

Государственное общеобразовательное
учреждение Ярославской области «Средняя школа
«Провинциальный колледж»

Исследовательское направление – информатика и ИКТ

Прогноз продаж торговой точки на основе методов машинного обучения

Курсовая работа

Выполнена ученицей
11 информационно-технологического
класса
ГООУ ЯО Средняя школа
«Провинциальный колледж»

Ситкиной Аленой

Ярославль, 2022

Оглавление

ВВЕДЕНИЕ	3
ГЛАВА 1. ОСНОВНЫЕ ПОНЯТИЯ	6
1.1. РАЗВИТИЕ МАШИННОГО ОБУЧЕНИЯ	6
1.2. ОСНОВНЫЕ ТИПЫ МАШИННОГО ОБУЧЕНИЯ	6
1.3. ПРЕДОБРАБОТКА ДАННЫХ	7
1.4. ПЕРЕОБУЧЕНИЕ И НЕДООБУЧЕНИЕ	8
1.5. ФУНКЦИОНАЛ ОШИБКИ И КРОСС-ВАЛИДАЦИЯ	8
1.6. ВИЗУАЛИЗАЦИЯ	9
ГЛАВА 2. ОСНОВНЫЕ МОДЕЛИ МАШИННОГО ОБУЧЕНИЯ	10
2.1. ЛИНЕЙНАЯ РЕГРЕССИЯ И ГРАДИЕНТНЫЙ СПУСК	10
2.1.1. Общее описание, функционал ошибки	10
2.1.2. Запись задачи в матричной форме	10
2.1.3. Решение задачи	11
2.1.4. Градиентный спуск	12
2.1.5. Переобучение, регуляризация, виды линейной регрессии	13
2.2. KNN	14
2.3. РЕШАЮЩИЕ ДЕРЕВЬЯ	16
2.4. СЛУЧАЙНЫЙ ЛЕС	17
2.5. ГРАДИЕНТНЫЙ БУСТИНГ	18
ГЛАВА 3. ИСПОЛЬЗУЕМЫЕ ИНСТРУМЕНТЫ	20
3.1. ЯЗЫК ПРОГРАММИРОВАНИЯ PYTHON	20
3.2. JUPITER NOTEBOOK	20
3.3. НЕОБХОДИМЫЕ БИБЛИОТЕКИ	20
ГЛАВА 4. РЕАЛИЗАЦИЯ ПРОГРАММЫ	23
ЗАКЛЮЧЕНИЕ	26
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ И ЛИТЕРАТУРЫ	27

Введение

Каждая организация стремится к оптимизации процесса закупок. В идеале, товар всегда должен быть в наличии, при этом его не должно быть в избытке. Прогнозирование того, сколько единиц каждого вида товара будет продано за определенный промежуток времени, очень важно для организации бизнеса. Анализ и прогноз продаж необходимы для корректных заказов товара, для уменьшения количества просроченного и списанного товара, для расширения и сбалансирования ассортимента торговой точки.

Актуальность исследования: отчасти алгоритм прогнозирования продаж торговой точки реализован в таких программах, как СБИС++ и 1С. Но, например, предприятия малого бизнеса зачастую не имеют доступа к полному пакету возможностей данных программ ввиду, во-первых, их стоимости, во-вторых, необходимости обращаться к услугам технических специалистов для их обслуживания. Поэтому прогноз прибыли предприятия зачастую осуществляется вручную. А это отнимает значительное количество времени, что влечет за собой дополнительные расходы в виде трудочасов (зарплаты). К тому же данный прогноз, во-первых, не является оперативным, во-вторых, из-за большого количества данных и медленного выполнения прогноза к моменту получения окончательных результатов он зачастую становится уже невостребованным.

Проект, разработанный в рамках данной курсовой работы, направлен на автоматизацию процесса прогнозирования продаж. Он позволит сэкономить рабочее время сотрудников, а также наиболее эффективно составить ассортимент, что приведет к привлечению новых покупателей и увеличению объема потребительской корзины существующих покупателей. Как итог, это позволит наиболее эффективно использовать оборотные средства предприятия, что повлечет за собой увеличение доходности торговой точки.

В работе мы будем опираться на данные торговой точки, расположенной в поселке Прибрежный, экспортированные в файл с расширением .xls из складской программы «1С: Торговое предприятие».

Цель исследования: разработать программу, прогнозирующую продажи торговой точки на основе технологии машинного обучения.

Задачи исследования:

1. Исследовать основные понятия, используемые в машинном обучении.
2. Исследовать принцип работы различных моделей.
3. Исследовать возможности языка Python и его дополнительных библиотек.
4. Сравнить качество работы различных моделей на исходных данных.

5. Создать графический интерфейс программы, осуществляющей прогноз продаж выбранной торговой точки.

Объект исследования: машинное обучение.

Предмет исследования: прогнозирование продаж сети торговых точек на основе машинного обучения.

Методы исследования:

▣ Анализ (необходимо для выявления принципов работы методов машинного обучения, принципов устройства различных функций потерь и способов оценки обобщающей способности)

▣ Обобщение (обобщение проделанной нами работы, подведение итогов)

▣ Синтез (использование и сочетание методов машинного обучения, функций потерь и различных инструментов в создаваемой нами программе)

▣ Эмпирический анализ (исследование информации о продажах выбранной торговой точки)

▣ Аналогия (используется при построении модели – закономерности, выявленные в обучающей выборке, будут применены к тестовым данным)

▣ Эксперимент (используется при тестировании написанного кода)

▣ Описание (описание принципа работы модели линейной регрессии)

▣ Моделирование (моделирование реальной ситуации – продаж торговой точки)

▣ Индукция (на основе определенной ограниченной выборки делаем выводы об общих закономерностях)

▣ Сравнение (сравнение моделей машинного обучения)

Обзор литературы и основных идей исследования:

Теоретический материал нашей работы построен на двух основных источниках – это книга Себастьяна Рашки и Вахида Мирджалили «Python и машинное обучение: машинное и глубокое обучение с использованием Python, scikit-learn и TensorFlow 2» и обучающий курс от НИУ ВШЭ на портале openedu.ru. В работе нами произведен обзор основных понятий (типы задач в машинном обучении, вопросы, связанные с предобработкой данных, переобучение и недообучение, функционал ошибки и кросс-валидация), а также принципов работы основных алгоритмов машинного обучения (kNN, линейная регрессия, решающие деревья, композиции алгоритмов). Для написания практической части нашей работы, а также глав, посвященных этому, мы обращались к документам языка Python и его библиотек (NumPy, Seaborn, CatBoost, LightGBM, Pandas, Scikit-learn). На основе изученных идей, принципов и методов была реализована программа, прогнозирующая продажи выбранной торговой точки.

Характеристика степени изученности темы:

Основные идеи и принципы работы моделей машинного обучения на данный момент изучены очень хорошо. И, как известно, машинное обучение находит широкое применение в сфере торговли. Но при этом предприятия малого бизнеса, ввиду описанных выше причин, зачастую не имеют доступа к программам, в которых реализованы необходимые алгоритмы машинного обучения.

Характеристика личного вклада автора в решение избранной проблемы и практическая значимость исследования:

В ходе данной работы нами была реализована программа, решающая задачу прогнозирования продаж торговой точки. Это решает первоначально стоящую проблему доступа предприятий малого бизнеса к данным инструментам.

Глава 1. Основные понятия

1.1. Развитие машинного обучения

Как пишут Себастьян Рашка и Вахид Мирджалили в книге «Python и машинное обучение: машинное и глубокое обучение с использованием Python, scikit-learn и TensorFlow 2», при нынешней зрелости современных технологий есть один ресурс, которого у нас вдоволь: значительный объем структурированных и неструктурированных данных. Во второй половине двадцатого века машинное обучение развивалось как подобласть искусственного интеллекта, включающая самообучающиеся алгоритмы, которые выводили знания из данных с целью вырабатывания прогнозов. Вместо того чтобы требовать от людей выводить правила вручную и строить модели путем анализа крупных объемов данных, машинное обучение предлагает более эффективную альтернативу для сбора знаний из данных, которая постепенно улучшает эффективность прогнозирующих моделей и принимает решения, управляемые данными. Важность машинного обучения возрастает не только в исследованиях, имеющих отношение к компьютерным наукам; его роль в нашей повседневной жизни становится еще больше. Машинное обучение позволяет нам пользоваться надежными фильтрами почтового спама, удобным программным обеспечением распознавания текста и речи и испытанными поисковыми механизмами. Кроме того, заметный прогресс был достигнут в медицинских приложениях.

1.2. Основные типы машинного обучения

Данный раздел написан по материалам книги С.Рашки и В.Мирджалили «Python и машинное обучение: машинное и глубокое обучение с использованием Python, scikit-learn и TensorFlow 2», страницы 31-38.

Теперь введем некоторую терминологию. Свойства объектов будем называть «признаками» (features), а переменная, значение которой мы учимся предсказывать – «целевая переменная» (target). Основные виды признаков (и целевых переменных): бинарные, категориальные, числовые.

Модели (методы) машинного обучения используются для решения множества самых разных задач. Опишем основные из них.

1. Обучение с учителем

При обучении с учителем у нас есть размеченные обучающие данные. Цель - обучить модель с их использованием, чтобы она могла предсказывать значение целевой переменной для не встречавшихся ей ранее объектов.

Классификация – это подкатегория обучения с учителем, где целевой переменной является категориальная метка класса нового объекта. Метки представляют собой дискретные неупорядоченные значения.

Регрессия (регрессионный анализ) – это подкатегория обучения с учителем, в которой модель прогнозирует непрерывный результат.

2. Обучение без учителя

В данном случае мы работаем с неразмеченными данными; целью является поиск неизвестной структуры в исходных данных.

Кластеризация – это методика анализа данных, целью которой является разбиение объектов на подгруппы (кластеры). Объекты из одного кластера при этом должны обладать определенной степенью подобия, но при этом они должны быть менее похожи на объекты других кластеров.

Понижение размерности – это подобласть обучения без учителя, целью которой является устранение из данных шума и сжатие данных в подпространство меньшей размерности с сохранением большей части существенной информации.

3. Обучение с подкреплением

Тип машинного обучения, целью которого является разработка системы (агента), улучшающей свои характеристики на основе взаимодействия со средой.

Так как мы работаем с задачей обучения с учителем (говоря более конкретно, с задачей регрессионного анализа), далее мы будем рассматривать именно ее.

1.3. Предобработка данных

Перед тем, как обучать нашу модель на данных, необходимо произвести их предобработку. Перечислим основные причины, почему это стоит делать.

- ☐ Поступившие к нам данные могут иметь форму, непригодную для дальнейшей работы;
- ☐ В данных могут быть пропуски. Не все модели поддерживают работу с данными с пропусками (например, дерево решений не поддерживает);
- ☐ В данных могут быть выбросы (объекты, имеющие значение целевой переменной, резко не соответствующее общей закономерности в данных); соответственно, нужно найти их (либо удостовериться, что их нет) и либо устранить, либо выбрать модель, устойчивую к ним;

- ❑ В данных может быть избыточная информация; соответственно, следует производить отбор признаков, прежде чем обучать на этих данных модель;
- ❑ Иногда бывает необходима стандартизация¹ данных для достижения наибольшей эффективности модели;
- ❑ Некоторые модели (например, линейная регрессия) не поддерживают работу с категориальными признаками, которые обычно встречаются практически в любом датасете. В таких случаях обычно используется one-hot кодирование: вместо одного признака мы формируем несколько (количество соответствует количеству уникальных значений в исходной колонке), один из них равен 1 (соответствует значению, которое было записано в исходном признаке), а остальные равны 0.

1.4. Переобучение и недообучение

Переобучение – это очень распространенная проблема в машинном обучении. Она возникает тогда, когда модель показывает хороший результат на тренировочной выборке, но плохо обобщается на тестовую.² При недообучении, напротив, модель недостаточно сложна для выявления структуры в данных.³ Эти проблемы решаются разными способами в зависимости от того, какую модель машинного обучения мы используем. Переобучение и недообучение, а также способы решения этих проблем будут подробнее рассмотрены в последующих главах при описании используемых нами моделей.

1.5. Функционал ошибки и кросс-валидация

При создании модели машинного обучения важно понимать, насколько хорошо она предсказывает целевую переменную. Функционал ошибки – это своего рода мера эффективности модели. Он выбирается в зависимости от решаемой задачи.

Заметим, что ошибка модели складывается из трех компонент: шум (noise), смещение (bias) и разброс (variance). Шум – это характеристика сложности и противоречивости данных. Он является свойством решаемой задачи. Смещение модели – это её способность приблизить лучшую среди всех возможных моделей. Разброс модели характеризует её

¹ Стандартизация – процесс придания данным характеристик стандартного нормального распределения: нулевого среднего и единичной дисперсии. (источник – тот же, что указан ниже, С. 75-76)

² Рашка С., Мирджалили В. Python и машинное обучение: машинное и глубокое обучение с использованием Python, scikit-learn и TensorFlow 2/ Пер. с англ. СПб.: ООО "Диалектика", 2020. С.109

³ Там же.

устойчивость к изменениям в обучающей выборке. Данная концепция пригодится нам в дальнейшем при изучении некоторых моделей.⁴

Для оценки качества модели нам необходимо разбить выборку на две части: тренировочную и тестовую. Но у такого подхода есть некоторые недостатки. При таком подходе мы подбираем гиперпараметры⁵ модели так, чтобы она показывала максимально хороший результат на тестовой выборке. То есть мы «подгоняем» гиперпараметры под тестовую выборку. Это значит, модель переобучается. Кросс-валидация – это способ, позволяющий бороться с данной проблемой. Мы по-прежнему будем разбивать выборку на тренировочную и тестовую, но при обучении тренировочная выборка разбивается на k блоков (обычно $k=5$) и каждый по очереди делается тестовым (правильнее сказать, валидирующим), остальные при этом формируют обучающую выборку. Так мы подбираем наилучшие гиперпараметры. Далее тестируем модель на тестовой (отложенной вначале) выборке.

1.6. Визуализация

Визуализация занимает очень важное место в машинном обучении. Зачастую она помогает увидеть новые закономерности в данных, обосновать введение нового признака. Также визуализировать можно модель и результаты ее работы. Это может подсказать нам, например, что модель переобучена. В языке Python существуют специальные библиотеки, предназначенные для визуализации: `matplotlib` и `Seaborn`, они будут подробнее рассмотрены в Главе 3.

В последующих главах описывается принцип работы, а также некоторые особенности разных моделей машинного обучения. Данные главы написаны на основе курса «Основы машинного обучения» от НИУ ВШЭ.⁶

⁴ Соколов Е.А., Зимовнов А.В., Ковалев Е.И., Кохтев В.М., Рысьмятова А.А., Филатов А.А. Основы машинного обучения: [Электронный ресурс]// Образовательная платформа «Открытое образование». URL: https://courses.openedu.ru/courses/course-v1:hse+INTRML+fall_2020/course/ (Дата обращения: 04.11.2021).

⁵ Гиперпараметры – это параметры модели, которые подбираются человеком вручную.

⁶ Там же, что и 4.

Глава 2. Основные модели машинного обучения

2.1. Линейная регрессия и градиентный спуск

2.1.1. Общее описание, функционал ошибки

Так как мы решаем задачу регрессии (предсказываем количество единиц товара, которое будет продано), мы будем использовать модель, которая чаще всего применяется для данного типа задач: линейную регрессию.

Итак, пусть мы рассматриваем объект x , который обладает d признаками. Тогда наша модель $a(x)$ будет иметь следующий вид: $a(x) = w_0 + w_1x_1 + \dots + w_dx_d$, где w_i - вес i -ого признака. Заметим, что данную формулу можно переписать в виде скалярного произведения векторов: $a(x) = w_0 + \langle w, x \rangle$. В общем случае можно считать, что у нас есть признак, который всегда равен единице, тогда формула переписывается в виде $a(x) = \langle w, x \rangle$.

Теперь опишем решаемую нами задачу. Для этого введем функционал ошибки. Мы будем использовать один из самых распространенных – MSE (mean squared error – среднеквадратичная ошибка). Она также является универсальной и может применяться при использовании практически любого алгоритма машинного обучения. Так как одной из наших задач является сравнение различных моделей, нами была выбрана именно эта функция потерь. Вычисляется она следующим образом: $\frac{1}{l} \sum_{i=1}^l (\langle w, x^i \rangle - y_i)^2$, где l – количество объектов в обучающей выборке. Тогда наша задача записывается следующим образом: $\frac{1}{l} \sum_{i=1}^l (\langle w, x^i \rangle - y_i)^2 \rightarrow \min_w$.

2.1.2. Запись задачи в матричной форме

Для удобства запишем нашу задачу в матричной форме. Введем матрицу объекты-признаки. Она содержит в себе всю обучающую выборку и записывается следующим

образом: $X = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{l1} & \dots & x_{ld} \end{pmatrix} \in \mathbb{R}^{l \times d}$, где строки соответствуют объектам, а столбцы –

признакам (считаем, что все они числовые или булевы). Также введем вектор весов $w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^d$ и вектор ответов (значений целевой переменной) для обучающей выборки $y =$

$\begin{pmatrix} y_1 \\ \vdots \\ y_l \end{pmatrix} \in \mathbb{R}^l$. Теперь заметим, что мы можем умножить матрицу X на вектор w (данная операция допустима). Сделаем это: $Xw = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{l1} & \cdots & x_{ld} \end{pmatrix} * \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^d w_i x_{1i} \\ \vdots \\ \sum_{i=1}^d w_i x_{li} \end{pmatrix} =$

$\begin{pmatrix} \langle w, x_1 \rangle \\ \vdots \\ \langle w, x_l \rangle \end{pmatrix} \in \mathbb{R}^l$. Заметим, что мы получили вектор предсказаний нашей модели для объектов

обучающей выборки. Теперь получим вектор отклонений предсказаний модели от ответов:

$Xw - y = \begin{pmatrix} \langle w, x_1 \rangle - y_1 \\ \vdots \\ \langle w, x_l \rangle - y_l \end{pmatrix}$. Вспомним такое понятие из линейной алгебры, как евклидова

норма. Для произвольного вектора z она записывается следующим образом: $\|z\|^2 = \sum_{j=1}^n z_j^2$.

Посчитаем евклидову норму для нашего вектора отклонений, умножим ее на $\frac{1}{l}$ и заметим, что мы как раз получим формулу для среднеквадратичной ошибки: $\frac{1}{l} \|Xw - y\|^2 = \frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2$. Таким образом, мы записали нашу задачу в матричной форме.

2.1.3. Решение задачи

Чтобы решить поставленную задачу, введем понятие градиента. Градиент – это вектор частных производных функции (обобщение понятия производной на функции нескольких переменных).

Нашей задачей является нахождение минимума для функции MSE. Она является строго выпуклой, следовательно, экстремум у неё ровно один – это и есть искомый глобальный минимум. При этом, как известно из высшей математики, в точке экстремума градиент функции равен 0. Это значит, что для получения ответа нам необходимо посчитать градиент MSE, приравнять его к нулю и решить полученную систему линейных уравнений.

Проделаем

эти

операции:

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l (y_i - a(x_i))^2 \rightarrow \min_w$$

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l (y_i - \langle w, x_i \rangle)^2 \rightarrow \min_w$$

$$\nabla_w Q = \begin{pmatrix} \frac{dQ}{dw_1} & \cdots & \frac{dQ}{dw_d} \end{pmatrix}$$

$$\begin{aligned}
\frac{dQ}{dw_j} &= \frac{1}{l} \sum_{i=1}^l \frac{d(y_i^2 - 2y_i \langle w, x_i \rangle + (\langle w, x_i \rangle)^2)}{dw_j} \\
&= \frac{1}{l} \sum_{i=1}^l \frac{dy_i^2}{dw_j} - \frac{2 * dy_i (w_1 x_{i1} + \dots + w_d x_{id})}{dw_j} + \frac{d(w_1 x_{i1} + \dots) (w_1 x_{i1} + \dots)}{dw_j} \\
&= \frac{1}{l} \sum_{i=1}^l 0 - 2y_i x_{ij} + 2x_{ij} \langle w, x_i \rangle = \frac{2}{l} \sum_{i=1}^l x_{ij} (\langle w, x_i \rangle - y_i)
\end{aligned}$$

Перепишем полученный результат в матричном виде:

$$\nabla \frac{1}{l} \|Xw - y\|^2 = \frac{2}{l} X^T (Xw - y)$$

Приравняем к нулю и решим полученную систему уравнений:

$$\begin{aligned}
\frac{2}{l} X^T (Xw - y) &= 0 \\
w &= (X^T X)^{-1} X^T y
\end{aligned}$$

Итак, мы в явном виде получили решение нашей задачи.

2.1.4. Градиентный спуск

Полученное нами выше решение имеет ряд существенных недостатков, как-то:

- ❑ Если матрица $X^T X$ – вырожденная, то обратной к ней не существует, и решение по данной формуле мы получить не сможем;
- ❑ Если матрица $X^T X$ не является вырожденной, но является почти вырожденной (определитель близок к нулю) будут проблемы с вычислением обратной матрицы;
- ❑ Нахождение обратной матрицы к матрице размера $n \times n$ выполняется за $O(n^3)$ операций; для больших n эти вычисления будут производиться очень долго;
- ❑ Решение в аналитическом виде существует не для всех возможных функций потерь; мы решили задачу только для одной выбранной нами функции.

Данные проблемы решаются применением градиентного спуска. Опишем алгоритм его работы.

Для начал введем понятие антиградиента. Антиградиент – это градиент, взятый с противоположным знаком. Он указывает на направление наискорейшего убывания функции.

Алгоритм градиентного спуска выглядит следующим образом:

1. Стартуем из случайной точки w^0 – это инициализация весов (может быть взята, например, из стандартного нормального распределения);

2. Делаем шаг спуска (сдвигаемся по антиградиенту): $w^t = w^{t-1} - \eta \nabla Q(w^{t-1})$, где η – размер шага;
3. Повторяем пункт 2, пока не окажемся достаточно близко к искомой точке. Можно рассматривать различные варианты определения момента останова, вот основные из них:

$$\square \quad \|w^t - w^{t-1}\| < \varepsilon$$

$$\square \quad \|\nabla Q(w^t)\| < \varepsilon$$

Также важно сделать несколько замечаний:

- Рассматриваемая функция потерь может, помимо глобального минимума, иметь несколько локальных минимумов. Чтобы получить корректный результат, используется мультистарт (мы запускаем алгоритм несколько раз, каждый раз выбирая разные стартовые точки);
- Длина шага является гиперпараметром: если шаг слишком большой, это может помешать сходимости алгоритма, а если слишком маленький, алгоритм может выполняться долго;
- Также помогает сходимости и ускоряет процесс масштабирование признаков;
- Из-за того, что данный алгоритм может работать долго (особенно на достаточно больших выборках), часто используют его модификацию – стохастический градиентный спуск. Идея заключается в том, что на каждом шаге мы берем градиент по нескольким случайно выбранным объектам, а не по всей выборке. Но стохастический градиентный спуск бывает нестабильным, при его использовании нужно аккуратно выбирать длину шага, иначе алгоритм может не сойтись.

Теперь мы можем записать, используя результаты, полученные ранее, как будет выглядеть формула для шага градиентного спуска, если мы используем MSE в качестве функции потерь. Она пригодится нам при реализации данного алгоритма:

$$w^t = w^{t-1} - \eta_t \frac{2}{l} X^T (Xw^{t-1} - y)$$

2.1.5. Переобучение, регуляризация, виды линейной регрессии

1. Гребневая регрессия

Известно, что одним из самых явных признаков переобучения модели линейной регрессии являются большие веса. Тогда можно рассмотреть следующую идею: будем штрафовать модель за большие веса (регуляризовать модель). Такая модель называется

гребневой регрессией (используемый здесь регуляризатор ещё иногда называют L2-нормой).

Теперь приведем строгую запись данной идеи:

$$\frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 + \lambda \|w\|^2 \rightarrow \min_w$$

Здесь λ – коэффициент регуляризации, он является гиперпараметром модели. Используя идеи, приведенные выше (можно показать, что данный функционал также строго выпуклый), решим поставленную задачу для регуляризованного функционала:

$$w = (X^T X + \lambda I)^{-1} X^T y$$

Здесь I – единичная матрица.⁷ Можно заметить, что добавление слагаемого λI зачастую будет даже из вырожденной матрицы делать невырожденную, и мы сможем найти решение задачи, используя аналитически выведенную формулу.

Осталось сделать лишь последнее замечание: в регуляризатор нельзя включать свободный коэффициент, так как он нужен для того, чтобы порядок предсказываемых значений был сопоставим с порядком целевой переменной; если мы потребуем все веса быть небольшими, а целевая переменная будет достаточно большой, модель не сможет качественно решить задачу.

2. LASSO

LASSO (Least Absolute Shrinkage and Selection Operator – в дословном переводе «оператор наименьшей абсолютной усадки и выбора»; используемый здесь регуляризатор ещё иногда называют L1-нормой) – это метод, идея которого заключается в том, что при регуляризации идет суммирование не квадратов весов, а модулей весов.

Приведем строгую запись данной идеи:

$$\frac{1}{l} \sum_{i=1}^l (\langle w, x_i \rangle - y_i)^2 + \lambda \sum_{j=1}^d |w_j| \rightarrow \min_w$$

Можно заметить, что при данном подходе некоторые веса обнуляются, что приводит к отбору признаков.

Итак, мы рассмотрели первую модель машинного обучения. Но можно заметить, что линейная регрессия будет показывать не самые лучшие результаты, если зависимость в данных не является линейной. Можно обучать модель с полиномиальными признаками, но тогда задача обучения модели становится на порядки сложнее и данных нужно на порядки больше. Поэтому необходимо рассмотреть и другие модели.

⁷ Единичная матрица – это матрица с единицами на главной диагонали и нулями на остальных позициях.

2.2. kNN

Принцип работы метода k ближайших соседей (k nearest neighbors) (напомним, мы разбираем только случай задачи регрессии) основан на гипотезе непрерывности. Согласно этой гипотезе, мы предполагаем, что близким объектам соответствуют близкие ответы. Обучение данной модели состоит в запоминании обучающей выборки.

Опишем алгоритм работы данной модели:

1. Пусть нам дан объект x . Для начала мы сортируем объекты обучающей выборки по расстоянию до объекта x : $\rho(x, x_1) \leq \rho(x, x_2) \leq \dots \leq \rho(x, x_l)$ (здесь l – это размер обучающей выборки, ρ – некоторая функция для вычисления расстояния между объектами (метрика));
2. Выбираем k ближайших соседей (здесь k – гиперпараметр);
3. Усредняем значения целевой переменной по выбранным ближайшим соседям:

$$a(x) = \frac{1}{k} \sum_{i=1}^k y_i.$$

Теперь остановимся подробнее на том, что такое метрика и какими бывают метрики. Метрика – это обобщение понятия “расстояние” на многомерные пространства. Метрикой называют функцию, принимающую два аргумента и удовлетворяющую следующим требованиям:

- ☐ $\rho(x, z) = 0 \iff x = z$;
- ☐ $\rho(x, z) = \rho(z, x)$;
- ☐ $\rho(x, z) \leq \rho(x, v) + \rho(v, z)$ ($\forall v$) – неравенство треугольника.

Приведем несколько примеров самых распространенных метрик (для числовых признаков):

- ☐ Евклидова метрика: $\rho(x, z) = \sqrt{\sum_{j=1}^d (x_j - z_j)^2}$;
- ☐ Манхэттенская метрика: $\rho(x, z) = \sum_{j=1}^d |x_j - z_j|$ (полезна в случаях, когда в данных много выбросов);
- ☐ Метрика Минковского: $\rho(x, z) = \sqrt[t]{\sum_{j=1}^d (x_j - z_j)^t}$ (здесь t – гиперпараметр).

Как можно заметить, для категориальных признаков эти формулы неприменимы. Существует много способов измерить расстояние между объектами, основываясь на категориальных признаках, но их рассмотрение выходит за рамки нашей курсовой работы. Как уже было сказано, мы можем применить one-hot кодирование и считать категориальные признаки числовыми.

Также на практике часто используют kNN с весами: более близкие объекты больше влияют на результат (имеют больший вес). Для расчета весов используют парзеновское окно: $w_i = K(\frac{\rho(x, x_i)}{h})$, где $K(z)$ – ядро (функция, которая не возрастает и положительна на отрезке $[0;1]$), h – ширина окна (гиперпараметр; он позволяет регулировать, с какого расстояния объекты практически перестают учитываться). Приведем некоторые примеры:

$$\boxtimes \text{ Гауссовское ядро: } K(z) = (2\pi)^{-0.5} \exp(-\frac{1}{2}z^2);$$

$$\boxtimes w_i = 1/\rho(x, x_i) \text{ (используется в библиотеке sklearn).}$$

Таким образом, формула для расчета значения нового объекта x будет выглядеть следующим образом (её называют формулой Надарая-Ватсона):

$$a(x) = \frac{\sum_{i=1}^k w_i y_i}{\sum_{i=1}^k w_i}$$

Можно заметить, что модель kNN показывает хорошие результаты, если обучающая выборка довольно большая. Также модель легко обучается. Но у неё есть несколько определенных недостатков, а именно: модель может работать достаточно долго (осуществляя поиск k ближайших соседей) и подбор гиперпараметров для неё довольно трудоемок. Поэтому необходимо также рассмотреть и другие модели.

2.3. Решающие деревья

Как понятно из названия модели, она основывается на построении графа. В его внутренних вершинах записываются предикаты, в листьях – прогнозы модели. Предикаты могут быть сколь угодно сложными (они могут содержать линейные модели, деревья и т. д.), но, как мы увидим позже, даже с использованием простейших предикатов (пороговых) можно построить очень хорошие (и, соответственно, сложные) модели. Также, модели, использующие сложные предикаты, сложны в обучении. Поэтому далее мы будем рассматривать только простейшие предикаты.

Опишем алгоритм построения дерева. Дерево строится жадно: на каждом шаге мы определяем наилучший предикат и разбиваем по нему обучающую выборку. Так мы продолжаем действовать до тех пор, пока не достигаем критерия останова. Его можно выбрать различными способами: можно ограничивать глубину дерева, количество листьев в нем, число объектов в листовой вершине и т. д. Обычно, выбор критерия останова не оказывает значительного влияния на качество модели (в отличие от величины ограничения, являющейся гиперпараметром). Далее прогнозы в листовых вершинах определяются как среднее арифметическое значений целевой переменной объектов обучающей выборки,

попавших в данную вершину: $c_v = \frac{1}{|R_v|} \sum_{(x_i, y_i) \in R_v} y_i$ (здесь мы вводим некоторые обозначения, которыми будем пользоваться в дальнейшем: R_v – это множество объектов в вершине v , c_v – прогноз модели для объектов листа v).

Теперь посмотрим, как определять наилучший предикат для каждой из вершин. Пусть на данный момент в рассматриваемой вершине содержится множество объектов R . Хаотичностью вершины будем называть дисперсию ответов в ней:

$$H(R) = \frac{1}{|R|} \sum_{(x_i, y_i) \in R} (y_i - y_R)^2, \text{ где } y_R = \frac{1}{|R|} \sum_{(x_i, y_i) \in R} y_i - \text{среднее арифметическое значений}$$

целевой переменной объектов множества R . Пусть каждый предикат разбивает множество R на два других множества – R_l и R_r (здесь l и r – индексы дочерних вершин). Тогда наилучшим будем считать предикат, который минимизирует взвешенную сумму хаотичностей в двух итоговых вершинах R_l и R_r : $Q(R, j, t) = \frac{|R_l|}{|R|} H(R_l) + \frac{|R_r|}{|R|} H(R_r) \rightarrow \min_{j, t}$,

где j – признак, t – порог разбиения.

Итак, мы рассмотрели алгоритм построения решающего дерева. Но у данной модели есть существенный недостаток – она очень неустойчива. То есть при незначительном изменении обучающей выборки итоговая модель может довольно сильно меняться. Для устранения данного недостатка используют композиции моделей, которые будут подробно рассмотрены в следующих главах.

2.4. Случайный лес

Для начала рассмотрим общие идеи, которые используют композиции моделей. Пусть у нас есть N базовых моделей $b_1(x), b_2(x) \dots b_N(x)$ (при этом к ним предъявляется формальное требование: вероятность их ошибки должна быть ниже 50%). Тогда прогноз композиции строится как среднее арифметическое прогнозов базовых моделей. При этом можно рассматривать разные принципы построения базовых моделей. Мы рассмотрим два основных варианта. Первый заключается в том, что мы будем обучать модели независимо, на разных подвыборках. Второй заключается в последовательном обучении моделей, причем строятся они так, чтобы каждая корректировала ошибки предыдущей. Первая идея лежит в основе принципа работы бэггинга и случайного леса и будет рассмотрена в этой главе, вторая же является основой бустинга и рассматривается в соответствующей главе.

Принцип работы бэггинга (bagging – bootstrap aggregating) основан на бутстрэпе. Метод бутстрэпа заключается в следующем. Пусть имеется выборка размера k . Возьмем из выборки k объектов с возвращением. Это означает, что мы будем k раз выбирать

произвольный объект выборки (считаем, что каждый объект «достаётся» с одинаковой вероятностью), причем каждый раз мы выбираем из всех исходных объектов. Отметим, что из-за возвращения среди выбранных объектов окажутся повторы. Далее мы обучаем базовые модели на подвыборках, построенных бутстрэпом.⁸

Заметим, что композиция, построенная таким образом, будет иметь такое же смещение, как и базовые модели. При этом её разброс зависит от ковариации базовых моделей (ковариация – это мера линейной зависимости двух случайных величин)⁹ и вычисляется следующим образом:

$$a(x) = \frac{\text{разброс}(b_n(x))}{n} + \text{ковариация}(b_n(x), b_m(x))$$

Итак, опишем алгоритм работы случайного леса. Для построения каждого из N деревьев с помощью бутстрэпа генерируется подвыборка и на ней обучается решающее дерево (по алгоритму, описанному ранее). При этом для уменьшения ковариации деревьев используется следующий прием: выбор оптимального предиката в очередной вершине очередного дерева осуществляется не по всем признакам, а по некоторому подмножеству множества признаков, которое генерируется отдельно для каждой вершины (рекомендуемый размер для этого подмножества – треть от исходного количества признаков). Для получения итогового прогноза модели необходимо усреднить прогнозы всех N деревьев.

Можно заметить, что с ростом числа деревьев N ошибка случайного леса выходит на асимптоту и не растёт. Поэтому случайный лес практически не содержит гиперпараметров, в следствие чего его часто называют универсальным алгоритмом машинного обучения. Но у бэггинга есть ряд проблем. Если базовая модель окажется смещенной, то и композиция в итоге окажется смещенной. Обучение и применение базовых моделей происходит довольно долго, также необходим большой объем памяти для хранения этих моделей.

2.5. Градиентный бустинг

Идея, как уже было сказано ранее, заключается в том, чтобы строить модели, составляющие композицию, так, чтобы каждая следующая максимально корректировала ошибки уже построенных моделей. Обычно в бустинге используются простые базовые модели (как мы увидим в дальнейшем, даже с помощью простых базовых моделей можно

⁸ Открытый курс машинного обучения. Тема 5. Композиции: бэггинг, случайный лес: [Электронный ресурс]// IT-сообщество «Хабр». URL: <https://habr.com/ru/company/ods/blog/324402/> (Дата обращения: 04.11.2021).

⁹ Ковариация – Википедия: [Электронный ресурс]// Википедия. URL: <https://ru.wikipedia.org/wiki/%D0%9A%D0%BE%D0%B2%D0%B0%D1%80%D0%B8%D0%B0%D1%86%D0%B8%D1%8F> (Дата обращения: 04.11.2021).

добиться хороших результатов). Итак, опишем жадный алгоритм построения данной композиции моделей:

1. Пусть N – число моделей в композиции, $b_1(x), b_2(x) \dots b_N(x)$ – сами модели. Итоговым прогнозом композиции моделей будем считать сумму прогнозов базовых моделей: $a_N(x) = \sum_{n=1}^N b_n(x)$. Первую модель мы уже умеем строить: $\frac{1}{l} \sum_{i=1}^l L(y_i, b_1(x_i)) \rightarrow \min_{b_1(x)}$ (здесь L – функция потерь);
2. Пусть мы уже построили $n-1$ модель, тогда модель с номером n должна строиться следующим образом: $\frac{1}{l} \sum_{i=1}^l L(y_i, a_{n-1}(x_i) + b_n(x_i)) \rightarrow \min_{b_n(x)}$. Но если используются недифференцируемые базовые модели, то пока мы не умеем обучать их на такой функционал;
3. Возьмем антипроизводную по второму аргументу функции потерь в точке, равной текущему значению прогноза композиции: $s_i^n = -\frac{dL(y_i, z)}{dz} \big|_{z=a_{n-1}(x_i)}$ (эти величины называют сдвигами). Тогда будем обучать новую базовую модель следующим образом: $\frac{1}{l} \sum_{i=1}^l (b_n(x_i) - s_i^n)^2 \rightarrow \min_{b_n(x)}$.

Теперь применим данный алгоритм для MSE (при подсчете сдвигов для простоты мы домножим производную на 0,5):

$$\text{Сдвиги: } s_i^n = -\frac{1}{2} * \frac{dL(y_i, z)}{dz} \big|_{z=a_{n-1}(x_i)} = -\frac{d}{dz} \frac{1}{2} (z - y_i)^2 \big|_{z=a_{n-1}(x_i)} = y_i - a_{n-1}(x_i);$$

$$\text{Итог: } \frac{1}{l} \sum_{i=1}^l (b_n(x_i) - (y_i - a_{n-1}(x_i)))^2 \rightarrow \min_{b_n(x)}.$$

Сделаем несколько замечаний, касающихся бустинга. Можно заметить, что при таком способе построения композиции моделей их смещение уменьшается, а разброс может увеличиваться. Поэтому в качестве базовых моделей следует брать модели, имеющие небольшой разброс, например, неглубокие решающие деревья. Гиперпараметрами данной композиции являются глубина базовых деревьев (можно считать, что глубиной мы регулируем их сложность, но в общем случае тут могут быть и другие варианты) и число деревьев (причем, в отличие от случайного леса, бустинг переобучается, если деревьев слишком много). Но у бустинга может возникнуть следующая проблема: вследствие своей простоты, базовые модели могут не справиться с приближением направления, указанного сдвигами. В качестве решения можно использовать регуляризацию: $a_n(x) = a_{n-1}(x) + \eta b_n(x)$, где $\eta \in (0; 1]$ – длина шага (чем меньше длина шага, тем больше деревьев потребуется для построения качественной композиции). Также для улучшения качества композиции можно для очередной базовой модели использовать рандомизацию по признакам (то есть

Ситкина А.Н. Прогноз продаж торговой точки на основе методов машинного обучения, 2021
обучать модель с использованием некоторого случайного подмножества признаков) и
объектам.

Глава 3. Используемые инструменты

3.1. Язык программирования Python

Python – это высокоуровневый язык программирования, который на данный момент является одним из самых популярных в мире. Он предоставляет эффективные высокоуровневые структуры данных, а также простой, но эффективный подход к объектно-ориентированному программированию. Его простой синтаксис и динамическая типизация наряду с тем, что он является легко читаемым, делают его идеальным языком для написания сценариев и быстрой разработки приложений в различных областях и на большинстве платформ.¹⁰

Также Python предоставляет множество инструментов для работы с данными и машинного обучения, из-за чего последнее время он стал одним из основных языков, используемых в этой области. По этой же причине он был выбран и нами.

Мы будем пользоваться последней версией Python 3, так как поддержка версии Python 2 недавно прекратилась.

3.2. Jupiter Notebook

Jupiter Notebook – это браузерный инструмент для интерактивного создания документов, сочетающих в себе пояснительный текст, математику, вычисления и их мультимедийный вывод.¹¹ Он является одним из самых популярных среди тех, кто занимается машинным обучением. За счет своей структуры он позволяет разбивать код на блоки и сразу запускать его, что обеспечивает быстрый и последовательный процесс разработки. Именно поэтому он был выбран как один из основных инструментов для выполнения работы.

3.3. Необходимые библиотеки

☐ Pandas

Pandas (обычно импортируется под названием `pd`) – это библиотека, используемая для анализа данных. Основными структурами данных являются Series (серия (колонка с

¹⁰ Документация языка Python. [Электронный ресурс]// Документация языка Python 3.10.0. URL: <https://docs.python.org/3/> (Дата обращения: 04.11.2021).

¹¹ Описание Jupiter Notebook. [Электронный ресурс]// Документация Jupiter Notebook. URL: <https://jupyter-notebook.readthedocs.io/en/latest/notebook.html> (Дата обращения: 04.11.2021).

данными)) и DataFrame (датафрейм (таблица с данными, состоит из серий)).¹² Это очень распространенная, продвинутая и быстроразвивающаяся библиотека, она предоставляет широкий набор инструментов для работы с данными (в основном с табличными), почему и была выбрана.

▣ NumPy

NumPy (обычно импортируется под названием np) – это библиотека, используемая для машинного обучения и научных вычислений. Она предоставляет широкий функционал для работы с многомерными массивами, в ней реализованы различные математические функции.¹³ Поэтому она и была выбрана.

▣ Datetime

Datetime – это модуль, содержащий широкий и удобный функционал для работы с информацией о датах и времени. В работе нами были использованы атрибут x.day, в котором записан день (x – дата, принадлежит классу date) и метод x.weekday(), возвращающий день недели для даты x (число от 0 до 6 соответственно).¹⁴

▣ Matplotlib и seaborn

Matplotlib – это библиотека визуализации. На ее основе работают функции визуализации в библиотеке pandas и вся библиотека seaborn, являющаяся зачастую более удобной, чем сама matplotlib.

Seaborn (обычно импортируется под названием sns), как было сказано ранее, также библиотека визуализации, основанная на matplotlib, но более удобная, чем matplotlib, почему и была выбрана.¹⁵

▣ Scikit-learn (sklearn)

Scikit-learn (sklearn) – это библиотека, в которой реализованы основные алгоритмы машинного обучения. Она является очень распространенным, продвинутым и удобным инструментом для реализации программ, решающих задачи машинного обучения, почему и была выбрана. Из этой библиотеки нами были взяты классы моделей машинного обучения, функции метрик а также некоторые другие инструменты.¹⁶

¹² Документация библиотеки Pandas. [Электронный ресурс]// Официальный сайт библиотеки Pandas. URL: https://pandas.pydata.org/docs/user_guide/index.html (Дата обращения: 04.11.2021).

¹³ Документация библиотеки NumPy. [Электронный ресурс]// Официальный сайт библиотеки NumPy. URL: <https://numpy.org/doc/stable/> (Дата обращения: 04.11.2021).

¹⁴ Документация языка Python. [Электронный ресурс]// Документация языка Python 3.10.0. URL: <https://docs.python.org/3/> (Дата обращения: 04.11.2021).

¹⁵ Документация библиотеки Seaborn. [Электронный ресурс]// Официальный сайт библиотеки Seaborn. URL: <https://seaborn.pydata.org> (Дата обращения: 04.11.2021).

¹⁶ Документация библиотеки Scikit-learn. [Электронный ресурс]// Официальный сайт библиотеки Scikit-learn. URL: https://www.sklearn.org/user_guide.html (Дата обращения: 04.11.2021).

☐ LightGBM

LightGBM – это одна из самых популярных современных библиотек, в которых реализован градиентный бустинг. У неё есть ряд преимуществ: высокая скорость обучения, снижение использования памяти, улучшенная точность и т. д. Поэтому она и была выбрана нами.¹⁷

☐ CatBoost

CatBoost – это алгоритм для градиентного бустинга на деревьях решений. Он разработан исследователями и инженерами Яндекса. Его отличительными чертами являются: хорошее качество работы с гиперпараметрами, используемыми по умолчанию, повышенная точность и быстрое прогнозирование.¹⁸ В этой библиотеке используется особый вид решающих деревьев – oblivious decision trees (симметричные решающие деревья). Главная их особенность заключается в том, что на каждом их уровне используется один и тот же предикат, что значительно ускоряет время обучения моделей.¹⁹

☐ tkinter

tkinter – это графическая библиотека, являющаяся стандартным интерфейсом Python для инструментария Tcl/Tk GUI.²⁰ Главное её преимущество заключается в том, что она является кроссплатформенной, почему и была выбрана нами.

¹⁷ Документация библиотеки LightGBM. [Электронный ресурс] // Официальный сайт библиотеки LightGBM. URL: <https://lightgbm.readthedocs.io/en/latest/index.html> (Дата обращения: 04.11.2021).

¹⁸ Документация библиотеки CatBoost. [Электронный ресурс] // Официальный сайт библиотеки CatBoost. URL: <https://catboost.ai/> (Дата обращения: 04.11.2021).

¹⁹ Соколов Е.А., Зимовнов А.В., Ковалев Е.И., Кохтев В.М., Рысьмятова А.А., Филатов А.А. Основы машинного обучения: [Электронный ресурс] // Образовательная платформа «Открытое образование». URL: https://courses.openedu.ru/courses/course-v1:hse+INTRML+fall_2020/course/ (Дата обращения: 04.11.2021).

²⁰ Документация языка Python. [Электронный ресурс] // Документация языка Python 3.10.0. URL: <https://docs.python.org/3/> (Дата обращения: 04.11.2021).

Глава 4. Реализация программы

В данном разделе в общем виде разбираются принципы работы нашей программы. Более подробные комментарии вместе с кодом программы находятся в Приложениях 1 и 2.

1. Исходные данные

После подключения необходимых библиотек происходит загрузка данных. Они записаны в файл с расширением .xls. В нем содержится информация о продажах товаров из категорий «Молочная продукция» и «Хлеб и хлебобулочные изделия» (взят несезонный товар, чтобы данная выборка была репрезентативна²¹) продуктового магазина в поселке Прибрежный за ноябрь 2020 года (см. Рис. 1).

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Наименование	Количество	Сумма продаж	Прибыль	Дата										
Аренка домашняя 250гр	2,000	130.00	46.00	02.11.2020										
Аццедофени 2.7% 200г БЗМК	4,000	71.81	9.85	02.11.2020										
Батон ГОРЧАЧНЫЙ 300 г. Хлебозавод №1 +	21,000	419.90	81.38	02.11.2020										
Батон Луковый 300г Орион	2,000	57.93	11.93	02.11.2020										
Батон НАРЕЗНОЙ 400 г. Хлебозавод №1 +	41,000	901.19	192.71	02.11.2020										
Батон Нарезной 400г Переславский х-д	5,000	164.98	38.48	02.11.2020										
Батон НАРЕЗНОЙ НАРЕЗКА 400 г. Хлебозавод №1 +	17,000	407.48	89.92	02.11.2020										
Батон Нарезной новый в нарезке в/с 0.2 кг	3,000	68.81	14.21	02.11.2020										
Батон Подмоховский Новый 370г Переславский х-д	5,000	154.92	32.42	02.11.2020										
Батон СОЛНЕЧНЫЙ в/с 300г Хлебозавод №1	17,000	339.77	74.23	02.11.2020										
Батон Творческий 300г Переславский х-д	5,000	175.00	32.50	02.11.2020										
Батончик к чаю в/с 200г Хлебозавод №1	9,000	152.94	35.58	02.11.2020										
Батончик к чаю 300г Орион	3,000	71.97	17.97	02.11.2020										
Беляш 100гр. Орион	10,000	489.88	119.88	02.11.2020										
Бифидок 3.2% 500 г Яромоллод БЗМК	2,000	78.00	13.22	02.11.2020										
Булка Городская в/с 200г Хлебозавод №1	9,000	134.75	26.66	02.11.2020										
Булочка в Дорогу 200г Переславский х-д	4,000	83.91	17.91	02.11.2020										
Булочка Веснушечка 4 шт/0.05 Переславль	2,000	77.59	14.59	02.11.2020										
Булочка для Гамбургеров 150г Переславский х-д	2,000	46.00	13.20	02.11.2020										
Ватрушка Венгерская 90г Орион	5,000	174.74	39.74	02.11.2020										
Ватрушка с творогом 75г Орион	7,000	181.77	41.77	02.11.2020										
Йогурт Активна натуральный клубника 2.9 % 150г БЗМК	1,000	48.00	9.64	02.11.2020										
Йогурт Даниссмо Фантазия с хруст. шариками и ягодным вкусом 6.9% БЗМК100г +	1,000	64.98	14.89	02.11.2020										
Йогурт клубника-малина 1.5% 500 г Яромоллод БЗМК	4,000	171.74	33.26	02.11.2020										

Рис. 1 Исходные данные

2. Предобработка

Далее мы выводим некоторые характеристики наших данных (розничная цена, день недели), создаем новые колонки (они содержат данные, которые в дальнейшем мы будем использовать как признаки), вручную обрабатываем признак category, так как иначе (из-за большого количества возможных значений) есть риск переобучения модели.

3. Визуализация

Визуализируем распределение целевой переменной и ее зависимости от признаков (чтобы убедиться, что данные зависимости действительно наблюдаются). В нашей работе были визуализированы распределение целевой переменной quantity, зависимость

²¹ Репрезентативная выборка – это выборка, с хорошей точностью моделирующая генеральную совокупность (множество всех рассматриваемых объектов).

количества продаж товара (quantity) от его розничной цены (price), график зависимости числа продаж (quantity) по дням (day), гистограмма зависимости числа продаж (quantity) от дня недели (week_day) и гистограмма зависимости числа продаж (quantity) от категории товара (bread, hotcake, milk_or_sour_milk, cheese_cottage_cheese, others). Так мы смогли убедиться в наличии зависимостей между признаками и целевой переменной.

4. Подготовка данных, разбиение на обучающую и тестовую выборки

Подготовим данные: уберем ненужные колонки (содержащие излишнюю информацию: целевую переменную и данные, напрямую с ней связанные, излишне подробные категории), применим one-hot кодирование для категориального признака week_day. Отделим колонку с целевой переменной, разобьем данные на обучающую и тестовую выборки (мы делаем это вручную, так как в данных как признак используется день недели; иначе при стандартном разбиении 80/20 или 70/30 в тестовую выборку может попасть либо очень маленькое количество объектов, у которых week_day=0, либо они могут не попасть вовсе; мы делаем следующее разбиение: первые три недели месяца – обучающая выборка, последняя неделя - тренировочная).

5. Обучение и сравнение моделей

Далее мы переходим к обучению моделей. Нами были обучены следующие модели: линейная регрессия (в том числе с регуляризацией, с использованием L-1 нормы и с сочетанием L-1 и L-2 норм, с использованием реализованного нами градиентного спуска), kNN, случайный лес и градиентный бустинг (здесь нами было проведено сравнение трех его реализаций в разных библиотеках: sklearn, CatBoost, LightGBM), а также были рассмотрены блендинг (среднее взвешенное) бустинга и случайного леса, блендинг бустинга и kNN. Ниже приводится таблица результатов, которые показывает каждая из моделей (с наилучшими гиперпараметрами).

Табл. 1 Сравнение моделей

Модель	MSE	MAE
Линейная регрессия	21.918	2.678
Гребневая регрессия	21.868	2.671
Lasso	22.562	2.694
ElasticNet	22.421	2.710
Вручную реализованный градиентный спуск	33.108	3.483
kNN	20.283	2.292
Случайный лес	17.681	2.194
Бустинг (CatBoostRegressor)	15.908	2.202
Бустинг и случайный лес	15.669	2.117

Бустинг и kNN

15.615

2.082

Таким образом, в качестве наилучшей модели нами был выбран блендинг градиентного бустинга и kNN.

6. Создание графического интерфейса программы

Кратко опишем функционал созданной нами программы и её интерфейса. При открытии окна программы (после закрытия приветственного окна, см. Рис. 2) можно увидеть, что можно ввести данные о товаре и увидеть прогноз нашей модели (см. Рис. 3). Также программа умеет обрабатывать файлы (список допустимых форматов можно увидеть, кликнув по кнопке «Помощь», там же можно увидеть и некоторые другие инструкции по работе с программой). Для создания графического интерфейса программы была использована библиотека tkinter. С полным кодом программы можно ознакомиться в

Приложении 2.

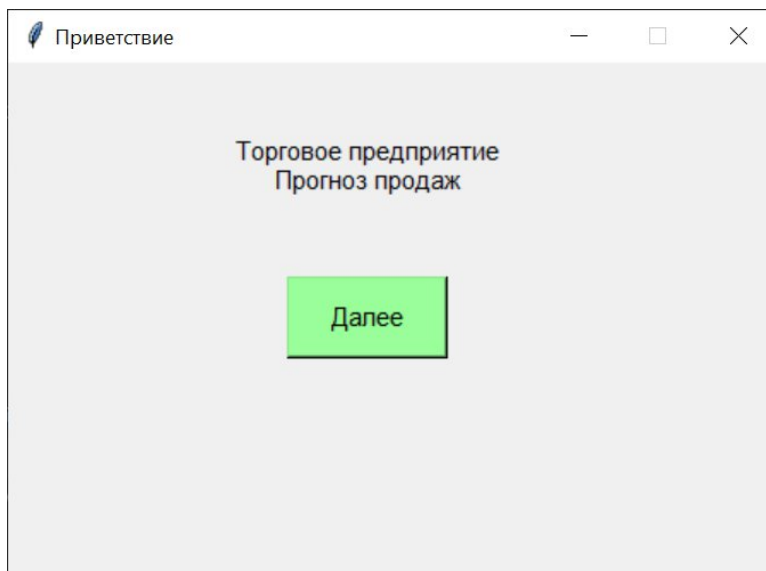


Рис. 2 Окно приветствия

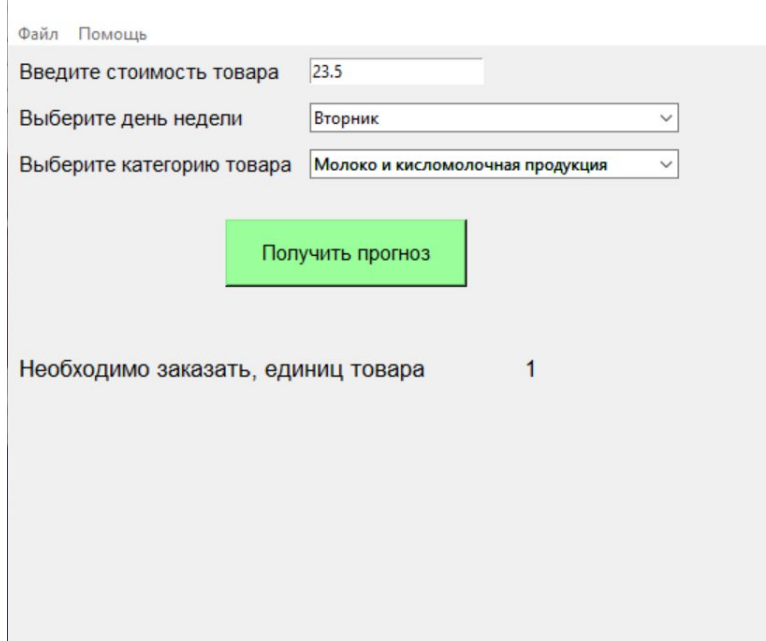


Рис. 3 Основное окно программы и демонстрация её работы

Заключение

В ходе данной исследовательской работы мы изучили основные понятия, лежащие в основе методов машинного обучения: виды задач, решаемых данными алгоритмами, проблемы переобучения и недообучения, понятия метрики и кросс-валидации. Далее, нами были исследованы различные модели машинного обучения: метод kNN (метод k ближайших соседей), линейная регрессия (составление линейного уравнения с весами), решающие деревья, случайный лес и градиентный бустинг (композиции решающих деревьев). Также нами были исследованы возможности языка Python и его дополнительных библиотек: matplotlib и Seaborn (визуализация данных), Pandas (работа с данными), NumPy (вычисления и линейная алгебра), sklearn (библиотека машинного обучения), datetime (работа с информацией о датах и времени) и tkinter (библиотека для создания графического интерфейса). Также были изучены библиотеки, в которых реализован градиентный бустинг, а именно CatBoost и LightGBM. В результате сравнения различных моделей было выявлено, что наилучший результат на наших данных показывает блендинг (среднее взвешенное) градиентного бустинга и kNN. После этого нами был реализован графический интерфейс для взаимодействия пользователя с программой (за основу взята выявленная ранее наилучшая модель).

Можно считать, что цель, как и все поставленные задачи, была выполнена. Созданный нами продукт решает проблему доступа предприятий малого бизнеса к программам, осуществляющим прогноз продаж. Также, наша программа может стать помощником для работников, занимающихся закупкой товара. В этом заключается практическая значимость нашей работы.

Список использованных источников и литературы

1. Документация библиотеки CatBoost. [Электронный ресурс] // Официальный сайт библиотеки CatBoost. URL: <https://catboost.ai/> (Дата обращения: 04.11.2021).
2. Документация библиотеки LightGBM. [Электронный ресурс] // Официальный сайт библиотеки LightGBM. URL: <https://lightgbm.readthedocs.io/en/latest/index.html> (Дата обращения: 04.11.2021).
3. Документация библиотеки NumPy. [Электронный ресурс] // Официальный сайт библиотеки NumPy. URL: <https://numpy.org/doc/stable/> (Дата обращения: 04.11.2021).
4. Документация библиотеки Pandas. [Электронный ресурс] // Официальный сайт библиотеки Pandas. URL: https://pandas.pydata.org/docs/user_guide/index.html (Дата обращения: 04.11.2021).
5. Документация библиотеки Scikit-learn. [Электронный ресурс] // Официальный сайт библиотеки Scikit-learn. URL: <https://scikit-learn.org/stable/index.html> (Дата обращения: 04.11.2021).
6. Документация библиотеки Seaborn. [Электронный ресурс] // Официальный сайт библиотеки Seaborn. URL: <https://seaborn.pydata.org> (Дата обращения: 04.11.2021).
7. Документация языка Python. [Электронный ресурс] // Документация языка Python 3.10.0. URL: <https://docs.python.org/3/> (Дата обращения: 04.11.2021).
8. Ковариация – Википедия: [Электронный ресурс]// Википедия. URL: <https://ru.wikipedia.org/wiki/%D0%9A%D0%BE%D0%B2%D0%B0%D1%80%D0%B8%D0%B0%D1%86%D0%B8%D1%8F> (Дата обращения: 04.11.2021).
9. Описание Jupiter Notebook. [Электронный ресурс] // Документация Jupiter Notebook. URL: <https://jupyter-notebook.readthedocs.io/en/latest/notebook.html> (Дата обращения: 05.01.2021).
10. Открытый курс машинного обучения. Тема 5. Композиции: бэггинг, случайный лес: [Электронный ресурс]// IT-сообщество «Хабр». URL: <https://habr.com/ru/company/ods/blog/324402/> (Дата обращения: 04.11.2021).
11. Рашка С., Мирджалили В. Python и машинное обучение: машинное и глубокое обучение с использованием Python, scikit-learn и TensorFlow 2/ Пер. с англ. СПб.: ООО "Диалектика", 2020. 848 с.
12. Соколов Е.А., Зимовнов А.В., Ковалев Е.И., Кохтев В.М., Рысьмятова А.А., Филатов А.А. Основы машинного обучения: [Электронный ресурс]// Образовательная

Ситкина А.Н. Прогноз продаж торговой точки на основе методов машинного обучения, 2021
платформа «Открытое образование». URL: https://courses.openedu.ru/courses/course-v1:hse+INTRML+fall_2020/course/ (Дата обращения: 04.11.2021).