

МИНОБРАЗОВАНИЯ РОССИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ
ВЫСШЕГО ОБРАЗОВАНИЯ
«ВОРОНЕЖСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ»

Факультет компьютерных наук

Кафедра цифровых технологий

Вероятностный подход к обработке последовательных данных

ВКР Бакалаврская работа

02.03.01 Математика и компьютерные науки

Распределенные системы и искусственный интеллект

Допущено к защите в ГЭК _____.2022__

Зав. кафедрой _____ *С. Д. Кургалин, д. ф.-м. н., профессор*
подпись

Обучающийся _____ *А. М. Гузенко, 4 курс, д/о*
подпись

Руководитель _____ *А. Ф. Клинских, д. ф.-м. н., профессор*
подпись *расшифровка подписи, ученая степень, звание, должность*

Воронеж 2022

Содержание

Введение

Применение анализа временных рядов и методов машинного обучения играют огромную роль в жизни человека и делает большой шаг в технологическом прогрессе. Машинное обучение применяется во многих сферах, от сетевых магазинов до промышленных производств. Где-то они носят рекомендательный характер для человека, а в каких-то уже могут заменить человека полностью. В пример можно привести систему автоматического торможения от Volvo, учитывая дистанцию и скорость, при помощи передней камеры, в режиме реального времени рассчитывая тормозной путь, система может принять решение об экстренном торможении, с реакцией, которая не сравнится с человеческой.

Анализ временных рядов и машинное обучение получили свое быстрое развитие в основном по двум причинам: увеличение количества информации для анализа и развитие вычислительных мощностей, но всему есть предел. Если для информации одна из основных проблем это правильный сбор и последующее интерпретирование информации, то для машинного обучения **актуален** вопрос эффективности алгоритмов. Для решения данной проблемы используются как другие вычислительные модели, например, квантовые вычисления, так и поиск более эффективных алгоритмов для обучения модели. Один из таких алгоритмов является вероятностный подход, или же Байесовский подход, названный в честь Томаса Байеса, который разработал одну из важнейших теорем в теории вероятности, теорему Байеса о зависимых вероятностях.

В последнее время вероятностных подход нашел отклик и стремительно развивается, созданы несколько библиотек для создания моделей, анализа данных. Решаются проблемы эффективности и точности моделей, созданием улучшенных методов для работы с вероятностными распределениями.

Главные **задачи** в данной работе:

1. Взять временной ряд, в котором требуется решить задачу классификации по одному или нескольким переменным.
2. Реализовать две модели для решения задач классификации.
3. Сравнить обычный подход в решении задачи классификации и вероятностный подход в эффективности на разных выборках.
4. Сравнить скорости обучения моделей на равном объеме данных.
5. Сделать выводы об областях применения данного подхода, его преимуществах и недостатках.

Литературный обзор

В первую очередь стоит отметить, что рассматриваемая область быстроразвивающаяся, и, соответственно, источники для исследования в данной области должны быть как можно более актуальными, но не стоит забывать о качестве этих источников, так как некоторые поспешные выводы могут быть ошибочны, но данный тезис в большей мере касается научных статей, результатов недавно проведенных исследований. Основопологающим для данной работы я выбрал три литературных источника, остальные же либо в какой-то мере дополняют их, либо более углубленно раскрывают некоторые темы.

Для анализа временных рядов за основу была взята книга «Практический анализ временных рядов: Прогнозирование со статистикой и машинное обучение» Эйлин Нильсен 2019 года. Эйлин кандидат прикладной физики Колумбийского университета, занималась исследованиями влияния развития технологий на право и медицину. В данный момент она занимается разработкой нейронной сети для прогнозирования в сфере финансов. В этой книге предлагается практическое знакомство с временными рядами, их анализом, обработкой и последующим применением.

Для Байесовского подхода за основу были взяты книги «Байесовские модели. Байесовская статистика на языке Python» Аллен Б. Дауни 2013 года и «Байесовский анализ на Python: введение в статистическое моделирование и вероятностное программирование с использованием PyMC3 и ArviZ» Освальдо Мартин 2018 года.

Аллен Б. Дауни профессор компьютерных наук в инженерном колледже штата Массачусетс, выпустил множество книг по обработке данных и программированию в целом. Книга по Байесовской статистике является руководством к применению языка программирования Python для данной сферы, без углубления в математические основы.

Освальдо Мартин ученый-исследователь агентства National Scientific and Technical Research Council (CONICET) в Аргентине, работающий в области структурной биоинформатики, является одним из основных разработчиков библиотек PyMC3 и ArviZ, преподаватель курса по Байесовскому анализу данных. Книга представляет собой введение в Байесовский анализ, излагается методический подход моделирования вероятностных моделей, применение теоремы Байеса для вывода логических следствий из используемых моделей и данных.

Дополнительно использовались следующие источники.

Для понимания математических процессов, то есть теории вероятности и математической статистики, идеально подойдут работы двух советских математиков, а именно «Курс теории вероятностей» Бориса Владимировича Гнеденко и «Теория вероятностей и математическая статистика» Владимира Ефимовича Гмурмана. Данные две книги дают систематические знания по теории вероятностей и математической статистике, предоставляют разбор реальных примеров и задач для самостоятельного решения.

Для понимания общего устройства и работы нейросетей подойдут некоторые разделы следующих книг «Глубокое обучение с точки зрения практика» Джош Паттерсон, «Pattern Recognition and Machine Learning» Christopher Michael Bishop, «Building Machine Learning and Deep Learning Models on Google Cloud Platform» Ekaba Bisong, «Глубокое обучение» Гудфеллоу Ян, «Прикладное глубокое обучение: подход к пониманию глубоких нейронных сетей на основе метода кейсов» Микелуччи Умберто.

Для лучшего математического понимания Байесовской статистики была использована книга «Bayesian Data Analysis» Andrew Gelman. Так же, данный автор описал в научном журнале алгоритм, который мы будем рассматривать в дальнейшем.

Для рассмотрения нового алгоритма Automatic Differentiation Variational Inference взята статья Alp Kucukelbir из научного журнал

Глава 1. Последовательные данные

1.1 Понятие последовательных данных

Последовательные данные (или временные ряды) – это серия точек данных, проиндексированных во временном или ином порядке. Чаще всего временной ряд представляет собой множество, полученное в последовательных равностоящих точках времени.

Данные временных рядов и их анализ приобретают все большее значение из-за получения таких данных посредством, например, интернет вещей, цифровизации здравоохранения, развития умных городов и так далее. По мере того, как непрерывный мониторинг и сбор данных становятся все более распространёнными, потребность в анализе временных рядов с использованием как статистических методов, так и методов машинного обучения, повысится. Анализ временных рядов часто сводится к вопросу о причинно-следственной связи, как прошлое повлияло на будущее.

Практическое применение анализа временных рядов и машинного обучения еще в 1980-ых годах включал широкий спектр сценариев:

- Специалисты по компьютерной безопасности задействовали их для выявления аномалий в качестве метода идентификации против взломов и вмешательств.
- Динамическое изменение масштаба времени, один из доминирующих методов измерения сходства между временными рядами. С увеличением вычислительной мощности стало возможно достаточно быстро вычислять «расстояние» между, например, аудиозаписями.
- Изобретены рекурсивные нейронные сети, которые показали свою эффективность для извлечения шаблонов из поврежденных данных.

Для анализа временных рядов используются различные графические отображения данных, такие как графики, сводные статистики, гистограммы, диаграммы рассеяния. Рассмотрим некоторые из них.

Во временных рядах, как и в данных, можно произвести операцию группировки, то есть выделить отдельные группы значений. Например, в перекрестных данных операции группировки позволяют определить средние значения для значений возраста, пола или места проживания респондентов. В анализе временных рядов так же востребованы групповые операции, позволяющие рассчитывать статистические показатели наборов, например, средимесячные или недельные медианы. Данные одного временного ряда становятся более понятными при разделении на несколько параллельных временных рядов.

Возьмем сведения о ежедневных ценах закрытия четырех основных европейских фондовых индексов, указанных с 1991 по 1998 годам, и отобразим их на графике.

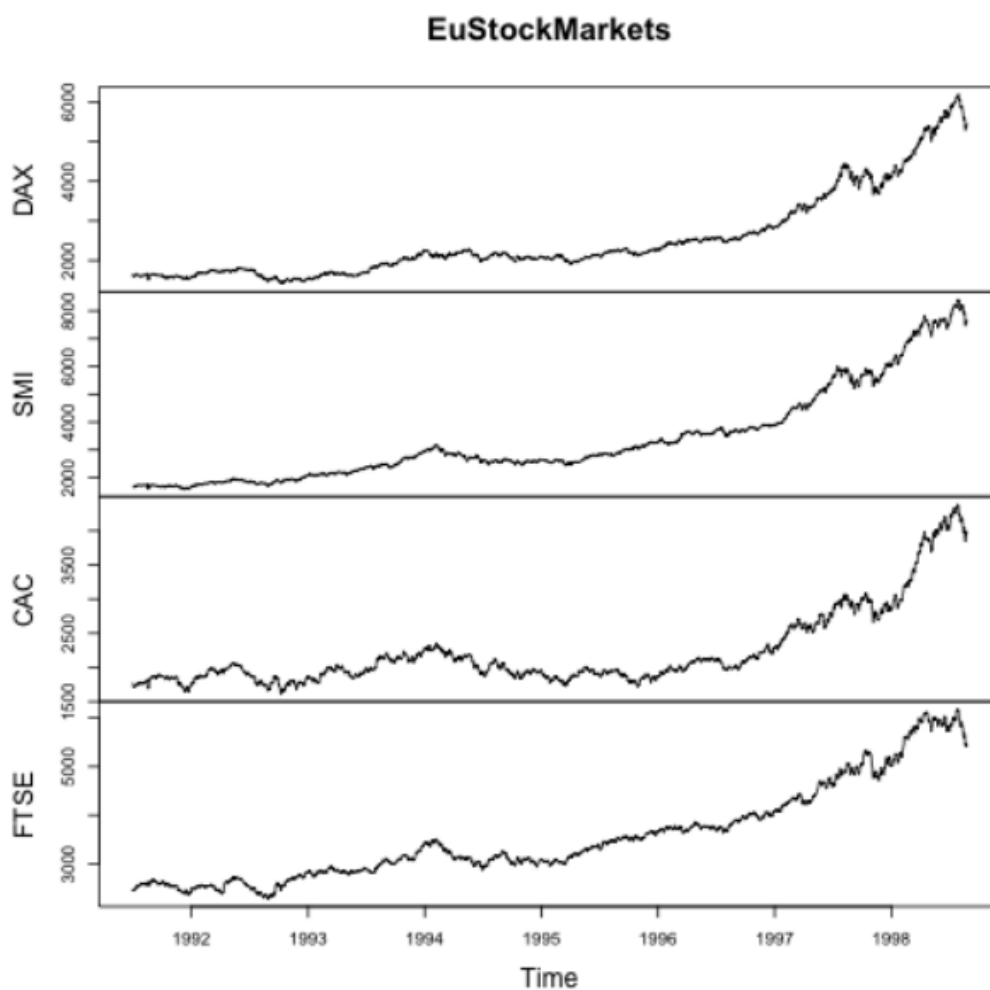


Рисунок 1 - График временного ряда

При помощи гистограммы мы можем оценить частоту появления данных, или, что более важно, проанализировать изменение значения во времени.

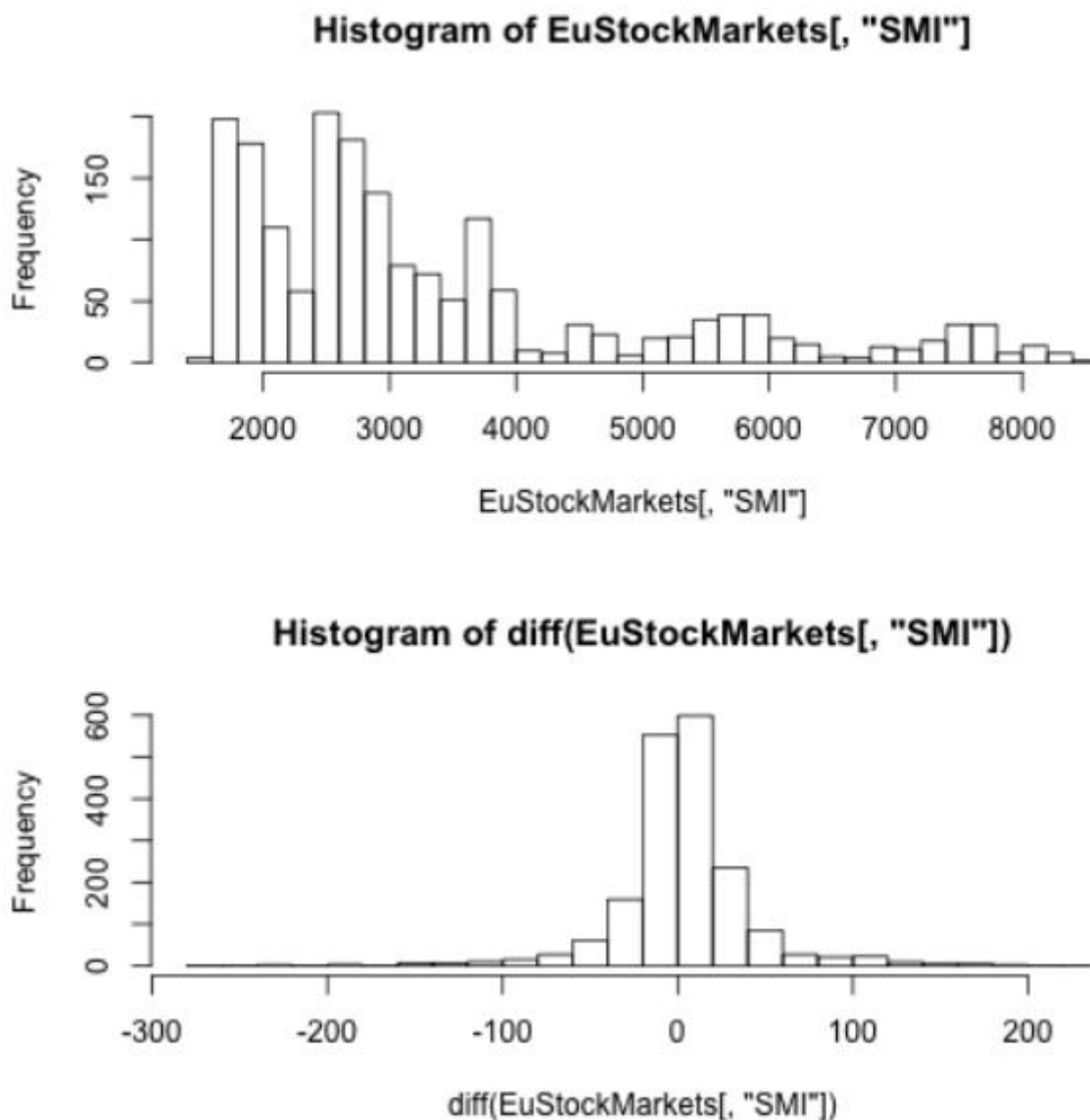


Рисунок 2 - Гистограмма временного ряда. Первая гистограмма – частота значений, вторая гистограмма – частота разностей

Диаграммы рассеяния мы можем использовать для определения взаимосвязи между ценами двух акций в отдельные моменты времени, а также отслеживания их временных изменений. На следующем рисунке присутствуют оба варианта.

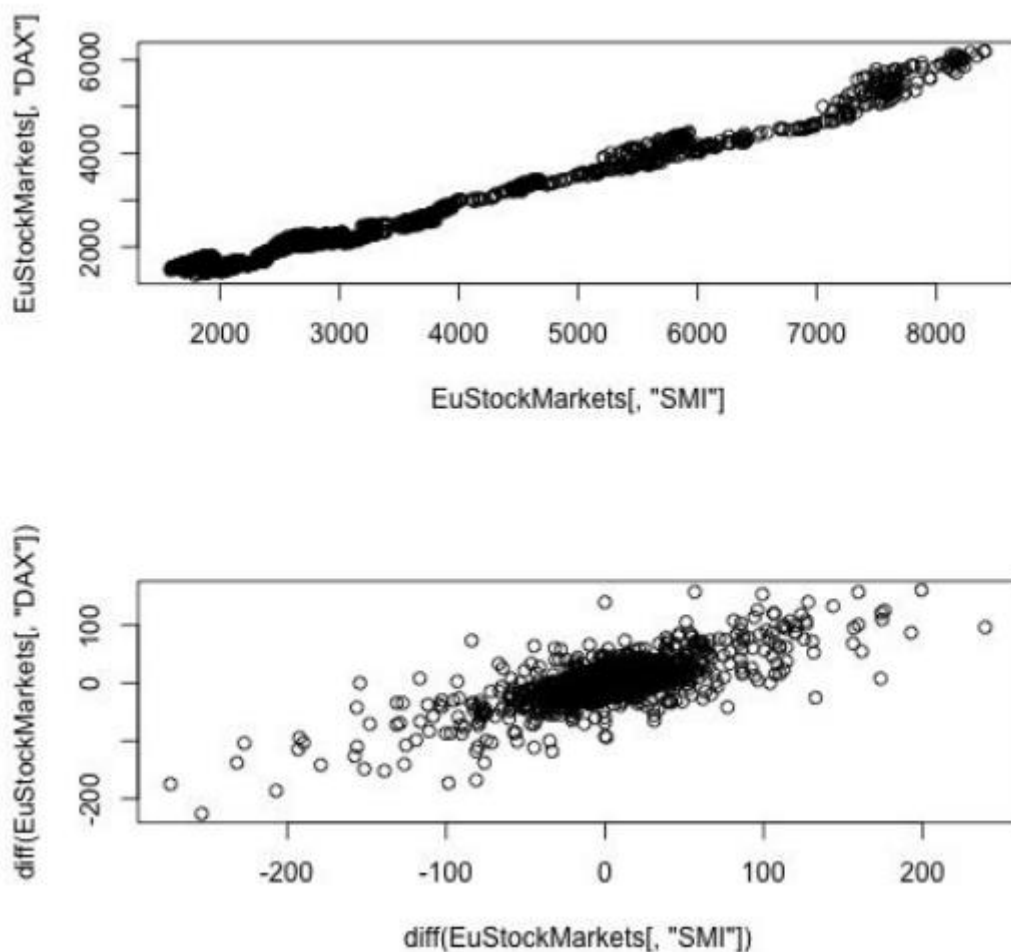


Рисунок 3 – Диаграмма рассеяния. На первой диаграмме – цена двух акций с течением времени, на второй – разница цен двух акций с течением времени

Помимо графического анализа стоит обращать внимание на некоторые характеристики временных рядов, такие как стационарность, корреляция, ложная корреляция.

Один из важных вопросов при анализе временного ряда – это какую систему он описывает, стабильную или изменчивую? Уровень стабильности, или стационарности, важен для оценки того, как долгосрочное поведение в прошлом отражает ее поведение в будущем. После оценки уровня стабильности мы пытаемся определить, присуще ли ему динамическое поведение. Простое определение стационарности процесса заключается в следующем: процесс считается стационарным, если для всех возможных смещений k распределение $y_t, y_{t+1}, \dots, y_{t+k}$ не зависит от t .

Самокорреляция – значение временного ряда в отдельные моменты времени могут коррелировать с его значениями в другие моменты времени.

Автокорреляция – общий случай самокорреляции, лишенному привязки к конкретному моменту времени. Автокорреляция сводится к решению задачи поиска взаимосвязи между любыми двумя точками общего временного ряда, расположенными на строго заданном расстоянии друг от друга, то есть автокорреляция дает представление о линейной взаимосвязи точек данных, полученных в разные моменты времени, как о функции разницы времени их получения.

Ложные корреляции – математическая зависимость, в которой переменные или события связаны, но в следствии совпадения. Ложные корреляции остаются важной проблемой, требующей самого пристального изучения и опровержения.

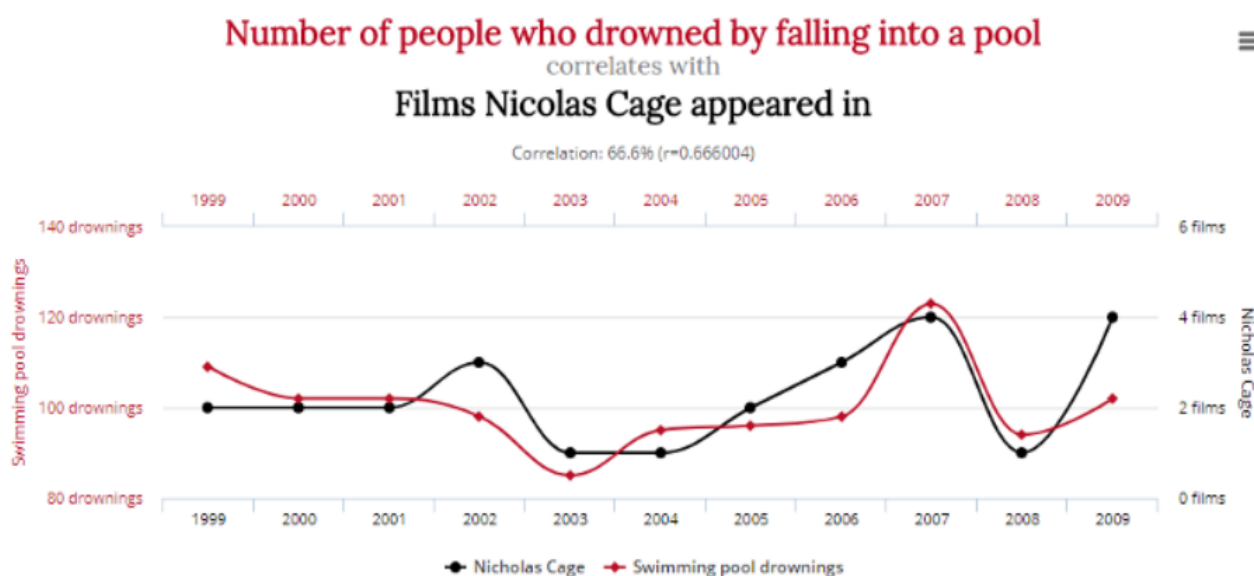


Рисунок 4 – График выхода фильмов с Николасом Кейджем и количество утонувших людей в бассейне

Моделирование данных выступает разновидностью анализа данных, который находит широкое применение при работе с временными рядами. Это следует из одного недостатка временных рядов: никакие две точки данных в одном и том же временном ряду не могут быть точно сопоставимы, поскольку они относятся к разному времени. Как только мы задумаемся о

том, что могло бы произойти в данный момент времени, мы придём к моделированию.

Моделирование временных рядов методами глубокого обучения – это новая, но весьма многообещающая дисциплина науки о данных. С помощью технологии глубокого обучения во временных рядах можно обнаружить сложные динамические процессы, зачастую скрытые для понимания даже опытных специалистов по анализу данных.

Моделирование и прогнозирование – схожие задачи. В обоих случаях сначала нужно сформулировать гипотезу о параметрах и поведении базовой системы, а затем экстраполировать имеющиеся данные для получения новых точек.

Но нужно четко понимать различия между моделированием и прогнозированием:

- Иногда качественные наблюдения проще обрабатывать методами моделирования, а не прогнозирования.
- Моделирование выполняется в определённом масштабе, что позволяет увидеть множество альтернативных сценариев, в то время как прогнозы составляются предельно точно.

1. 2 Вероятностный подход к анализу последовательных данных.

Для начала стоит напомнить, что такое вероятность в общем понимании. Вероятность – это число между 0 и 1, которые представляют собой уровень уверенности, что некоторый факт справедлив. Числом 1 представляется абсолютная уверенность, что некоторый факт справедлив. Числом 0 – абсолютная уверенность, что этот факт не справедлив. Промежуточные числа в этом интервале определяют степень уверенности. Число 0,5 означает, что предсказанное событие в одинаковой степени может как осуществиться, так и не осуществиться.

Распределение вероятностей – это математический объект, который описывает, насколько возможными являются различные события. Эти

события ограничены каким-либо образом, то есть представляют собой набор возможных событий, например набор возможных чисел $\{1, 2, 3, 4, 5, 6\}$ для игральных костей (исключая неопределенные случаи). Общепринятым и полезным представлением концепции в статистике является интерпретация данных как генерируемых из некоторого истинного распределения вероятностей с неизвестными параметрами. То есть нам нужно найти значения этих параметров с использованием только частичной выборки из истинного распределения вероятностей. В обобщенном случае у нас нет доступа к такому истинному распределению вероятностей, поэтому необходимо создать модель, чтобы попытаться аппроксимировать это распределение. Вероятностные модели создаются при помощи правильного объединения распределений вероятностей.

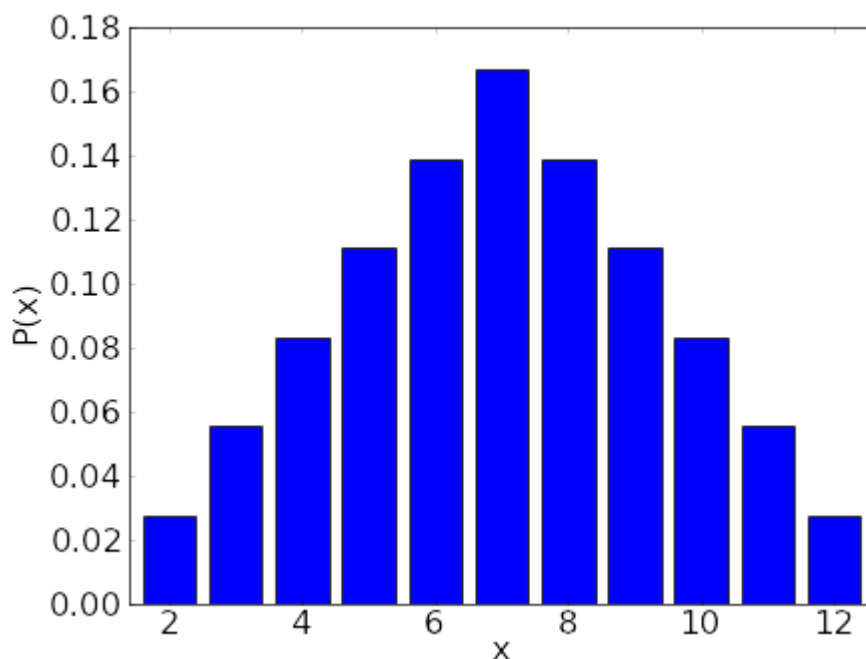


Рисунок 5 – Вероятностное распределение броска двух игральных костей

Для построения сложных моделей активно используется теорема Байеса

$$p(\theta, y) = \frac{p(y|\theta)p(\theta)}{p(y)} \quad (1.1)$$

Если заменить элемент θ на «предположение» (гипотезу), а элемент y на «данные», то теорема Байеса показывает, как вычислить вероятность предположения θ при наличии данных y . Для того, чтобы превратить предположение в некоторый объект, нужно использовать вероятностные распределения. В вероятностном подходе для анализа данных мы часто будем использовать теорему Байеса. Ниже приведем названия для элементов данной теоремы:

- $p(\theta)$ – априорная вероятность
- $p(y|\theta)$ – правдоподобие
- $p(\theta|y)$ – апостериорная вероятность
- $p(y)$ – предельное правдоподобие

Априорная вероятность должна соответствовать тому, что нам известно о значении параметра θ перед рассмотрением данных y . Если нам ничего не известно, то можно использовать постоянные фиксированные априорные вероятности, которые не содержат какого-либо значимого объема информации.

Правдоподобие определяет, как будут представлены данные в дальнейшем анализе. Это выражение правдоподобности данных с учетом принятых параметров.

Апостериорная вероятность – это результат байесовского анализа, которые отображает все, что известно о задаче (проблеме) с учетом имеющихся данных и используемой модели. Апостериорная вероятность – это распределение вероятностей для параметра θ в используемой модели. Такое распределение это баланс между априорной вероятностью и правдоподобием. С теоретической концептуальной точки зрения апостериорную вероятность можно воспринимать как обновленную вероятность в свете новых данных. Апостериорная вероятность, полученная в результате одного процесса анализа, может использоваться как априорная вероятность для нового процесса анализа. Это свойство делает байесовский

анализ особенно подходящим для анализа данных, которые становятся доступными в определенном последовательном порядке. Примерами могут служить системы раннего оповещения о природных катастрофах, которые обрабатывают в режиме онлайн данные, поступающие с метеорологических станций и спутников.

Правдоподобие – это вероятность исследуемых данных, усредненная по всем возможным значениям, которые могут принимать параметры. В любом случае мы не уделяем особого внимания предельному правдоподобию и будем считать его простым фактором нормализации. Такой подход принят потому, что при анализе распределения апостериорной вероятности нас будут интересовать только относительные, а не абсолютные значения параметров.

$$p(\theta, y) \propto p(y|\theta)p(\theta) \quad (1.2)$$

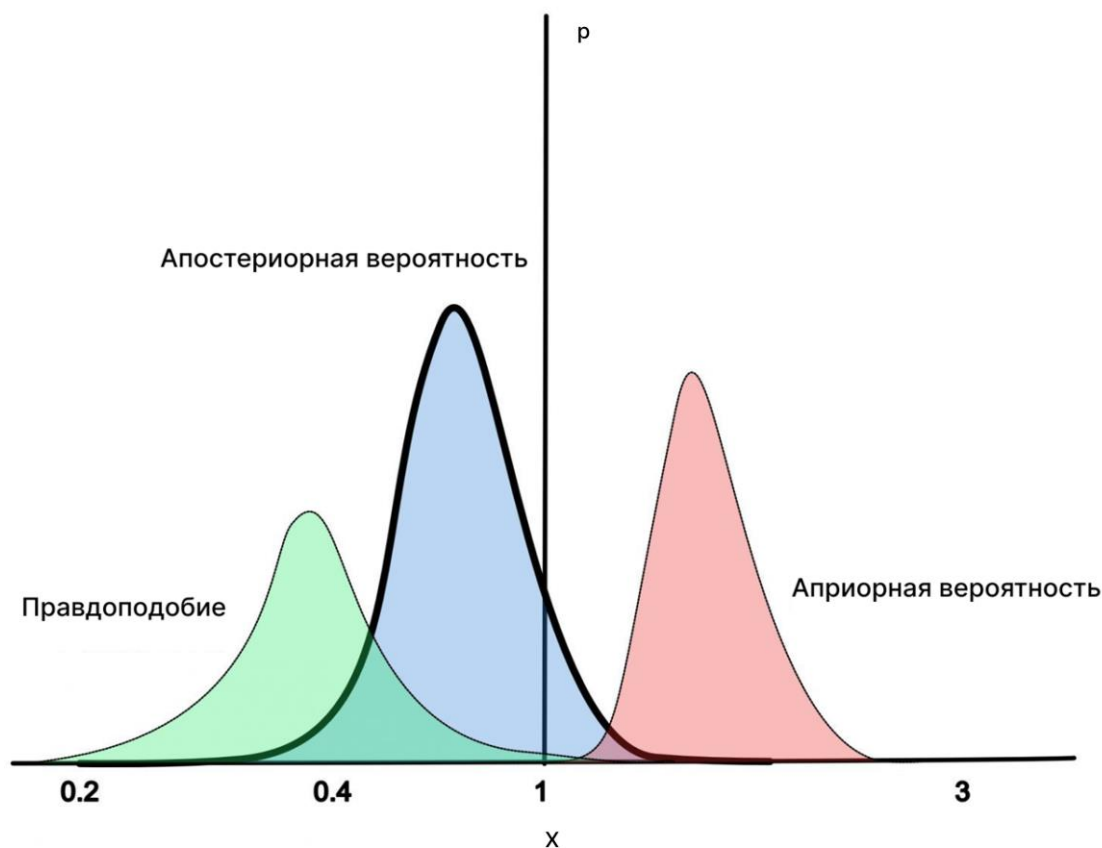


Рисунок 6 – Графический пример работы Байесовской теоремы

Глава 2. Нейронные сети

2.1 Нейронная сеть прямого распространения

Линейные модели для регрессии и классификации базируются на линейных комбинациях фиксированных нелинейных базисных функций $\phi_j(x)$ и имеют форму

$$y(x, w) = f\left(\sum_{j=1}^M \omega_j \phi_j(x)\right) \quad (2.1)$$

где $f(\cdot)$ нелинейная функция активации в случае классификации и идентично в случае регрессии. Наша цель расширить данную модель создав базисную функцию $\phi_j(x)$, которая будет зависеть от параметров и потом позволим этим параметрам изменяться вместе с коэффициентами $\{\omega_j\}$. Конечно, существует много способов построения параметрической нелинейной базисной функции. Нейронные сети используют базисные функции формы, подобной (2.1), поэтому каждая базисная функция является нелинейной функцией для линейных входных данных, где коэффициенты в линейной комбинации являются адаптивными параметрами.

Это приводит нас к базовой нейронной сети, которая может быть описана как серия функциональных преобразований. В начале мы строим M линейных комбинаций входных переменных x_1, \dots, x_D в форме

$$a_j = \sum_{i=1}^D \omega_{ji}^{(1)} x_i + \omega_{j0}^{(1)} \quad (2.2)$$

где $j = 1, \dots, M$ и верхний индекс (1) означает, что данные параметры находятся на «первом» уровне нейронной сети. Мы будем ссылаться на параметр $\omega_{ji}^{(1)}$ как на вес, а на параметр $\omega_{j0}^{(1)}$ как на погрешность. Переменные a_j

известны как «активатор». Потом каждая из них преобразуется с помощью дифференцируемой нелинейной «функции активации» $h(\cdot)$, что дает нам

$$z_j = h(a_j) \quad (2.3)$$

Переменные z_j соответствуют результатам в базисной функции (2.1), которые в контексте нейронных сетей, называются «скрытым узлом». Нелинейная функция $h(\cdot)$ обычно являются сигмоидальными функциями, такими как логическая сигмоидная или «tanh» функция. Учитывая (2.1) эти значения снова линейно комбинируются, чтобы получить выходной слой активации.

$$a_k = \sum_{j=1}^M \omega_{kj}^{(2)} z_j + \omega_{k0}^{(2)} \quad (2.4)$$

где $k = 1, \dots, K$ и K общее число результатов. Это преобразование соответствует второму уровню нейронной сети, и снова $\omega_{k0}^{(2)}$ является параметром погрешности. Наконец, выходной слой активации преобразуется, используя соответствующую функцию активации, чтобы дать набор выводов нейронной сети y_k . Выбор функции активации основывается на природе данных и предполагаемым распределением искомым данных и следует тем же соображениям, что и для линейных моделей. Таким образом, для стандартных регрессионных задач функция активации идентична, значит $y_k = a_k$. Аналогичным образом, для нескольких задач бинарной классификации, каждый выходной слой активации преобразуется, используя логическую сигмоидную функцию так, что

$$y_k = \sigma(a_k) \quad (2.5)$$

где

$$\sigma(a_k) = \frac{1}{1 + \exp(-a)} \quad (2.6)$$

Мы можем объединить различные этапы для того, чтобы получить общую сетевую функцию, которая для сигмоидной функции слоя активации имеет вид

$$y_k(x, w) = \sigma \left(\sum_{j=1}^M \omega_{kj}^{(2)} h \left(\sum_{i=1}^D \omega_{ji}^{(1)} x_i + \omega_{j0}^{(1)} \right) + \omega_{k0}^{(2)} \right) \quad (2.7)$$

где набор всех параметров веса и погрешности сгруппирован в один вектор w . Таким образом нейронная сеть — это нелинейная функция из множества входных $\{x_i\}$ и выходных $\{y_k\}$ значений, регулируемых вектором w , вектором регулируемых параметров.

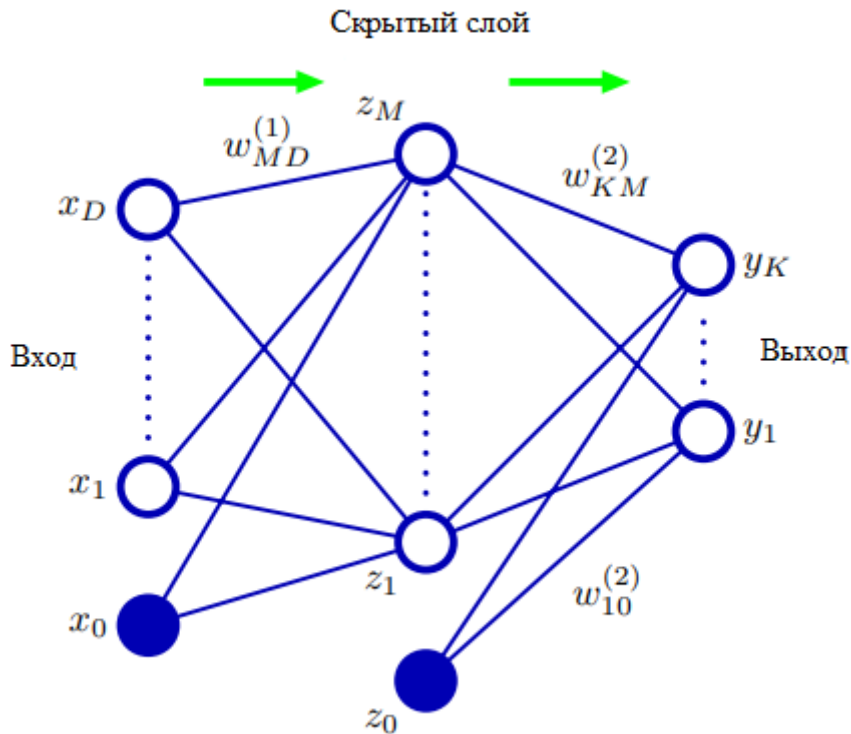


Рисунок 7 – Сетевая диаграмма для двухслойной нейронной сети (2.7)

Процесс расчета (2.7) может быть интерпретирована как прямое пространство данных через нейронную сеть.

Параметры погрешности в (2.2) могут быть включены во множество весовых параметров при помощи определения дополнительной входной переменной x_0 , определим ее значение как $x_0 = 1$, тогда (2.2) примет вид

$$a_j = \sum_{i=0}^D \omega_{ji}^{(1)} x_i \quad (2.8)$$

Мы можем аналогичным образом превратить веса второго слоя в погрешности второго слоя, так что общая сетевая функция примет вид

$$y_k(x, w) = \sigma \left(\sum_{j=0}^M \omega_{kj}^{(2)} h \left(\sum_{i=0}^D \omega_{ji}^{(1)} x_i \right) \right) \quad (2.9)$$

Как мы могли видеть на рисунке 2, модель нейронной сети содержит два этапа обработки данных, каждая из которых напоминает «персептрон» и по этой причине нейронная сеть известна как многослойный персептрон (MLP или multilayer perceptron). Ключевое различие с персептроном — это то, что MLP использует непрерывные сигмоидальные нелинейности в скрытом слое, в то время как персептрон использует кусочно-постоянные нелинейности. Это значит, что MLP дифференцируется по отношению к параметрам нейронной сети и это свойство имеет решающую роль в обучении нейронной сети.

Если функции активации скрытого слоя в нейронной сети приняты как линейные, то для каждой нейронной сети мы всегда можем найти эквивалентную нейронную сеть без скрытого слоя. Это следует из того факта, что композиция последовательных линейных преобразований и есть, само по себе, линейное преобразование. Однако, если число скрытых узлов меньше, чем количество узлов ввода и вывода, то преобразования этой нейронной

сети не являются наиболее общими возможными линейными преобразованиями ввода и вывода, так как данные теряются в уменьшении размерности скрытых узлов.

Архитектура нейронной сети, показанная на рисунке 2, является наиболее часто используемой на практике. Однако ее легко обобщить, рассматривая дополнительные слои обработки, каждая из которых состоит из взвешенной линейной комбинации вида (2.4) с последующим поэлементным преобразованием при помощи нелинейной функции активации. Обратим внимание, что существует некоторая путаница в терминологии, по отношению к подсчету количества слоев в таких сетях. Таким образом сеть на рисунке 2 может быть описана как трехуровневая сеть (которая считает количество узлов в слоях и трактует входные данные как узлы). Мы рекомендуем терминологию, в которой нейронная сеть на рисунке 2 называется двухуровневой нейронной сетью, так как количество слоев адаптивных весов важнее для определения свойств сети.

Другим обобщением сетевой архитектуры является включение соединения «пропускного слоя», каждое из которых связано с соответствующим адаптивным параметром. Для реализации, в двухслойной нейронной сети они будут идти напрямую от входным данным к выходным данным. В принципе, нейронная сеть с сигмоидными спрятанными узлами могут всегда может имитировать соединения пропускного слоя (для ограниченных входные значений) используя достаточно небольшой первый слой весов, что во всем рабочем диапазоне скрытые узлы являются линейными, а затем компенсируется большим значением весов скрытого слоя для выходных данных. На практике, однако, возможно эффективно включать пропускной слой явно.

Кроме того, сеть может быть разрозненной и не все возможные соединения присутствуют.

Поскольку существует прямое соответствие между диаграммой нейронной сети и ее математической функцией, мы можем разработать общие сетевые отображения, за счет рассмотрения более сложных диаграмм

нейронных сетей. Однако они должны быть ограничены архитектурой с прямым распространением, другими словами, для тех, у кого нет замкнутых направленных циклов, нужно убедиться, что выходные данные являются детерминированными функциями входных данных. Это продемонстрировано на простом примере на рисунке 3. Каждый (скрытый или входной) узел в нейронной сети вычисляется как

$$z_k = h \left(\sum_j \omega_{kj} z_j \right) \quad (2.10)$$

где сумма проходит по всем узлам, которые отправляют соединения с узлом k (и параметр погрешности включен в суммирование). Для заданного набора значений, примененных к входным данным нейронной сети, последовательное применение (2.10) позволяет активировать все узлы в нейронной сети, подлежащие расчету, включая выходные слои.

Нейронные сети называют универсальными аппроксиматорами. Для примера, двухслойная нейронная сеть с линейными выходными данными может равномерно аппроксимировать любую непрерывную функцию на компактной входной области с произвольной точностью при условии, что нейронная сеть имеет достаточное количество скрытых узлов. Этот результат справедлив для широкого для широкого спектра функций активации скрытых узлов, за исключением полиномиальных. Хотя такие теоремы обнадеживают, ключевая проблема заключается в том, как найти подходящие значения параметров из множества тренировочных данных.

Способность двухуровневой сети моделировать широкий спектр функций изображена на рисунке 8.

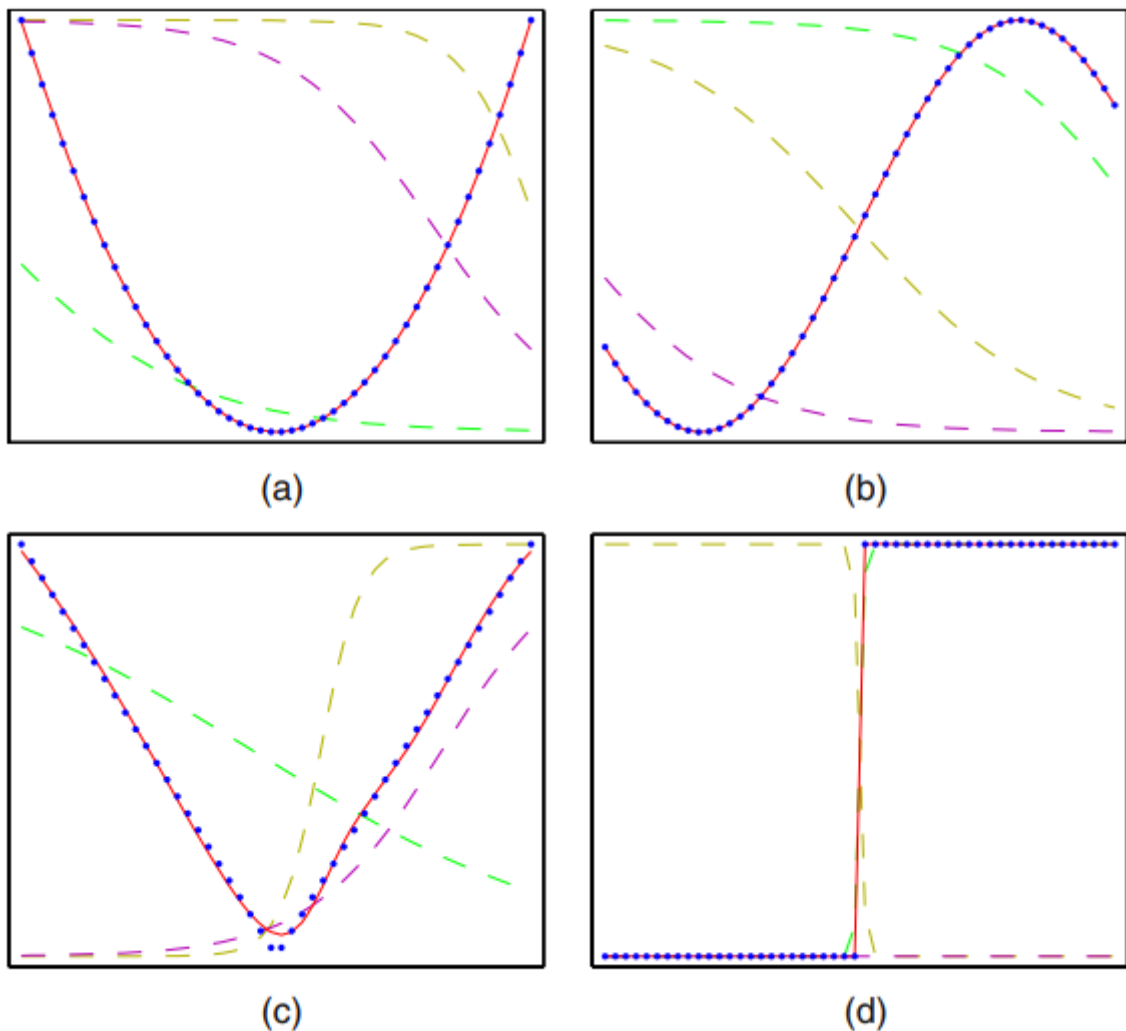


Рисунок 8 – Иллюстрация возможностей многослойного персептрона для аппроксимации четырех различных функций (a) $f(x) = x^2$ (b) $f(x) = \sin(x)$ (c) $f(x) = |x|$ (d) $f(x) = H(x)$, где $H(x)$ – ступенчатая функция Хэвисайда. В каждом примере $N = 50$ точки данных показаны синими точками, были выбраны равномерно в интервале $(-1,1)$ и соответствующие вычисленные значения $f(x)$. Эти точки в дальнейшем используются для обучения двухуровневой сети, имеющей скрытые блоки с функциями активации «tanh» и линейными выходными блоками. Результирующие сетевые функции показаны красными кривыми, а выходы трех скрытых блоков показаны тремя пунктирными кривыми.

2. 2 Обучение сети

До сих пор мы рассматривали нейронные сети как общий класс параметрических нелинейных функций из вектора x (входных параметров) в вектор y (выходных параметров). Простой подход к проблеме определения параметров нейронной сети состоит в том, чтобы минимизировать функцию

ошибки суммы квадратов. Учитывая обучающий набор, состоящий из входных параметров вектора x_n , где $n = 1, \dots, N$ вместе с соответствующим набором целевых векторов t_n , мы минимизируем функцию ошибки.

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2 \quad (3.1)$$

Однако мы можем предоставить гораздо более общий взгляд на обучение нейронной сети, в первую очередь дав вероятностную интерпретацию выходным данным сети. Нам это даст более четкую мотивацию, как для выбора нелинейности выходных узлов, так и для выбора функции ошибки. Мы начнем с обсуждения проблем регрессии, а пока рассмотрим единственную целевую переменную t , которая может принимать любое реальное значение. Мы предполагаем, что t имеет Гауссово распределение с зависимыми от x средним значением, которое задается выходными значениями нейронной сети, так что

$$p(t|x, w) = N(t|y(x, w), \beta^{-1}) \quad (3.2)$$

где β точность (обратная дисперсии) Гауссовского шума. Для условного распределения (3.2) достаточно в качестве функции активации выходного устройства взять тождество, так как такая сеть может аппроксимировать любую непрерывную функцию от x до y . Учитывая набор данных из N независимых, одинаково распределенных наблюдений $X = \{x_1, \dots, x_N\}$, вместе с соответствующими целевыми значениями $t = \{t_1, \dots, t_N\}$, мы можем построить соответствующую функцию правдоподобия

$$p(t|X, w, \beta) = \prod_{n=1}^N p(t_n, x_n, w, \beta) \quad (3.3)$$

Взяв отрицательный логарифм, мы получим

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi) \quad (3.4)$$

с помощью которых мы можем узнать параметры w и β . Начнем с рассмотрения w . Максимизация функции правдоподобия эквивалентная минимизации функции ошибки суммы квадратов, определяемая как

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 \quad (3.5)$$

где отброшены аддитивные и мультипликативные константы. Значение w находится минимизацией $E(w)$, которое будет обозначено w_{ML} , так как оно соответствует решению с максимальным правдоподобием. На практике нелинейность функции нейронной сети $y(x_n, w)$ приводит к «не выпуклости» ошибки $E(w)$, и поэтому на практике локальные максимумы правдоподобия могут быть найдены путем соответствия локальным минимумам функции ошибки.

Найдя w_{ML} , значение β может быть найдено за счет минимизации отрицательной логарифмической вероятности, в результате получим

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{ML}) - t_n\}^2 \quad (3.6)$$

Обратим внимание, что это можно рассчитать после завершения итеративной оптимизации, необходимой для поиска w_{ML} . Если у нас есть несколько целевых переменных и мы предполагаем, что они независимые, обусловленные x и w с общей точностью шума β , то условное распределение целевых значений определяется выражением

$$p(t|x, w) = N(t|y(x, w), \beta^{-1}I) \quad (3.7)$$

Следуя тому же аргументу, что и для одной целевой переменной, мы видим, что максимальные веса правдоподобия определяются путем минимизации функции ошибок суммы квадратов (3.1). Тогда точность шума определяется выражением

$$\frac{1}{\beta_{ML}} = \frac{1}{NK} \sum_{n=1}^N \|y(x_n, w_{ML}) - t_n\|^2 \quad (3.8)$$

где K – количество целевых переменных. От предположения о независимости можно отказаться при помощи немного более сложной задачи оптимизации.

Мы можем рассматривать нейронную сеть, как имеющую на выходе функцию активации, которая тождественна, так что $y_k = a_k$. Соответствующая функция ошибок суммы квадратов имеет свойство

$$\frac{\partial E}{\partial a_k} = y_k - t_k \quad (3.9)$$

которое мы будем использовать при обсуждении обратного распространения ошибок в разделе 4. Обратное распространение ошибки.

Теперь рассмотрим случай бинарной классификации, в которой у нас есть одна целевая переменная t , такая что $t = 1$ обозначает класс C_1 , а $t = 0$ обозначает класс C_2 . Мы рассмотрим нейронную сеть, имеющую один выход, который описан функцией активации в виду логической сигмоиды.

$$y = \sigma(a) \equiv \frac{1}{1 + \exp(-a)} \quad (3.10)$$

тогда $0 \leq y(x, w) \leq 1$. Мы можем интерпретировать $y(x, w)$ как условную вероятность $p(C_1|x)$, где $p(C_2, x)$ задается уравнением $1 - y(x, w)$. Условное распределение целевых переменных при входных данных является распределением Бернулли в виде

$$p(t|x, w) = y(x, w)^t \{1 - y(x, w)\}^{1-t} \quad (3.11)$$

Если мы рассматриваем обучающий набор независимых наблюдений, тогда функция ошибки, которая дается отрицательным логарифмом функции правдоподобия, будет кросс-энтропийной функцией ошибок формы

$$E(w) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (3.12)$$

где y_n обозначает $y(x_n, w)$. Отметим, что нет аналога шумовой точности β , так как предполагается, что целевые значения правильно помечены. Однако модель легко расширяется, чтобы учесть ошибки маркировки. Simard et al. (2003) обнаружил, что использование кросс-энтропийной функции ошибок вместо суммы квадратов для задачи классификации приводит к более быстрому обучению, а так же к лучшему обобщению.

Если у нас есть K отдельных двоичных классификаций для выполнения, то мы можем использовать нейронную сеть, имеющую K выходов, каждая из которых представлена логической сигмоидной функцией активации. С каждой выходом связана метка двоичного класса $t_k \in \{0,1\}$, где $k = 1, \dots, K$. Если мы предположим, что метки классов независимы, учитывая входной вектор, то условное распределение целевых данных будет иметь вид

$$p(t|x, w) = \prod_{k=1}^K y_k(x, w)^{t_k} [1 - y_k(x, w)]^{1-t_k} \quad (3.13)$$

Отрицательный логарифм соответствующей функции правдоподобия дает следующую функцию ошибок

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K \{t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk})\} \quad (3.14)$$

где y_{nk} обозначает $y_k(x_n, w)$. Снова, производная функций ошибок со спектром активации для конкретных выходных узлов принимает вид (3.9), как и в случае регрессии.

Предположим, что мы используем стандартную двухуровневую нейросеть показанную. Мы видим, что весовые параметры в первом слое

сети распределяются между различными выходами, тогда как в линейной модели каждая классификация решается независимо. Первый уровень сети может рассматриваться как выполнение нелинейного извлечения признаков, и совместное использование извлечения между различными выходными данными может сэкономить на вычислениях, а также может привести к улучшению обобщения.

Наконец, мы рассмотрим стандартную задачу много классовой классификации, в которой каждый вход относится к одному из K взаимоисключающих классов. Бинарные целевые переменные $t_k \in \{0,1\}$ имеет схему кодирования 1 из K с указанием класса, а выходы нейронной сети интерпретируются как $y_k(x, w) = p(t_k = 1|x)$, в итоге имею следующую функцию ошибки

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(x_n, w) \quad (3.15)$$

Мы видим, что функция активации узлов выхода, соответствующая канонической ссылке, задается функцией «SoftMax».

$$y_k(x, w) = \frac{\exp(a_k(x, w))}{\sum_j \exp(a_j(x, w))} \quad (3.16)$$

что удовлетворяет $0 \leq y_k \leq 1$ и $\sum_k y_k = 1$. Отметим, что $y_k(x, w)$ не изменяются, если ко всем $a_k(x, w)$ добавляется константа, в результате чего функция ошибок остается постоянной для некоторых направлений в пространстве весов. Это вырождение снимается, если к функции ошибок добавить соответствующий член регуляризации.

Еще раз, производная функции ошибок относительно активации для конкретного выходного узла принимает знакомую форму (3.9).

Таким образом, существует естественный выбор как функции активации узла вывода, так и функции ошибки согласования, в зависимости от типа решаемой проблемы. Для регрессии мы используем линейные выходы и

ошибку суммы квадратов, для (нескольких независимых) двоичных классификаций мы используем логистические сигмоидные выходы (функции активации для этих выходов) и функцию кросс-энтропийной ошибки, а для многоклассовой классификации мы для выходов используем «SoftMax», с соответствующей мультиклассовой кросс-энтропийной функцией ошибки. Для задач классификации, включающих два класса, мы можем использовать один логистический сигмоидный выход, или, как альтернатива, мы можем использовать нейронную сеть с двумя выходами, каждая из которых имеет функцию активации «SoftMax».

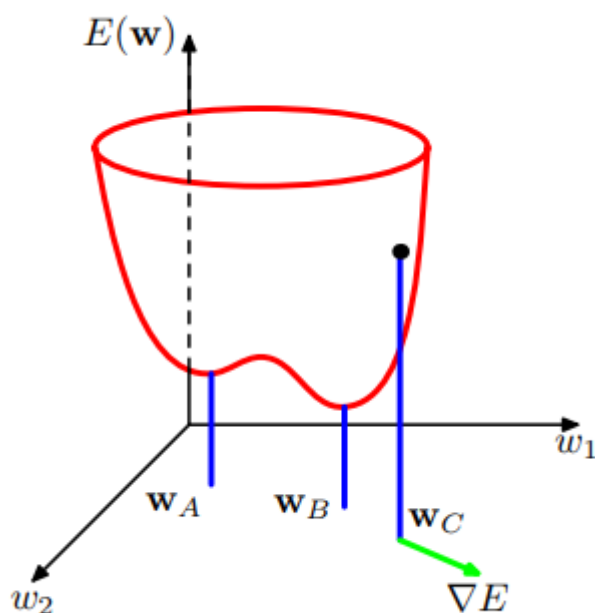


Рисунок 9. – Геометрический вид функции ошибки $E(w)$ как поверхности, расположенной над пространством весов. Точка w_A – это локальный минимум, а w_B – глобальный минимум. В любой точке w_C локальный градиент ошибки задается вектором ∇E .

2.3 Обратное распространение ошибки

Наша цель в этом разделе – найти эффективный метод расчета градиента функции ошибок $E(w)$ для нейронной сети с прямой связью.

Мы покажем, что это может быть достигнуто с помощью схемы передачи локальных сообщений, в которой информация пересылается поочередно вперед и назад по сети и известна как «Распространение ошибок» или проще «Обратное распространение».

Следует отметить, что термин обратное распространение используется в литературе для обозначения множества разных вещей. Например, многослойную архитектуру персептрона иногда называют сетью обратного распространения. Также термин обратного распространения используется для описания обучения многослойного персептрона с использованием градиентного спуска, применяемого к функции ошибок суммы квадратов. Для уточнения терминологии полезно более внимательно рассмотреть характер тренировочного процесса. Большинство обучающих алгоритмов включают итеративную процедуру минимизации функции ошибок, с корректировкой весов. Каждый шаг мы можем разбить на два этапа. На первом этапе необходимо вычислить производные функции ошибок по весам. Как мы увидим позже, важный вклад обратного распространения заключается в предоставлении эффективного метода для вычисления таких производных. Так как именно на этом этапе ошибки распространяются по нейронной сети в обратном направлении, мы будем использовать термин обратное распространение специально для расчета производных. На втором этапе производные затем используются для вычисления корректировок, которые необходимо внести в веса. Самый простой из таких приемов – градиентный спуск. Важно понимать, что эти две стадии различны. Таким образом, первый этап, а именно распространение ошибок назад по сети для расчета производных, может быть применен ко многим другим типам сетей, а не только к многослойному персептрону. Его также можно применить к функциям ошибок, отличающихся от суммы квадратов, и к вычислению других производных, таких как матрицы Якобиана и Гессе. Аналогично, второй этап регулировки веса с помощью вычисленных

производных могут быть решены с использованием различных схем оптимизации, многие из которых значительно более мощные, чем простой градиентный спуск.

2.4 Расчет производных функции ошибок

Теперь мы выведем алгоритм обратного распространения ошибок для общей сети, имеющей произвольную топологию прямой связи, произвольные дифференцируемые функции активации и широкий класс функции ошибок. Затем полученные формулы будут проиллюстрированы, используя многослойную нейронную сеть, имеющей один слой сигмоидальных скрытых узлов вместе с функцией ошибки суммы квадратов.

Многие функции ошибок, представляющий практический интерес, например те, которые определяются максимумом функции правдоподобия для набора независимых и одинаково распределенных данных, содержат сумму членов, по одному для каждой точки данных в обучающем наборе, поэтому

$$E(w) = \sum_{n=1}^N E_n(w) \quad (4.1)$$

Здесь мы рассмотрим задачу вычисления градиента ∇E для одного такого члена в функциях ошибок. Это может быть использовано непосредственно для последовательной оптимизации, или результаты могут накапливаться по обучающей выборке в случае пакетных методов.

Рассмотрим сначала простую линейную модель, в которой выводы y_k являются линейными комбинациями x_i , так что

$$y_k = \sum_i \omega_{ki} x_i \quad (4.2)$$

вместе с функцией ошибок, которая для конкретного шаблона n принимает форму

$$E_n = \frac{1}{2} \sum_k (y_{nk} - t_{tk})^2 \quad (4.3)$$

где $y_{nk} = y_k(x_n, w)$. Градиент этой функции ошибок относительно веса ω_{ji} определяется выражением

$$\frac{\partial E_n}{\partial \omega_{ji}} = (y_{nj} - t_{nj})x_{ni} \quad (4.4)$$

который может быть интерпретирован как «локальное» вычисление включающие результат «сигнала ошибки» $y_{nj} - t_{nj}$ связанного с выходным концом линии ω_{ji} и переменной x_{ni} связанного с входным концом линии.

В общей сети с прямой связью каждый узел вычисляет взвешенную сумму своих входных данных в форме

$$a_j = \sum_i \omega_{ji} z_i \quad (4.5)$$

где z_i – это активация узла или входа, который отправляет соединение узлом j и ω_{ji} вес, который связан этим подключением. В разделе 2 мы видели, что смещения можно включить в эту сумму, введя дополнительную единицу или вход, с фиксированной функцией активацией +1. Поэтому нам не нужно явно разбираться с предубеждениями явно. Сумма в (4.4) преобразуется нелинейной функцией активации $h(\cdot)$, чтобы дать активацию z_j узла j в виде

$$z_j = h(a_j) \quad (4.6)$$

Обратим внимание, что одна или несколько переменных z_i в сумме (4.5) могут быть входными, и аналогично, узлы j могут быть выходными. Для каждого шаблона в обучающем наборе мы будем предполагать, что мы предоставили соответствующий входной вектор в сеть и вычислили функции активации всех скрытых и входных узлов в сети путем последовательного применения (4.5) и (4.6). Этот процесс часто называют прямым распространением, так как его можно рассматривать как прямой поток данных через сеть.

Теперь рассмотрим расчет производной от E_n по отношению к весу ω_{ji} . Выходы различных узлов будут зависеть от конкретного входного шаб-

лона n . Однако, чтобы не загромождать обозначения, мы будем опускать индекс n в сетевых переменных. Прежде всего отметим, что E_n зависит только от веса ω_{ji} только через суммированный вход a_j в узел j . Следовательно, мы можем применять цепное правило для частных производных, чтобы получить

$$\frac{\partial E_n}{\partial \omega_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial \omega_{ji}} \quad (4.7)$$

теперь введем полезное обозначение

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} \quad (4.8)$$

где δ ошибка по причинам, которую мы скоро увидим. Используя (4.6), мы можем написать

$$\frac{\partial a_j}{\partial \omega_{ji}} = z_i \quad (4.9)$$

подставляя (4.8) и (4.9) в (4.7), мы получим

$$\frac{\partial E_n}{\partial \omega_{ji}} = \delta_j z_i \quad (4.10)$$

Уравнение (4.10) говорит нам, что требуемая производная получается простым умножением значения δ для веса конечного узла выхода и z для веса конечного узла входа (где $z = 1$ в случае смещения). Обратим внимание, что уравнение примет ту же форму, что и для простой линейной модели, рассмотренной во втором разделе. Таким образом, чтобы рассчитать производные, нам нужно только вычислить значение δ_j для каждого скрытого и выходного узла в сети, а затем применить в (4.9)

Как мы могли видеть для выходных узлов у нас есть

$$\delta_k = y_k - t_k \quad (4.11)$$

при условии, что мы используем каноническую ссылку в качестве функции активации в выходном узле. Чтобы рассчитать δ для скрытых узлов, мы снова используем цепное правило для частных производных.

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \quad (4.12)$$

где сумма пробегает все узлы k , к которым узел j отправляет соединения. Расположение узлов и весов проиллюстрировано на рисунке 6. Обратим внимание, что узлы, помеченные k , могут включать другие скрытые узлы и/или узлы входа. В записи (4.11) мы используем тот факт, что вариации a_j вызывают вариации функции ошибок только через вариации переменных a_k . Если мы теперь подставим определение δ из (4.7) в (4.11) и воспользуемся формулами (4.4) и (4.5), мы получим следующую формулу обратного распространения ошибок.

$$\delta_j = h'(a_j) \sum_k \omega_{kj} \delta_k \quad (4.13)$$

которая говорит нам, что значение δ для конкретного скрытого узла может быть получено путем распространения δ в обратном направлении от узлов выше в сети, как и проиллюстрировано на рисунке 6. Отметим, что суммирование в (4.13) ведется по первому индексу ω_{kj} (соответствуя обратному распространению данных по сети), тогда как в уравнении прямого распространения (2.10) оно проходит по второму индексу. Поскольку мы уже знаем значения δ для выходных узлов отсюда следует, что, рекурсивно применяя (4.13), мы можем рассчитать δ для всех скрытых узлов в сети прямого распространения, независимо от ее топологии.

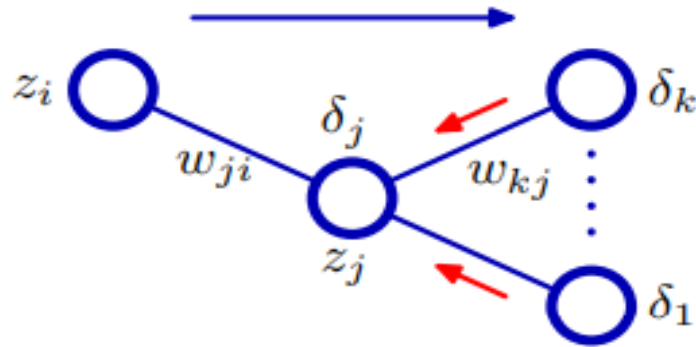


Рисунок 9 – Иллюстрация расчета δ_j для скрытых узлов j с помощью обратного распространения δ от единиц k , к которым блок j отправляет соединения. Синяя стрела обозначает направление данных при прямом распространении, а красные стрелки указывают на обратное распространение данных ошибки.

Таким образом процедуру обратного распространения можно резюмировать следующим образом

1. Примените входной вектор x_n к сети и распространите его вперед по сети, используя (4.4) и (4.5), чтобы найти активации всех скрытых и входных узлов.

2. Рассчитайте δ_k для всех выходных узлов, используя (4.11)

3. Выполните обратное распространение δ_k , используя (4.13), чтобы получить δ_j для каждого скрытого узла в сети.

4. Используйте (4.10) для расчета требуемых производных.

Для пакетных методов производная полной ошибки E может быть затем получена путем повторения вышеуказанных шагов для каждого шаблона в обучающем наборе и последующего суммирования через все шаблоны

$$\frac{\partial E}{\partial \omega_{ji}} = \sum_n \frac{\partial E_n}{\partial \omega_{ji}} \quad (4.13)$$

В приведенном выше выводе мы неявно предполагали, что каждый скрытый и входной узел в сети имеет одну и ту же функцию активации $h(\cdot)$. Однако вывод легко обобщается, что позволит различным узлам иметь индивидуальные функции активации, просто отслеживая, какая форма $h(\cdot)$ соответствует какому узлу.

Список использованных источников

1. Мартин, Освальдо. Байесовский анализ на Python : введение в статистическое моделирование и вероятностное программирование с использованием PyMC3 и ArviZ / Освальдо Мартин ; перевод с английского А. В. Снастина. – Москва : ДМК Пресс, 2020. - 339 с.

2. Дауни, Аллен Б. Байесовские модели : байесовская статистика на языке программирования Python / Аллен Б. Дауни ; перевод с английского В. А. Яроцкого. - Москва : ДМК Пресс, 2018. - 181 с.

3. Нильсен, Эйлин. Практический анализ временных рядов : прогнозирование со статистикой и машинное обучение: перевод с английского / Нильсен, Эйлин. – Москва : Диалектика ; Санкт-Петербург Диалектика, 2021. - 538 с.

4. Гнеденко, Борис Владимирович. Курс теории вероятностей : [Учеб. Для мат. спец. ун-тов] / Б. В. Гнеденко. - 6-е изд., перераб. и доп. - М. : Наука, 1988. - 446 с.

5. Гмурман, Владимир Ефимович. Теория вероятностей и математическая статистика : учебное пособие для студентов вузов / В. Е. Гмурман. - 12-е изд. / перераб. - Москва : Юрайт : Высш. образование, 2009. - 478 с.

6. Паттерсон, Джош. Глубокое обучение с точки зрения практика / Джош Паттерсон, Адам Гибсон ; пер. с англ. А. А. Слинкина. - Москва : ДМК Пресс, 2018. - 415 с.

7. Christopher M. Bishop. Pattern Recognition and Machine Learning / Christopher Michael Bishop - New York : Springer New York, cop. 2006. - 738 p.
8. Ekaba Bisong. Building Machine Learning and Deep Learning Models on Google Cloud Platform / Ekaba Bisong - Berkeley, CA : Apress, cop. 2019. – 709 с.
9. Гудфеллоу Я. Глубокое обучение / Я. Гудфеллоу, И. Бенджио, А. Курвилль ; [пер. с англ. А. А. Слинкина]. - 2-е цв. изд., испр. - Москва : ДМК Пресс, 2018. - 651 с.
10. Микелуччи У. Прикладное глубокое обучение : подход к пониманию глубоких ннейронных сетей на основе метода кейсов / У. Микелуччи ; перевод с английского Андрея Логунова. - Санкт-Петербург : БХВ-Петербург, 2020. - 368 с.
11. Bayesian Data Analysis / Andrew Gelman, John B. Carlin, Hal S. Stern and etc.- third edition - Oxford, United States : CRC Press, cop. 2013. – 675 с.
12. Andrew Gelman. The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo / Andrew Gelman // Journal of Machine Learning Research. – 2014. – №15. – С. 1593-1623.
13. Alp Kucukelbir. Automatic Differentiation Variational Inference / Alp Kucukelbir // Journal of Machine Learning Research. – 2017. - №18. – С. 1-45