131         AICs_c_for_tex_line = ''
132
133         print('Признаковое пространство', space_ind, ':', space)
134         print('Мощность выборки n = ', len(self.y))
135
136         for link_fun_index in range(self.count_of_param_fun):
137
138             result = self.fit_one(space, link_fun_index)
139
140             if get_tex_code:
141                 thetas_for_tex_line +=
142                     ↪ self.round_and_format_params(result['parameters'],
143                     ↪ is_tex_output=True) + ' & '
144                 AICs_c_for_tex_line += '$' + str(result['metrics']['AIC_c'])
145                     ↪ + '$' + ' & '
146
147             result['parameters'] =
148                 ↪ self.round_and_format_params(result['parameters'])
149             results.append(result)
150             print(result)
151
152         print()
153
154         if get_tex_code:
155             self.thetas_for_tex.append(thetas_for_tex_line[:-3])
156             self.AICs_c_for_tex.append(AICs_c_for_tex_line[:-3])
157
158         if get_tex_code:
159             self.print_tex()

```

---



## ПРИЛОЖЕНИЕ Б

### Веб-сервис 'Метод максимального правдоподобия'.

В данной работе был разработан программный комплекс, реализующий работу с методом максимального правдоподобия. В отличие от программы из приложения А данный комплекс имеет пользовательский интерфейс, его не нужно скачивать, поскольку он является веб-ресурсом.

Рассмотрим пример взаимодействия с сервисом. На главной странице сайта пользователю предлагается загрузить в систему набор данных. Загрузить файл можно либо перетаскив его в соответствующий блок, либо выбрать из всплывающего окна по соответствующей кнопке. При этом данный файл должен иметь расширение csv и быть не более 50 Мб. После валидации загруженного файла системой кнопка 'Далее' становится активной, пользователю предлагается перейти на следующий шаг. На каждом последующем шаге есть возможность вернуться на предыдущий этап с помощью навигационного меню, расположенного слева, при этом если на нем внести изменения, последующие этапы станут недоступны их будет необходимо пройти снова.

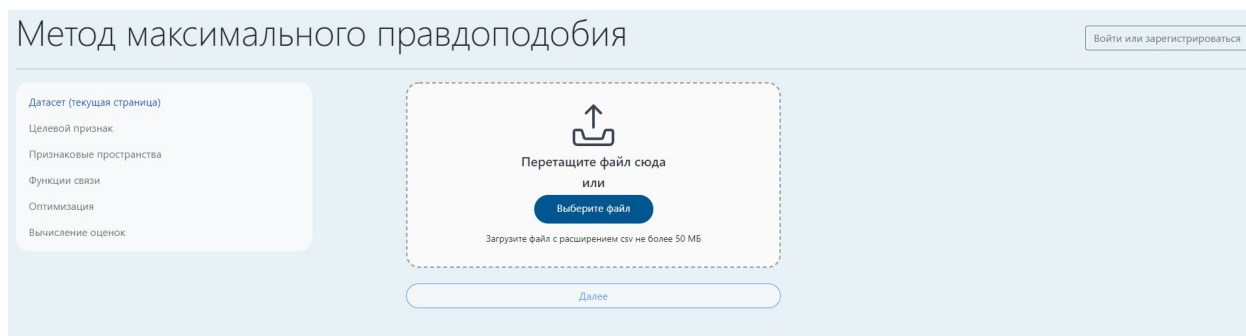


Рисунок А.1 Главная страница веб-сервиса

На следующем шаге необходимо выбрать из выпадающего списка целевой признак. Признаки были получены в результате разбора загруженного пользователем файла. Также необходимо выбрать распределение целевого признака из доступных вариантов. На момент написания данной работы поддерживается 2 возможных распределения: Пуассоновское, геометрическое.

Рисунок А.2 Введение целевого признака

На 3 шаге пользователю необходимо составить набор признаковых пространств. Для этого в соответствующем блоке справа выбираются необходимые признаки и вводится название данного признакового пространства. После его добавления оно отображается в центральном блоке. Для компактности в центральном блоке пишутся не названия признаков, а их порядковые номера, порядок которых определен в соответствующем списке справа.

Рисунок А.3 Определение признаковых пространств

После задания признаковых пространств для случаев Пуассоновской или геометрической регрессии необходимо указать функции связи. Изначально есть 2 предложенных пользователю варианта ( $\exp(t)$ ,  $\exp(-t^2)$ ). В поле для ввода также есть список логарифмов функций связи, определенных в списке выше, их необходимость обусловлена техническими особенностями. Поскольку в функции логарифмического правдоподобия может присутствовать логарифм от функции связи, то предпочтительней использо-

вать их упрощенные записи, например  $-t^2$ , вместо  $\ln(\exp(-t^2))$ . Пользователь может составить список логарифмов функций связи 'обернув' каждую функцию из списка выше логарифмом, однако для некоторых случаев метод оптимизации может не сойтись.

The screenshot shows a web interface titled "Метод максимального правдоподобия" (Method of maximum likelihood). On the left is a sidebar with navigation links: "Датасет", "Целевой признак", "Признаковые пространства", "Функции связи (текущая страница)", "Оптимизация", and "Вычисление оценок". The main area has a heading "Функции связи:" and a text input field containing the following code:

```
link_funs = [
    exp(t),
    exp(-t * t),
]

log_link_funs = [
    t,
    -t * t,
]
```

Below the input field is a "Далее" (Next) button. In the top right corner, there is a link "Войти или зарегистрироваться" (Login or register).

Рисунок А.4 Определение функций связи

На данном этапе пользователю предлагается выбрать один из методов оптимизации из выпадающего списка.

The screenshot shows the same web interface as Figure A.4, but at the "Оптимизация" (Optimization) step. The sidebar navigation is the same. The main area has a heading "Метод оптимизации:" and a dropdown menu with "SHGO" selected. Below the dropdown is a "Далее" (Next) button. The "Войти или зарегистрироваться" link remains in the top right corner.

Рисунок А.5 Выбор метода оптимизации

После ввода необходимых данных происходит вычисление оценок и показателей качества моделей. Пользователь может наблюдать за прогрессом вычислений по соответствующему индикатору выполнения.

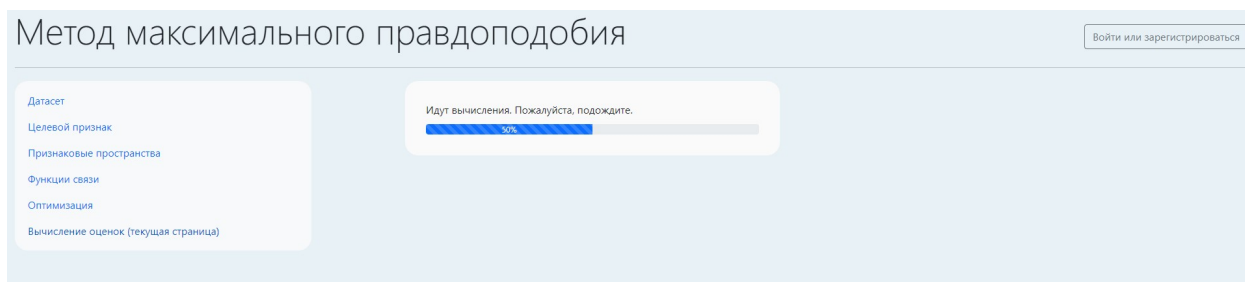


Рисунок А.6 Вычисление оценок моделей

По окончании вычислений пользователю предлагается получить результаты в одном (или нескольких) форматах. На момент написания данной работы поддерживается 4 формата: json, tex, txt, pdf.

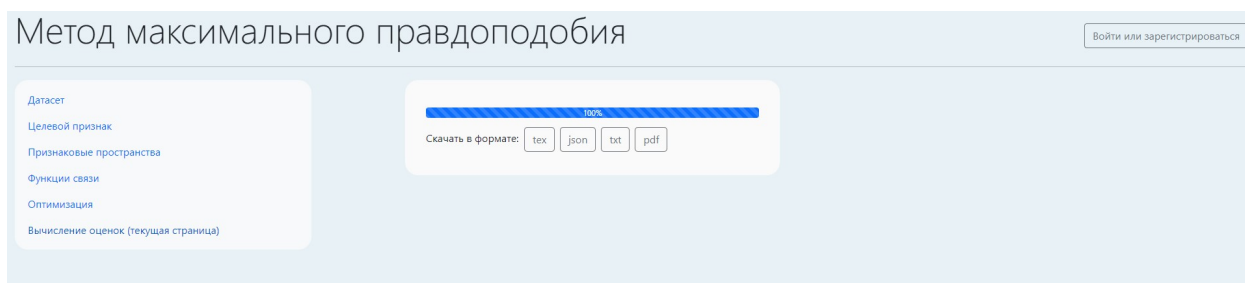


Рисунок А.7 Выгрузка вычисленных оценок

## Схема работы веб-сервиса

Сервис состоит из 3-х частей. Каждая часть – отдельный сервер. Рассмотрим все 3 части:

### 1. Сервер №1 (фронтенд).

Данная часть отвечает за визуальную составляющую приложения, а также за составление HTTP запросов к бекенд серверу. В качестве базового фронтенд-фреймворка был использован Angular, поэтому данный сайт реализует шаблон single-page-application, т.е. при посещении сайта пользователь загружает всего одну страницу, которая может динамически меняться. Верстка осуществлялась при помощи HTML, CSS, а также Bootstrap, благодаря которому сайт получился адаптивным под любые устройства. Макет приложения составлялся в Figma. Для написания скриптов был использован ЯП TypeScript.

## 2. Сервер №2 (бекенд).

Данная часть отвечает за внутреннюю логику приложения, скрытую от внешнего доступа и реализует архитектурный паттерн REST API. Рассмотрим пример выполнения основного запроса (вычисление параметров и оценивание моделей) к веб-сервису. После того, как пользователь ввел все необходимые данные фронтенд сервер формирует запрос к бекенд серверу. Бекенд сервер принимает данный запрос, вычисляет количество моделей, которые нужно обучить и распределяет их на группы. Для каждой группы асинхронно вызывается программа на Python, которая обучает модели и возвращает массив json объектов с информацией об обученных моделях. После завершения вычислений внутри группы на фронтенд сервер отправляется сообщение с информацией о количестве обученных моделей, после чего происходит изменение индикатора прогресса. Для двунаправленного общения между бекенд сервером и фронтенд сервером был использован протокол WebSocket.

Более подробно рассмотрим устройство бекенд части:

- Java - основной язык, на котором написана серверная часть;
- для обучения моделей использовался язык программирования Python, программы которого вызываются из под Java;
- в основе приложения лежит фреймворк Spring Boot, в качестве контейнера сервлетов использовался Tomcat;
- для создания веб-части использовался фреймворк Spring Web;
- в качестве реализации протокола WebSocket был использован Spring WebSocket;
- для реализации системы аутентификации и авторизации, а также настройки безопасности приложения был использован Spring Security;
- для сбора данных о времени работы отдельных модулей, а также дополнительном логировании были использован фреймворк Spring AOP;
- система логирования была написана при помощи библиотеки Log4J;
- юнит тестирование, а также интеграционное тестирование было реализовано при помощи JUnit и Mockito;

- для связи с базой данных были использованы Spring Data, а также MyBatis;
- в качестве системы контроля версий баз данных использовался FlyWay.  
Таким образом при клонировании проекта не обязательно копировать базу данных, достаточно запустить DDL скрипты с FlyWay;
- для создания pdf файлов с информацией об обученных моделях была использована библиотека Itext7.

### 3. Сервер №3 (система управления базой данных).

В качестве СУБД использовался PostgreSQL. Базы данных были необходимы для хранения регистрационной информации о пользователях. Также пользователи, прошедшие авторизацию имеют возможность сохранять результаты своих вычислений в системе.