

ОБРАЗОВАТЕЛЬНАЯ МАНГА



МАШИННОЕ ОБУЧЕНИЕ

СОТРУДНИКУ ГОРОДСКОЙ АДМИНИСТРАЦИИ КИЁХАРА КАДЗУМА ПОРУЧЕНО ЗАДАНИЕ, КОТОРОЕ БЕЗ МАШИННОГО ОБУЧЕНИЯ НЕ ВЫПОЛНИТЬ. ПОД РУКОВОДСТВОМ СВОЕЙ ДАВНЕЙ ЗНАКОМОЙ МИЯНО САЯКА ОН ОСВАИВАЕТ ПРЕМУДРОСТИ РАБОТЫ С ИСКУССТВЕННЫМ ИНТЕЛЛЕКТОМ – ОТ САМЫХ АЗОВ ДО ГЛУБОКОГО ОБУЧЕНИЯ.

ВМЕСТЕ С ГЕРОЯМИ МАНГИ ЧИТАТЕЛИ УЗНАЮТ О ТОМ, ЧТО ТАКОЕ РЕГРЕССИЯ И КАК ПРОВОДИТЬ КЛАССИФИКАЦИЮ, ОЗНАКОМЯТСЯ С ПРИНЦИПАМИ ОЦЕНКИ ТЕСТОВЫХ ДАННЫХ И ОСОБЕННОСТЯМИ РАБОТЫ НЕЙРОННЫХ СЕТЕЙ. В ЗАКЛЮЧИТЕЛЬНОЙ ЧАСТИ ИЗЛАГАЮТСЯ МЕТОДЫ ОБУЧЕНИЯ БЕЗ УЧИТЕЛЯ.

МАНГА ПРЕДНАЗНАЧЕНА ДЛЯ ТЕХ, КТО НАЧИНАЕТ ЗНАКОМСТВО С МАШИНЫМ ОБУЧЕНИЕМ И ОСВОИЛ МАТЕМАТИКУ НА УРОВНЕ ПЕРВЫХ КУРСОВ УНИВЕРСИТЕТА.

Интернет-магазин:
www.dmkpress.com

Оптовая продажа:
КТК «Галактика»
books@aliens-kniga.ru

DMK
издательство
www.dmk.ru

ISBN 978-5-97060-830-2



9 785970 608302 >

ЗАНИМАТЕЛЬНАЯ

МАНГА

МАШИННОЕ ОБУЧЕНИЕ

Араки Масахиро
Ватари Макана

ЗАНИМАТЕЛЬНАЯ МАНГА
МАШИННОЕ ОБУЧЕНИЕ



Араки Масахиро
Ватари Макана
Office SAWA, Ltd.



Ohmsha

DMK
издательство

Занимательное машинное обучение

Манга

マンガでわかる

機械学習

荒木 雅弘／著
渡 まかな／作画
ウェルテ／制作



OHM
Ohmsha

ОБРАЗОВАТЕЛЬНАЯ МАНГА

ЗАНИМАТЕЛЬНОЕ

МАШИННОЕ ОБУЧЕНИЕ

Араки Масахиро
Художник Ватари Макана

Перевод
А. С. Слащевой



ДМК
издательство

Москва
ДМК Пресс, 2020

УДК 004.4
ББК 32.972
А79

Араки М., Ватари М.

А79 Занимательная манга. Машинное обучение: манга / Араки Масахиро (автор), Ватари Макана (худ.); пер. с яп. А. С. Слащевой. — М.: ДМК Пресс, 2020. — 214 с. : ил. — (Серия «Образовательная манга»). — Доп. тит. л. яп.

ISBN 978-5-97060-830-2

Сотруднику городской администрации Киёхара Кадзума поручено задание, которое без машинного обучения не выполнить. Под руководством своей давней знакомой Мияно Саяка он осваивает премудрости работы с искусственным интеллектом – от самых азов до глубокого обучения.

Вместе с героями манги читатели узнают о том, что такое регрессия и как проводить классификацию, ознакомятся с принципами оценки тестовых данных и особенностями работы нейронных сетей. В заключительной части излагаются методы обучения без учителя.

Издание предназначено для тех, кто начинает знакомство с машинным обучением и освоил математику на уровне первых курсов университета.

УДК 004.4
ББК 32.972

Manga de Wakaru Kikai Gakushu (Manga Guide: Machine Learning)

By Araki Masaxiro (Author), Illustration by Vatari Makana

Published by Ohmsha, Ltd.

Russian language edition copyright © 2020 by DMK Press

Все права защищены. Никакая часть этого издания не может быть воспроизведена в любой форме или любыми средствами, электронными или механическими, включая фотографирование, ксерокопирование или иные средства копирования или сохранения информации, без письменного разрешения издательства.

ISBN 978-4-274-22244-3 (яп.)
ISBN 978-5-97060-830-2 (рус.)

Copyright © 2018 by and Office sawa, Ltd.
© Издание, перевод, ДМК Пресс, 2019

ПРЕДИСЛОВИЕ

В этой книге я представил несколько репрезентативных методов машинного обучения и попытался по возможности просто изложить их суть. Ее предполагаемая аудитория – те, кто только начинает знакомство с машинным обучением и уже владеет математикой на уровне первых курсов университета. Но если вы не дружите с математикой, то можете ознакомиться с разъяснениями в конце каждой главы и примерно понять, какие задачи решаются с помощью этих методов.

Особенность данной книги в том, что в начале каждой главы ставится задача, а затем постепенно объясняются методы машинного обучения, необходимые для ее решения. В таблице ниже перечислены задачи и методы, которые будут рассматриваться в каждой главе.

Глава	Задача	Метод
1	Прогноз количества участников мероприятия	Линейная регрессия
2	Определение вероятности заболевания диабетом	Логистическая регрессия, решающее дерево
3	Оценка результатов обучения	Метод проверки на резервированных данных, перекрестная проверка
4	Сортировка винограда	Сверточная нейронная сеть
5	Определение вероятности заболевания диабетом (повтор)	Ансамблевые методы
6	Рекомендация события	Кластерный анализ, матричное разложение

В каждой главе будет предложено лишь введение в тот или иной метод. Если вы хотите применить его на практике для решения какой-либо задачи, я советую обратиться к учебникам, которые указаны в списке рекомендованной литературы в конце книги.

Наконец, я благодарю всех сотрудников издательства Ohmsha за возможность написать эту книгу. Я также благодарен г-же Ватари Макана и всем сотрудникам Уильтэ, которые превратили мою рукопись в веселую мангу.

Июль 2018 года,
Араки Масахиро

СОДЕРЖАНИЕ

ПРЕДИСЛОВИЕ	V
-------------------	---

Пролог

ПОГОВОРИМ О МАШИННОМ ОБУЧЕНИИ	1
-------------------------------------	---

В кабинете у Саяка (1). Саяка и старшекласница Ай	14
---	----

Глава 1

ЧТО ТАКОЕ РЕГРЕССИЯ	15
---------------------------	----

1.1. Сложности с прогнозом	16
----------------------------------	----

1.2. Определяем зависимые и независимые переменные	17
--	----

1.3. Находим функцию линейной регрессии	20
---	----

1.4. Регуляризация результата	22
-------------------------------------	----

В кабинете у Саяка (2). Математическое повторение (1)	34
---	----

Глава 2

КАК ДЕЛАТЬ КЛАССИФИКАЦИЮ?	39
---------------------------------	----

2.1. Приводим данные в порядок	46
--------------------------------------	----

2.2. Определяем класс данных	47
------------------------------------	----

2.3. Логистическая регрессия	49
------------------------------------	----

2.4. Классификация по решающему дереву	55
--	----

В кабинете у Саяка (3). Математическое повторение (2)	74
---	----

Глава 3

ОЦЕНКА РЕЗУЛЬТАТОВ 77

3.1. Без проверки тестовых данных никак нельзя.....	82
3.2. Обучающая, тестовая и проверочная выборки	83
3.3. Метод перекрестной проверки (кросс-валидации)	85
3.4. Доля правильно предсказанных объектов, точность, полнота и F-мера.....	87
В кабинете у Саяка (4). Математическое повторение (3).....	95

Глава 4

ГЛУБОКОЕ ОБУЧЕНИЕ..... 97

4.1. Нейронная сеть	103
4.2. Обучение методом обратного распространения ошибок	107
4.3. Вызовы глубокого обучения.....	111
4.3.1. Проблема глубокой нейронной сети	112
4.3.2. Хитрости многоступенчатого обучения 1. Метод предварительного обучения	113
4.3.3. Хитрости многоступенчатого обучения 2. Функция активации.....	115
4.3.4. Хитрости многоступенчатого обучения 3. Как избежать переобучения.....	117
4.3.5. Нейронные сети со специализированной структурой.....	118
В кабинете у Саяка (5). Математическое повторение (4)	134

Глава 5

АНСАМБЛЕВЫЕ МЕТОДЫ..... 139

5.1. Бэггинг	146
5.2. Случайный лес.....	149
5.3. Бустинг	152
В кабинете у Саяка (6). Математическое повторение (5)	160

Глава 6	
ОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ	165
6.1. Кластеризация.....	172
6.1.1. Иерархическая кластеризация	173
6.1.2. Разделяющая кластеризация.....	175
6.2. Разложение матрицы.....	179
В кабинете у Саяка (7). Математическое повторение (6)	191
 ЭПИЛОГ	 197
 ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ	 205

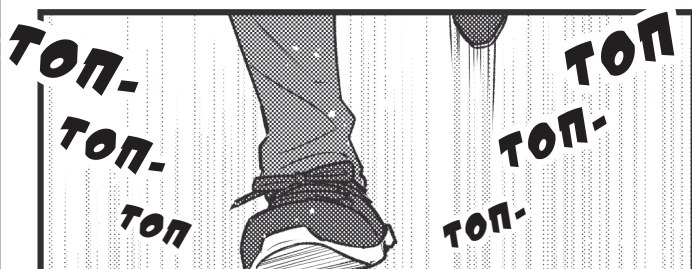
ПРОЛОГ

ПОГОВОРИМ О МАШИННОМ ОБУЧЕНИИ

ЗАЧЕМ НУЖНО
МАШИННОЕ ОБУЧЕНИЕ?



В одном университете

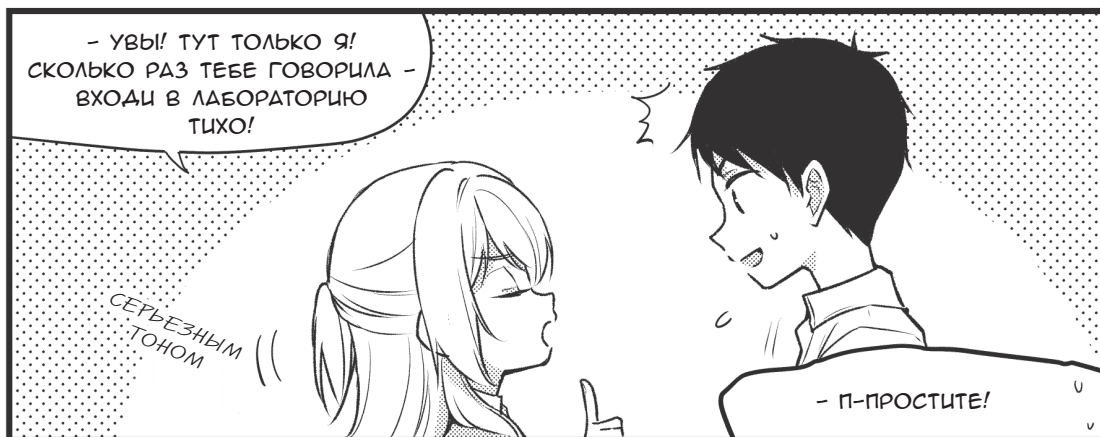




Киёхара Кадзума

Вот уже год как работает в городской администрации.

Изучал машинное обучение в университете на факультете компьютерных технологий, но так ничему толком не выучился.



Сэмпай – название старших по курсу/классу учеников или студентов. Антоним – слово «кохай», которое обозначает младшего по курсу. – Прим. перев.

КИЁХАРА-КУН СОВСЕМ
НЕ ПОМЕНЯЛСЯ,
ХОТЬ ТЕПЕРЬ
И НАШЕЛ РАБОТУ...



Мияно Саяка

Сэмпай Киёхары.
Учится на втором курсе
магистратуры.



ВСЕ ТАКОЙ ЖЕ
ШУМНЫЙ, БЕСТОЛКО-
ВЫЙ, НЕ СЛЫШИШЬ,
ЧТО ТЕБЕ ГОВОРЯТ,
И ПОТОМ...



САЯКА-СЭМПАЙ,
ВЫ КАК ОБЫЧНО...

ЧТО КАК ОБЫЧНО,
КАК ОБЫЧНО
ПРИДИРАЮСЬ...

НЕТ, НЕТ!

КАК ОБЫЧНО,
МИЛАЯ...



Угу
Ну мне
так кажется...

НЕТ, ЭТО НЕ ТАК.
А КОГДА ВЕРНЕТСЯ
ПРОФЕССОР НАМИГОЭ?



ВОТ
ПОДАРОК

ОН ВМЕСТЕ С ДРУГИМИ
ЛАБОРАНТАМИ В КОМАНДИРОВКЕ
ЗА ГРАНИЦЕЙ, РАНЬШЕ, ЧЕМ
ЧЕРЕЗ ДВА МЕСЯЦА,
НЕ ВЕРНЕТСЯ.

СПАСИБО!



ЧТО?!

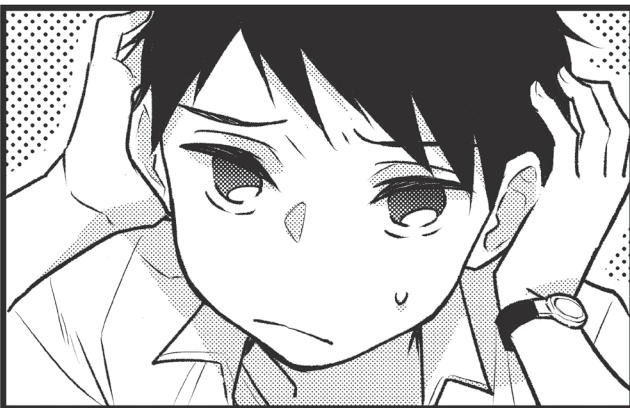


БЛИН, МНЕ ЭТО
НИКАК НЕ ПОДХОДИТ!

ЧТО-ТО
СЛУЧИЛОСЬ?

ТРЯСЕТСЯ

В
УЖАСЕ



Я ХОТЕЛ ПОПРОСИТЬ
ПРОФЕССОРА НАМИГОЗ
РАССКАЗАТЬ МНЕ ПРО
МАШИННОЕ ОБУЧЕНИЕ...

МАШИННОЕ
ОБУЧЕНИЕ?



КИЁХАРА, ТЫ ЖЕ РАБОТАЕШЬ У СЕБЯ
В МЕСТНОЙ АДМИНИСТРАЦИИ.
ТЫ ЖЕ ГОВОРИЛ: "ХОЧУ СПОКОЙНО
РАБОТАТЬ ДОМА И ЖИТЬ В СВОЕ
УДОВОЛЬСТВИЕ". ЗАЧЕМ ТЕБЕ
МАШИННОЕ ОБУЧЕНИЕ?

Киёхара-кун
в университете

О-ХО-ХО...

Что? Работу искать?
Хочу не слишком
напряженно работать
в хорошем месте



У ТЕБЯ ЧТО-ТО
СТЯСЛОСЬ? МОЖЕТ,
МНЕ РАССКАЖЕШЬ?

У НАС ЕСТЬ
КОНСУЛЬТАНТ, КОТОРЫЙ ПОРУЧИЛ
МНЕ ПОРАБОТАТЬ С ЕГО ЦИ-ПРОГРАМ-
МОЙ, ПРОГНОЗИРУЮЩЕЙ КОЛИЧЕСТВО
ГОСТЕЙ НА ПУБЛИЧНЫХ МЕРОПРИЯТИЯХ
НАШЕЙ АДМИНИСТРАЦИИ.

ЭТОТ КОНСУЛЬТАНТ -

К МОМЕНТУ НАСТУПЛЕНИЯ
СИНГУЛЯРНОСТИ ПОЯВИТСЯ
ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ, МОЗГ,
СРАВНИМЫЙ С ЧЕЛОВЕЧЕСКИМ
ИЛИ ДАЖЕ ПРЕВОСХОДЯЩИЙ ЕГО,
КОТОРЫЙ СДЕЛАЕТ ВСЮ ЧЕЛОВЕ-
ЧЕСКУЮ РАБОТУ НЕНУЖНОЙ!

ВООДУШЕВЛЕННО

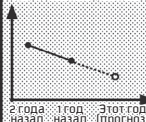
...ПРИМЕРНО ТАКОЙ ТИП...
Я ПОДСЧИТАЛ КОЛИЧЕСТВО ГОСТЕЙ
ЧЕРЕЗ ЦИ-ПРОГРАММУ, КОТОРУЮ
ОН РАЗРАБОТАЛ...

СУДЯ ПО ПРОГНОЗАМ ЦИ,
КОЛИЧЕСТВО ГОСТЕЙ
УМЕНЬШИТСЯ!

СПАСИБО!

УФ!

Количество гостей



ЦИФРЫ БЫЛИ ПОДОЗРИ-
ТЕЛЬНЫМИ, И КОГДА
Я ИХ ПРОВЕРИЛ...

ОЙ!

...ТО УВИДЕЛ,
ЧТО КОЛИЧЕСТВО ГОСТЕЙ
ДВА ПОСЛЕДНИХ ГОДА
СНИЖАЕТСЯ ПО ПРЯМОЙ!

Я ДОЛОЖИЛ ОБ ЭТОМ
ОТВЕТСТВЕННОМУ
ЗА РЕКЛАМУ, НО ТОТ
НЕ ОБРАТИЛ
ВНИМАНИЯ...

ПОЭТОМУ Я ПОДУМАЛ, ЧТО
ЕСЛИ БЫ СМОГ, ИСПОЛЬЗУЯ
ДАННЫЕ ЗА ДЕСЯТЬ ЛЕТ,
СДЕЛАТЬ ПРОГНОЗ
ПРИ ПОМОЩИ МАШИННОГО

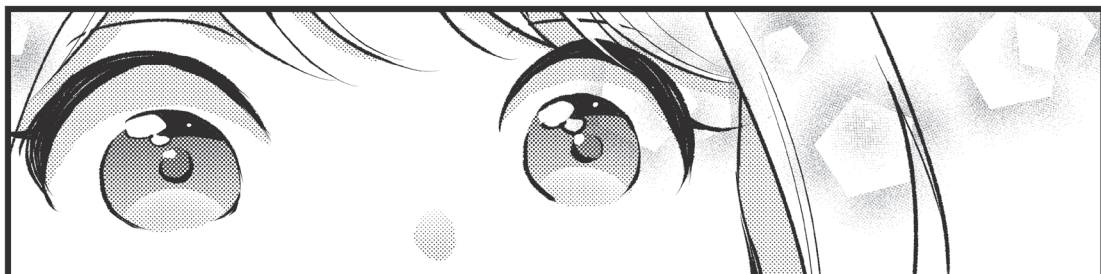
ОБУЧЕНИЯ, ТО У МЕНЯ
ПОЛУЧИЛОСЬ БЫ ЕГО
УБЕДИТЬ...

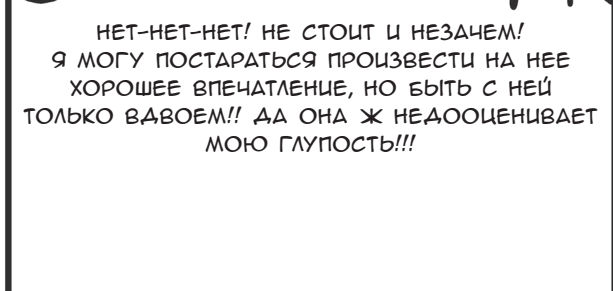
МОЖЕТ,
НА ЭТО КОЛИЧЕСТВО
ОСАДКОВ В СЕЗОН
ДОЖДЕЙ ВЛИЯЕТ?..

НО Я
НИЧЕГО УМЕЮ...

ПОЭТОМУ Я ПРИШЕЛ
К ПРОФЕССОРУ
НАМИГОЗ, ЧТОБЫ ОН
РАССКАЗАЛ МНЕ
О МАШИННОМ
ОБУЧЕНИИ...

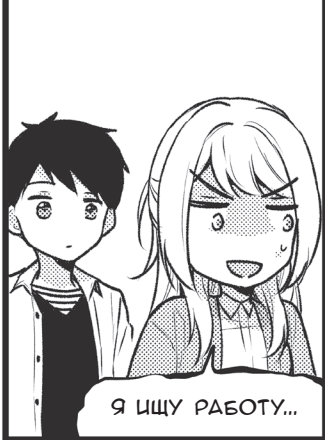








КСТАТИ, САЯКА-СЭМПАЙ,
А ПОЧЕМУ ВЫ
НЕ В КОМАНДИРОВКЕ
С УЧИТЕЛЕМ НАМИГОЭ?



Я ИЩУ РАБОТУ...



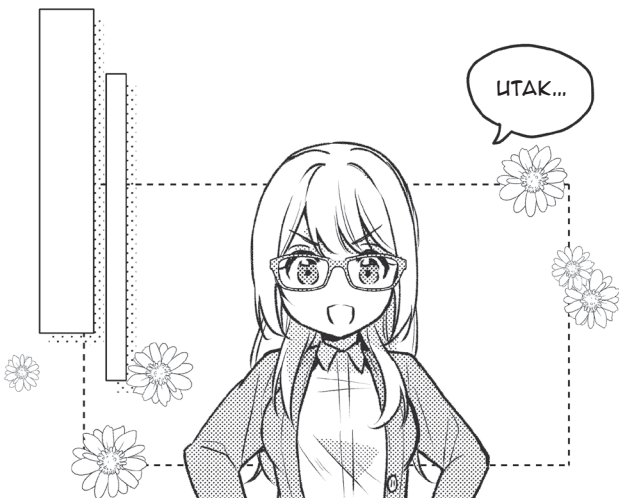
А РАЗВЕ ВЫ
НЕ ДОЛЖНЫ ХОДИТЬ
НА СОБЕСЕДОВАНИЯ?

ПОМОЛЧИ-КА!



Я ПРОСТО НЕ МОГУ ПРЕПОДАВАТЬ
ТЕМ, КТО ГОВОРIT ТАКИЕ ВЕЩИ!

ПРОСТИТЕ! БУДЬТЕ
СНИСХОДИТЕЛЬНЫ!



ИТАК...



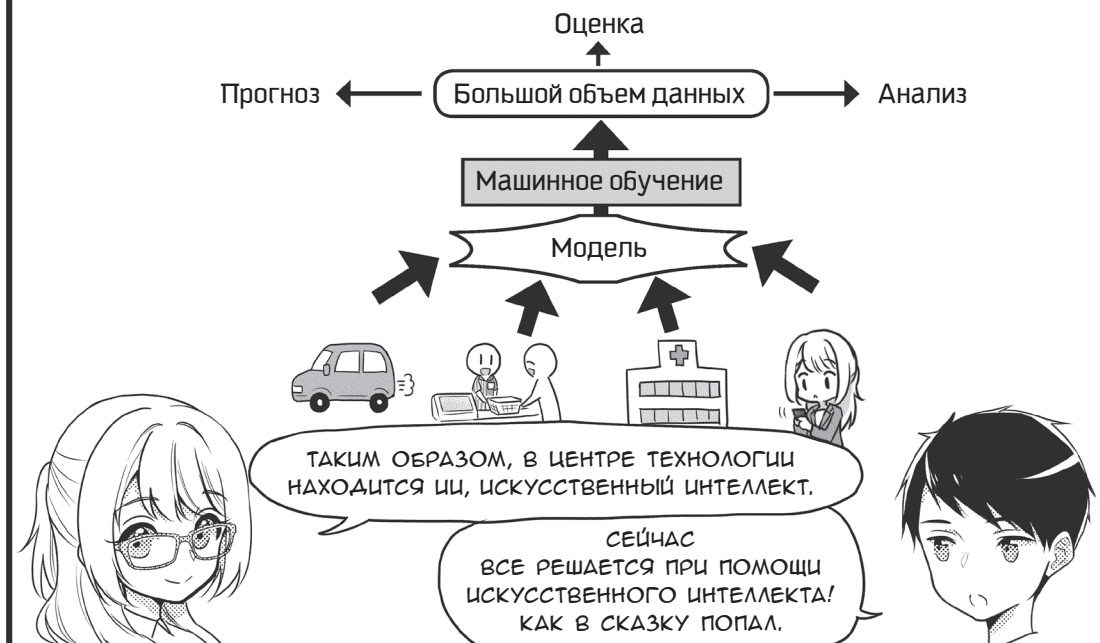
ПРЕЖДЕ ВСЕГО, КИЁХАРА-КУН,
ДАВАЙ ПРОВЕРИМ, ЧТО ТЫ ЗНАЕШЬ
О МАШИННОМ ОБУЧЕНИИ.

НУ... ЭТО КОГДА ОН АНАЛИЗИРУЕТ
БОЛЬШОЙ ОБЪЕМ ДАННЫХ
И ДАЕТ ОТВЕТ?

ОЧКИ...

НУ ЧТО, КИЁХАРА...
ИТАК, МАШИННОЕ ОБУЧЕНИЕ - ЭТО

ПОСТРОЕНИЕ НА ОСНОВАНИИ БОЛЬШОГО ОБЪЕМА ДАННЫХ
МОДЕЛИ, КОТОРАЯ МОЖЕТ ОЦЕНИВАТЬ И ДЕЛАТЬ ПРОГНОЗЫ.



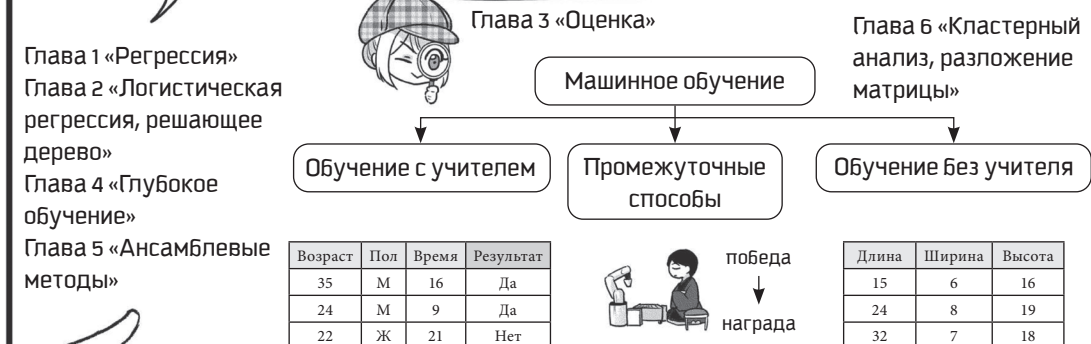
ИИ СЕЙЧАС ДОСТАТОЧНО РАСПРОСТРАНЕН! В ЦЕЛОМ СЧИТАЕТСЯ, ЧТО ОН ВСКОРЕ ЗАМЕНИТ ЛЮДЕЙ ПРИ ВЫПОЛНЕНИИ НЕКОТОРЫХ ПРОСТЫХ ИНТЕЛЛЕКТУАЛЬНЫХ ЗАДАНИЙ, НО В ДРУГИХ СИТУАЦИЯХ ОН ПОМОЖЕТ РАСШИРИТЬ ВОЗМОЖНОСТИ ЧЕЛОВЕЧЕСКОГО УМА.

ОГО.

БОЛЕЕ ТОГО, "МАШИННОЕ ОБУЧЕНИЕ" ТЕСНО СВЯЗАНО С ТЕХНОЛОГИЕЙ DATA MINING (ДОБЫЧА ДАННЫХ), КОТОРАЯ ПОЗВОЛЯЕТ ПОЛУЧИТЬ НУЖНЫЕ ДАННЫЕ В РЕЗУЛЬТАТЕ АНАЛИЗА ОГРОМНОГО, НЕПРЕДСТАВИМОГО ДЛЯ ЧЕЛОВЕЧЕСКОГО РАЗУМА ОБЪЕМА ДАННЫХ.

Data mining – метод обнаружения скрытых паттернов в огромном объеме данных при помощи статистики и математических методов.

ПОСКОЛЬКУ С ПОМОЩЬЮ МАШИННОГО ОБУЧЕНИЯ РЕШАЕТСЯ ОГРОМНОЕ КОЛИЧЕСТВО ЗАДАЧ, ЕГО ОБОЗРЕТЬ В ЦЕЛОМ ТРУДНО, ОДНАКО МЕТОДЫ ЛЕГКО РАЗДЕЛИТЬ ПРИМЕРНО НА ТРИ ГРУППЫ: **ОБУЧЕНИЕ С УЧИТЕЛЕМ; ПРОМЕЖУТОЧНЫЕ СПОСОБЫ; ОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ.**





НАПРИМЕР, КИЁХАРА, ТВОЯ ЗАДАЧА ОПРЕДЕЛЕНИЯ КОЛИЧЕСТВА ГОСТЕЙ, КОТОРЫЕ ПРИДУТ НА МЕРОПРИЯТИЕ, НА ОСНОВАНИИ ИМЕЮЩИХСЯ ДАННЫХ НАЗЫВАЕТСЯ ЗАДАЧЕЙ РЕГРЕССИИ.

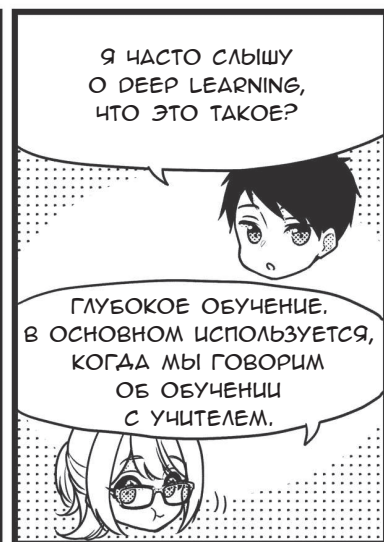
А ЕСЛИ ТЫ ХОЧЕШЬ НАЙТИ ОТВЕТ НА ВОПРОС, КУПЯТ ЛИ КАКОЙ-НИБУДЬ ТОВАР, ТО ЭТО ЗАДАЧА КЛАССИФИКАЦИИ.

Данные для задачи регрессии

Кол-во комнат	Время ходьбы от станции, мин	Возраст дома	Арендная плата
1	15	6	48 000
2	2	2	60 000
3	20	25	50 000

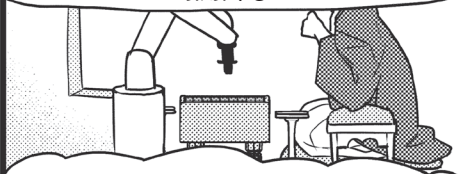
Данные для задачи классификации

Возраст	Пол	Время	Купит?
35	М	16	Да
24	М	9	Да
22	Ж	21	Нет



ПРОЛОГ. ПОГОВОРИМ О МАШИННОМ ОБУЧЕНИИ

ЗАТЕМ ИДУТ ПРОМЕЖУТОЧНЫЕ МЕТОДЫ. КИЁХАРА-КУН, ТЫ СЛЫШАЛ КОГДА-НИБУДАЬ НОВОСТИ О ТОМ, ЧТО ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ ОБЫГРАЛ ЧЕЛОВЕКА В ШАХМАТЫ ИЛИ ГО?



АГА, ЭТО БЫЛО ВОСХИТИТЕЛЬНО. ХОТЯ РАНЬШЕ КАЗАЛОСЬ, ЧТО НИ ЭТОГО НЕ СМОЖЕТ.

МЕТОД, ИСПОЛЗУЕМЫЙ ИСКУССТВЕННЫМ ИНТЕЛЛЕКТОМ ДЛЯ ИГРЫ В ГО ИЛИ ШАХМАТЫ, НАЗЫВАЕТСЯ ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ.



王将	将	将	将	将	将	将	将	王将
	王将						王将	
歩兵	歩兵	歩兵	歩兵	歩兵	歩兵	歩兵	歩兵	歩兵
	角行						飛車	
龍馬	桂馬	銀将	金将	玉将	金将	銀将	桂馬	龍馬

Победа –
положительное
подкрепление

Проигрыш –
отрицательное
подкрепление

После каждого хода неизвестно, какой следующий ход лучше сделать*

ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ НЕ ДАЕТ ОТВЕТА НА ВОПРОС, КАКОЙ ХОД ЛУЧШЕ СДЕЛАТЬ, НО ВМЕСТО ЭТОГО ВЫДАЕТ ПОДКРЕПЛЕНИЕ В ЗАВИСИМОСТИ ОТ ПОБЕДЫ ИЛИ ПОРАЖЕНИЯ.

ЛЮДЯМ ТОЖЕ МОЖНО ДАВАТЬ ПОДКРЕПЛЕНИЕ.



АА, ЭТИ СЛАДОСТИ – КАК РАЗ ПОДКРЕПЛЕНИЕ ДЛЯ МЕНЯ.

НА ОСНОВАНИИ ПОДКРЕПЛЕНИЯ ОПРЕДЕЛЯЕТСЯ ОПТИМАЛЬНЫЙ ПОРЯДОК ДЕЙСТВИЙ. ТАК ОБУЧАЮТ РОБОТОВ ДЛЯ ВОЖДЕНИЯ АВТОМОБИЛЕЙ.



ОГО... КАК УДОБНО. ВОИСТИНУ XXI ВЕК!

ШУТКИ В СТОРОНУ!

* Приведенная таблица – расстановка для игры в сёги, японские шахматы. – Прим. перев.

ЧТО ЖЕ КАСАЕТСЯ ОБУЧЕНИЯ
БЕЗ УЧИТЕЛЯ, ТО ЖЕЛАЕМОГО ОТВЕТА
НА ВОПРОСЫ СРЕДИ ДАННЫХ НЕТ.

Обучение без учителя


Неразмеченные данные



Длина	Ширина	Высота
15	6	16
24	8	19
32	7	18

А КАК ТОГДА
ПРОХОДИТ ОБУЧЕНИЕ?

НУ, КАК ОБЫЧНО.
ЦЕЛЬ СИСТЕМ БЕЗ УЧИТЕЛЯ -
ЭТО ОБНАРУЖИТЬ В БОЛЬШОМ
ОБЪЕМЕ ДАННЫХ ЗНАНИЯ,
КОТОРЫЕ МОГУТ ПРИГОДИТЬСЯ
ЧЕЛОВЕКУ.




ОНИ ИСПОЛЗУЮТСЯ, НАПРИМЕР,
ДЛЯ РЕКОМЕНДАЦИИ ТОВАРОВ ПРИ ПОКУПКЕ
В ИНТЕРНЕТ-МАГАЗИНАХ ИЛИ ПРИ ПОИСКЕ
СТРАННОСТЕЙ В ИНФОРМАЦИИ
О ДЕЙСТВИЯХ МЕХАНИЗМОВ.

История покупок

ID	#1	#2	#3	#4
115		1		
124		1		1
232				1

Рекомендуемые
артисты



ПОКА ЧТО, НАДЕЮСЬ,
ТЫ ВСЕ ПОНЯЛ ОТНОСИТЕЛЬНО
МАШИННОГО ОБУЧЕНИЯ?

ДА...

ЭТО УЖЕ
ВОСЬМОЙ.



В кабинете у Саяка (1)

Саяка и старшеклассница Ай



Давно не виделись, Ай-тян. С тех пор как мы были у дедушки, да?

Возможно, тогда собирались все его внуки. Чем ты сегодня занималась?



Рассказывала младшему товарищу про машинное обучение. Он хоть и учился у нас, но, боюсь, очень мало что понял... Ай-тян, ты же хочешь учиться в математическом классе, может, тоже зайдешь послушать?

Машинное обучение? ИИ? Это, наверное, очень сложно!



Там используются математические модели, но это все делается компьютерами. А основы этих моделей может понять и старшеклассник.

Я писала программы для астрономических расчетов. Математика – мой любимый предмет, но пойму ли я?



Конечно, поймешь. Мы начнем с регрессии. Приходи послушать!

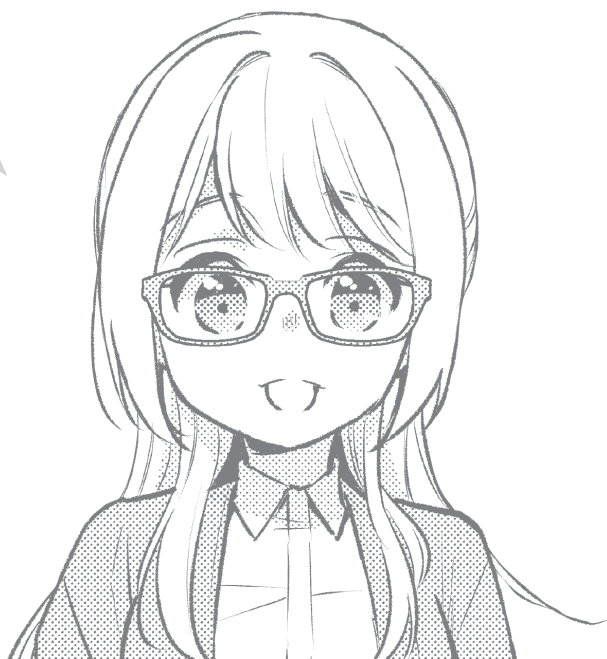
Хорошо, если что-то будет трудно, я буду задавать вопросы!



ГЛАВА 1

ЧТО ТАКОЕ РЕГРЕССИЯ

ЛИНЕЙНАЯ РЕГРЕССИЯ!
РЕГУЛЯРИЗАЦИЯ!



ДЛЯ НАЧАЛА
ПОГОВОРИМ
О РЕГРЕССИИ.

ПОЖАЛУЙСТА!



В КАЧЕСТВЕ ПРИМЕРА
ВОЗЬМЕМ СИТУАЦИЮ, КОГДА НАДО
ПОДСЧИТАТЬ КОЛИЧЕСТВО ГОСТЕЙ
НА ПИАР-МЕРОПРИЯТИИ, КОТОРОЕ
ОРГАНИЗОВЫВАЕТ ГОРОД.

НА НЕМ БУДЕТ ПОДАВАТЬСЯ СОК,
ВЫЖАТЫЙ ИЗ МЕСТНЫХ ФРУКТОВ,
ПОЭТОМУ НАДО КАК МОЖНО ТОЧНЕЕ
ПОДСЧИТАТЬ КОЛИЧЕСТВО
УЧАСТНИКОВ.

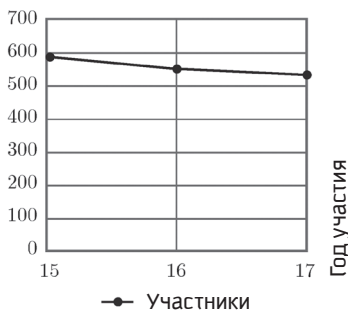
АГА!



1.1. СЛОЖНОСТИ С ПРОГНОЗОМ

ДЛЯ НАЧАЛА ПОСМОТРИМ
НА ЭТОТ ГРАФИК.

Количество участников

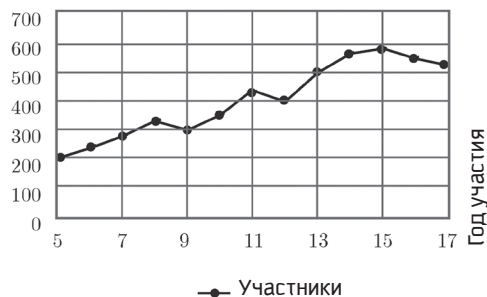


ЭТОТ ГРАФИК ПОКАЗЫВАЕТ
ИЗМЕНЕНИЕ ЧИСЛА УЧАСТНИКОВ
ЗА ТРИ ГОДА. ЕСЛИ МЫ
ПРЕДСКАЖЕМ КОЛИЧЕСТВО
УЧАСТНИКОВ НА ЕГО ОСНОВЕ,
ЧТО ПОЛУЧИТСЯ?

ЗА ТРИ ГОДА
ДЕМОНСТРИРУЕТСЯ
ТЕНДЕНЦИЯ К СНИЖЕНИЮ.

ДА, А ЕСЛИ МЫ ВЗГЛЯНЕМ
НА ГРАФИК УЧАСТНИКОВ
С 13-ГО ГОДА?

Количество участников



В ДОЛГОСРОЧНОЙ
ПЕРСПЕКТИВЕ ВИАНА
ТЕНДЕНЦИЯ К УВЕЛИЧЕНИЮ
ЧИСЛА УЧАСТНИКОВ.

ИМЕННО ТАК! В ЗАВИСИМОСТИ ОТ ТОГО, КАКОЙ ПЕРИОД МЫ ВЗЯЛИ, ТЕНДЕНЦИИ АБСОЛЮТНО ПРОТИВОПОЛОЖНЫЕ.



И ЭТО ТОТ СПОСОБ, КОТОРЫЙ ИСПОЛЬЗОВАЛ КОНСУЛЬТАНТ?

ДА.

ДАЖЕ ЕСЛИ ТЕНДЕНЦИЯ И ПРАВИЛЬНАЯ, ГРАФИК НЕОБЯЗАТЕЛЬНО БУДЕТ ПОХОДИТЬ НА ПРЯМУЮ.

ВОООООУЩЕВЛЕННО



1.2. ОПРЕДЕЛЯЕМ ЗАВИСИМЫЕ И НЕЗАВИСИМЫЕ ПЕРЕМЕННЫЕ

ЧТО Ж, КИЁХАРА, КАК МЫ ТЕПЕРЬ БУДЕМ ПРЕДСКАЗЫВАТЬ КОЛИЧЕСТВО ГОСТЕЙ?

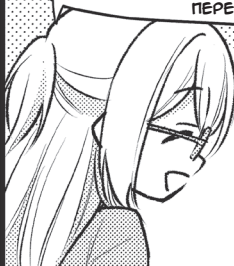
СТАРАТЕЛЬНО ЗАПИСЫВАЕТ



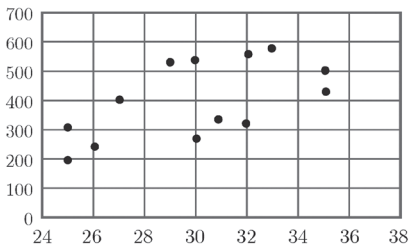
НУ... Я ДУМАЮ, ВЛИЯНИЕ ОКАЗЫВАЮТ ПОГОДА, ТЕМПЕРАТУРА, В ОСОБЕННОСТИ КОЛИЧЕСТВО ОСАДКОВ В СЕЗОН ДОЖДЕЙ...

А ТЕПЕРЬ ПОПРОБУЕМ ВЫБРАТЬ МЕТОД ПРЕДСКАЗАНИЯ В ЗАВИСИМОСТИ ОТ ФАКТОРОВ. ТЕ РЕЗУЛЬТАТЫ, КОТОРЫЕ МЫ ХОТИМ СПРОГНОЗИРОВАТЬ, НАЗЫВАЮТСЯ ЗАВИСИМЫМИ ПЕРЕМЕННЫМИ, А ФАКТОРЫ, ВЛИЯЮЩИЕ НА РЕЗУЛЬТАТ, - НЕЗАВИСИМЫМИ ПЕРЕМЕННЫМИ.

ДА!



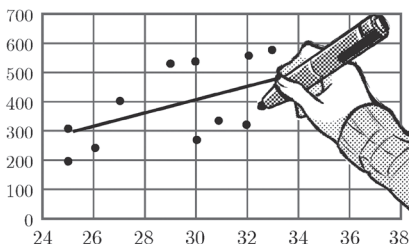
ПУСТЬ У НАС БУДЕТ ОДНА НЕЗАВИСИМАЯ ПЕРЕМЕННАЯ. РАССМОТРИМ ВЛИЯНИЕ ТЕМПЕРАТУРЫ В ДЕНЬ МЕРОПРИЯТИЯ НА КОЛИЧЕСТВО ГОСТЕЙ.



НА ЭТОМ ГРАФИКЕ ГОРИЗОНТАЛЬНАЯ ОСЬ - НЕЗАВИСИМАЯ ПЕРЕМЕННАЯ (ПОГОДА), А ВЕРТИКАЛЬНАЯ - КОЛИЧЕСТВО УЧАСТНИКОВ (ЗАВИСИМАЯ ПЕРЕМЕННАЯ). ОТМЕТИМ ТОЧКИ, РАВНЫЕ КОЛИЧЕСТВУ УЧАСТНИКОВ, И ПРОВЕДЕМ ПО НИМ ПРЯМУЮ.

ЭТО НЕВОЗМОЖНО?

КАК-ТО ТАК. ПО ВОЗМОЖНОСТИ ПРОВОДИМ ЕЕ НЕДАЛЕКО ОТ ВСЕХ ТОЧЕК. КОГДА МЫ ПРОВОДИМ ЛИНИЮ, ОНА ВЫГЛЯДИТ ТАК:



СКРИП СКРИП

ИТАК, ОБОЗНАЧИМ НАКЛОН ГРАФИКА w_1 ,
СДВИГ ГРАФИКА ОТНОСИТЕЛЬНО
ВЕРТИКАЛЬНОЙ ОСИ - w_0 , ТЕМПЕРАТУРУ ЗА x ,
А КОЛИЧЕСТВО УЧАСТНИКОВ - ЗА y

$$y = w_1x + w_0$$

- И ПОЛУЧИМ ЭТУ ФОРМУЛУ.

ЕСЛИ ТЕМПЕРАТУРА РАСТЕТ,
ТО И КОЛИЧЕСТВО УЧАСТНИКОВ
РАСТЕТ, А ЭТО НЕПОХОЖЕ
НА ПРАВДУ.

АГА.
ВИДИМО, НЕЗАВИСИМЫХ
ПЕРЕМЕННЫХ НЕСКОЛЬКО.

НЕСКОЛЬКО?
ЭТО УЖАСНО?

СЛОЖНОСТЬ В ТОМ, ЧТО СРЕДИ
НЕЗАВИСИМЫХ ПЕРЕМЕННЫХ ЕСТЬ ТАКИЕ,
КОТОРЫЕ ОКАЗЫВАЮТ ВЛИЯНИЕ
НА ЗАВИСИМЫЕ ПЕРЕМЕННЫЕ,
А ЕСТЬ И ТАКИЕ, ВЛИЯНИЕ КОТОРЫХ
НЕВЕЛИКО.

Влияние велико

Влияние мало

ПОГОДА

ВЛАЖНОСТЬ

ЧТОБЫ ПРОСТО СМОДЕЛИРОВАТЬ ТАКИЕ УСЛОВИЯ,
НЕОБХОДИМ СПОСОБ, КОТОРЫЙ БУДЕТ УПОРЯДОЧИВАТЬ
ЗАВИСИМЫЕ ПЕРЕМЕННЫЕ В ЗАВИСИМОСТИ ОТ УЧИТЫ-
ВАЕМОГО **ВЕСА** НЕЗАВИСИМЫХ ПЕРЕМЕННЫХ.

ЭТО
ЛИНЕЙНАЯ РЕГРЕССИЯ.

Температура × вес → учитывание
Влажность × вес → цена
Зависимая переменная

Что такое линейная регрессия?

АГА! ВОТ И РЕГРЕССИЯ!

ЕСЛИ НЕЗАВИСИМЫХ ПЕРЕМЕННЫХ ДВЕ,
ТО ОБОЗНАЧИМ ИХ КАК x_1 И x_2 ,
ИХ ВЕС КАК w_1 И w_2 СООТВЕТСТВЕННО,
А ИХ ВЗВЕШЕННАЯ СУММА ВМЕСТЕ
С ПОСТОЯННОЙ СОСТАВЛЯЮЩЕЙ w_0
ЗАПИСЫВАЕТСЯ...

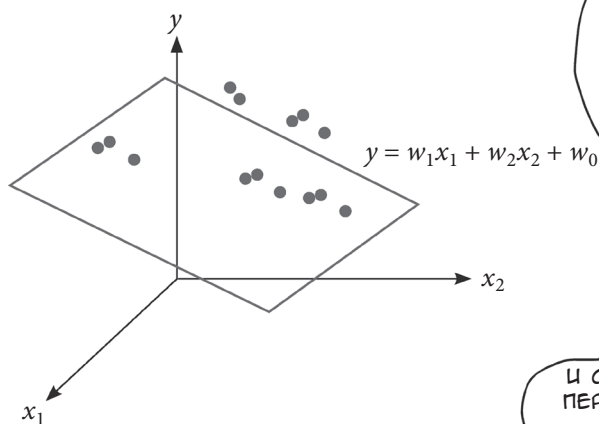
$$y = w_1x_1 + w_2x_2 + w_0$$

...ВОТ ТАК.

ОЧЕНЬ ПОХОЖЕ
НА НЕДАВНЮЮ ФОРМУЛУ.

СКРИП
СКРИП

МЫ ЖЕ СТРОИЛИ
ДВУМЕРНЫЙ ГРАФИК?



МОЖНО ПОСТРОИТЬ ТРЕХМЕРНЫЙ
ГРАФИК, КАК НА КАРТИНКЕ,
И СПРОГНОЗИРОВАТЬ КОЛИЧЕСТВО
УЧАСТНИКОВ. ПОЛУЧИТСЯ ТАКАЯ
ВОТ ПОВЕРХНОСТЬ.

ФУХ! 000

И С КАЖДОЙ НОВОЙ НЕЗАВИСИМОЙ
ПЕРЕМЕННОЙ БУДЕТ УВЕЛИЧИВАТЬСЯ
КОЛИЧЕСТВО ИЗМЕРЕНИЙ?

ВСЕ ТАК... ОБОБЩАЕМ: ЕСЛИ
В ДАННЫХ СУЩЕСТВУЕТ
 d НЕЗАВИСИМЫХ ПЕРЕМЕННЫХ,
ТО ТОЧКИ НА ГРАФИКЕ БУДУТ
РАСПОЛАГАТЬСЯ В $(d+1)$ -МЕРНОМ
ПРОСТРАНСТВЕ, И ПОЭТОМУ НАМ НАДО
РЕШИТЬ ЗАДАЧУ НАХОЖДЕНИЯ
 d -МЕРНОЙ ГИПЕРПЛОСКОСТИ.

ЧТО?

ТО ЕСТЬ ЕСЛИ НЕЗАВИСИМЫХ ПЕРЕМЕННЫХ
БУДЕТ 10, ТО НАМ НУЖНО БУДЕТ НАЙТИ
ДЕСЯТИМЕРНУЮ ГИПЕРПЛОСКОСТЬ?
А ЭТО ВООООЩЕ ВОЗМОЖНО?

ДА! ПРИДЕТСЯ ПОПРОСИТЬ
ПОМОЩИ У МАТЕМАТИКИ!

1.3. НАХОДИМ ФУНКЦИЮ ЛИНЕЙНОЙ РЕГРЕССИИ

Шаг 1

Пусть у нас есть d независимых переменных. Тогда обозначим d -мерный столбец-вектор \mathbf{x} . В соответствии с этим вес, отличный от постоянной w_0 , также будет обозначен d -мерным вектором-столбцом \mathbf{w} , и уравнение гиперплоскости примет такой вид:

$$y = \mathbf{w}^T \mathbf{x} + w_0, \quad (1.1)$$

где T – обозначение транспонирования.

Шаг 2

До этого мы обозначили за y величину взвешенной суммы независимых переменных, а теперь обозначим так величину зависимых переменных в данных для обучения. Взвешенную сумму независимых переменных мы обозначим как $\hat{c}(\mathbf{x})$. Значок над буквой означает, что мы не можем гарантировать правильность полученных данных. Таким образом уравнение (1.1) примет вид (1.2):

$$\hat{c}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0. \quad (1.2)$$



Тот факт, что функция линейной регрессии не слишком отличается от использованных данных, означает, что величина линейной функции $\hat{c}(\mathbf{x})$, куда входят независимые переменные \mathbf{x} , также мало отличается от величины зависимой переменной y . Цель в том, чтобы сделать эту разницу как можно меньше. Однако если эта разница проявляется в наборе данных, то случаи, где величина зависимой переменной выше величины линейной функции, накладываются на те, где величина зависимой переменной ниже величины линейной функции, и они компенсируют друг друга.

Шаг 3

Для определения «отклонения» линейной функции от имеющихся данных возводим в квадрат разницу между линейной функцией $\hat{c}(\mathbf{x})$ и зависимой переменной y ; т.е. находим **квадрат ошибки**. Уменьшение квадрата ошибки путем корректировки веса линейной функции называется обучением по **методу наименьших квадратов**.

Таким образом, добавив к \mathbf{x} в уравнении (1.2) 0-мерность и определив его величину как равную 1, а также добавив w_0 к 0-мерному \mathbf{w} , получим, что функция регрессии будет записываться как внутреннее произведение $(d + 1)$ -мерного вектора (1.3).

$$\hat{c}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}. \quad (1.3)$$

Шаг 4

Оценим коэффициенты \mathbf{w} этого уравнения, используя обучающие данные $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. Постараемся максимально возможно уменьшить разницу между величинами зависимой переменной y и функции $\hat{c}(\mathbf{x})$, рассчитанной по уравнению 1.3. Ошибка определяется значением коэффициентов \mathbf{w} в уравнении (1.3) обозначив ее $E(\mathbf{w})$, получим следующее уравнение:

$$E(\mathbf{w}) = \sum_{i=1}^n (y - \hat{c}(\mathbf{x}))^2, \quad (1.4)$$

$$= \sum_{i=1}^n (y - \mathbf{w}^T \mathbf{x}_i)^2. \quad (1.5)$$

Шаг 5

Чтобы избавиться от трудоемких вычислений суммы, представим независимые переменные матрицами, а зависимые – векторами. Обозначим матрицу, имеющую n позиций по вертикали, которая получилась путем транспонирования независимой переменной \mathbf{x} d -мерного вектора-столбца как \mathbf{X} , y – вектор-столбец величин независимой переменной y , а \mathbf{w} – вектор-столбец коэффициентов.

В итоге отклонение примет следующий вид:

$$E(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}). \quad (1.6)$$

Чтобы минимизировать отклонение, нужно найти такие величины коэффициентов \mathbf{w} , при которых производная функции ошибки равняется 0, то есть:

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0, \quad (1.7)$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (1.8)$$

где \mathbf{A}^{-1} – матрица, обратная матрице \mathbf{A} .



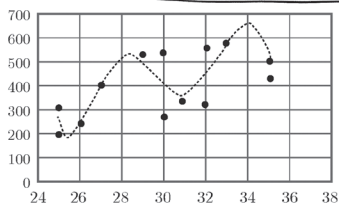
Таким образом, мы можем аналитически найти по обучающим данным веса \mathbf{w} при помощи минимизации суммы квадратов ошибки. При подстановке \mathbf{w} в уравнение (1.3) получится функция линейной регрессии $\hat{c}(\mathbf{x})$.

1.4. РЕГУЛЯРИЗАЦИЯ РЕЗУЛЬТАТА

ЕСЛИ КОЭФФИЦИЕНТЫ ЛИНЕЙНЫЕ, ТО ВЕСА МОГУТ БЫТЬ РАССЧИТАНЫ ТАКИМ ЖЕ ОБРАЗОМ, ДАЖЕ ЕСЛИ ЭТО УРАВНЕНИЯ ВЫСОКОГО ПОРЯДКА. ТАК МОЖНО ПРОВОДИТЬ ОБУЧЕНИЕ ПРИ ПОМОЩИ САМЫХ СЛОЖНЫХ УРАВНЕНИЙ РЕГРЕССИИ.



ЕСЛИ Я ЗНАЮ ЧИСЛЕННЫЙ ВЕС, ТО ОН МОЖЕТ БЫТЬ ПРОСТО ПОДАСТАВЛЕН В ЭТУ ФОРМУЛУ?



НЕ СОВСЕМ ТАК. В ЭТОЙ ФОРМУЛЕ ДАЖЕ ПРИ НЕБОЛЬШОМ ИЗМЕНЕНИИ ДАННЫХ НА ВХОДЕ РЕЗУЛЬТАТЫ НА ВЫХОДЕ МОГУТ СИЛЬНО ОТЛИЧАТЬСЯ, И НЕЛЬЗЯ БУДЕТ ПОЛУЧИТЬ ХОРОШИЙ РЕЗУЛЬТАТ ПРИ ИСПОЛЬЗОВАНИИ ДАННЫХ, ОТЛИЧНЫХ ОТ ОБУЧАЮЩИХ.

СЛОЖНОВАТО. А КАК СДЕЛАТЬ ХОРОШИЙ ПРОГНОЗ?



ЕСЛИ НЕБОЛЬШОЕ ИЗМЕНЕНИЕ КАКОГО-ЛИБО ПАРАМЕТРА НА ВХОДЕ ДАЕТ БОЛЬШОЕ ИЗМЕНЕНИЕ НА ВЫХОДЕ, ЭТО ОЗНАЧАЕТ, ЧТО КОЭФФИЦИЕНТ УРАВНЕНИЯ СЛИШКОМ ВЕЛИК, ПОЭТОМУ ЕГО НЕОБХОДИМО УМЕНЬШИТЬ.

ТАК БЫСТРО?

НО, ВПРОЧЕМ, ЕСТЬ И ДРУГАЯ ТОЧКА ЗРЕНИЯ НА ВЕС, В СЛУЧАЯХ, КОГДА ПРАВИЛЬНОСТЬ ПРОГНОЗА ВАЖНЕЕ, ЧЕМ ОБЪЯСНИМОСТЬ РЕЗУЛЬТАТОВ ОБУЧЕНИЯ.

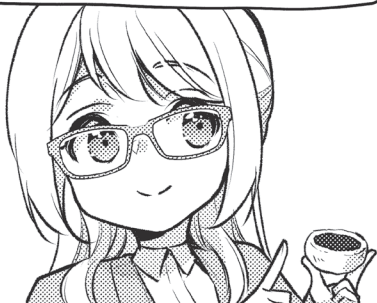


ОБЪЯСНИМОСТЬ ВАЖНЕЕ ТОЧНОСТИ?

С НАВИТЫМ РТОМ



ДОПУСТИМ, ВМЕСТО ТОГО ЧТОБЫ ПРОГНОЗИРОВАТЬ, КАКИЕ ФАКТОРЫ ПОВЛИЯЮТ НА КАЧЕСТВО ТОВАРА, НАДО НАЙТИ ТЕ, КОТОРЫЕ ВЛИЯЮТ НА НЕГО БОЛЬШЕ ВСЕГО.



ТАК МОЖНО НАЙТИ ТЕ ХАРАКТЕРИСТИКИ, КОТОРЫЕ ДЕЛАЮТ ТОВАР БРАКОВАННЫМ.

Пирожок
Составляющие:

- мука
- начинка
- сахар



Что влияет больше?

ЕСЛИ ГОВОРИТЬ КОНКРЕТНЕЕ, МОЖНО ПОДАСТАВИТЬ В КАЧЕСТВЕ ВЕСА ПЕРЕМЕННЫХ В ФОРМУЛУ ЛИНЕЙНОЙ РЕГРЕССИИ 0 И ПОСМОТРЕТЬ, ИЗМЕНИТСЯ ЛИ КОЛИЧЕСТВО ИЗМЕРЕНИЙ.



ДРУГИМИ СЛОВАМИ, СЛЕДУЕТ НАЙТИ СПОСОБ, ЧТОБЫ КОЭФФИЦИЕНТ w В УРАВНЕНИИ ЛИНЕЙНОЙ РЕГРЕССИИ УМЕНЬШИЛСЯ, ЕСЛИ ВЕЛИЧИНА ЕГО БОЛЬШАЯ, ИЛИ ЖЕ СТАЛ РАВЕН НУЛЮ.

И ЭТОТ СПОСОБ НАЗЫВАЕТСЯ РЕГУЛЯРИЗАЦИЕЙ.

Регуляризация

ОН НУЖЕН, ЧТОБЫ КОЭФФИЦИЕНТ НЕ БЫЛ СЛИШКОМ БОЛЬШИМ...

МЕТОД РЕГУЛЯРИЗАЦИИ, ПРИ КОТОРОМ МЫ УМЕНЬШАЕМ БОЛЬШИЕ КОЭФФИЦИЕНТЫ, НАЗЫВАЕТСЯ РИДЖ-РЕГРЕССИЕЙ,

А ЧТОБЫ УВЕЛИЧИТЬ КОЛИЧЕСТВО ВЕЛИЧИН, РАВНЫХ НУЛЮ, ИСПОЛЬЗУЕТСЯ ЛАССО-РЕГРЕССИЯ.

РИДЖ-РЕГРЕССИЯ

ЛАССО-РЕГРЕССИЯ

РЕГУЛЯРИЗАЦИЯ ОСУЩЕСТВЛЯЕТСЯ ДОБАВЛЕНИЕМ ДОПОЛНИТЕЛЬНОГО ЧЛЕНА К УРАВНЕНИЮ ОШИБКИ.

НАЧНЕМ С ОБЪЯСНЕНИЯ РИДЖ-РЕГРЕССИИ.

МЫ ДОБАВЛЯЕМ ДОПОЛНИТЕЛЬНЫЙ ЧЛЕН, КВАДРАТ ПАРАМЕТРА w .

Чтобы уменьшить это...

...надо увеличить это

$$E(w) = (y - Xw)^T (y - Xw) + \alpha w^T w$$

Регулируем баланс

Если величина веса слишком маленькая...

...то она далека от правильной

А α ОТКУДА ВЗЯЛАСЬ?

α - ЭТО ВЕС ДОПОЛНИТЕЛЬНОГО ЧЛЕНА РЕГУЛЯРИЗАЦИИ. ЕСЛИ ПАРАМЕТР БОЛЬШОЙ, ТО ЭФФЕКТ РЕГУЛЯРИЗАЦИИ СТАНОВИТСЯ ВАЖНЕЕ ЭФФЕКТИВНОСТИ, ЕСЛИ МАЛЕНЬКИЙ, ТО ЭФФЕКТИВНОСТЬ СТАНОВИТСЯ ВАЖНЕЕ.



ИСПОЛЬЗУЯ РИДЖ-РЕГРЕССИЮ, МЫ НАХОДИМ ВЕЛИЧИНУ w , КОГДА ГРАДИЕНТ ПО w ФУНКЦИИ ОШИБКИ РАВЕН 0, КАК И В СЛУЧАЕ НАХОЖДЕНИЯ ЭТОГО ПАРАМЕТРА МЕТОДОМ НАИМЕНЬШИХ КВАДРАТОВ.

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

I - единичная матрица.



КСТАТИ, А ПОЧЕМУ РИДЖ-РЕГРЕССИЯ ТАК НАЗЫВАЕТСЯ?



RIDGE ОЗНАЧАЕТ ГРЕБЕНЬ ГОРЫ, И ЕДИНИЧНАЯ МАТРИЦА НА НЕГО ПОХОЖА*.

* Есть другие версии.

ХМ...



КАК Я И ГОВОРИЛА, РИДЖ-РЕГРЕССИЯ - ЭТО РЕГУЛЯРИЗАЦИЯ ДЛЯ УМЕНЬШЕНИЯ ВЕЛИЧИНЫ ПАРАМЕТРА.

А ТЕПЕРЬ ПОГОВОРИМ О РЕГРЕССИИ "ЛАССО". ЭТО РЕГУЛЯРИЗАЦИЯ, ПРИ КОТОРОЙ w СТАНОВИТСЯ АБСОЛЮТНОЙ ВЕЛИЧИНОЙ.



RIDGE - ЭТО КВАДРАТ w , ЛАССО - ЭТО АБСОЛЮТНАЯ ВЕЛИЧИНА w .

А ЧТО ТАКОЕ ЛАССО?

ЛАССО - ЭТО ПЕТЛЯ ДЛЯ ЛОВЛИ КОГО-НИБУДЬ.



ВРОДЕ ЭТО СЛОВО ИСПОЛЬЗУЮТ, КОГДА ГОВОРЯТ О КОВЕОЯХ.

ПРЕДСТАВЬ, ЧТО ВО МНОЖЕСТВО ПАРАМЕТРОВ КИДАЮТ ЛАССО И ВЫБИРАЮТ САМЫЕ МАЛЕНЬКИЕ ИЗ НИХ.

ИЗНАЧАЛЬНО ЭТО АББРЕВИАТУРА ФРАЗЫ **LEAST ABSOLUTE SHRINKAGE AND SELECTION OPERATOR**.

НАДО ЖЕ.

УРАВНЕНИЕ ОЦЕНКИ ОШИБОК РЕГРЕССИИ "ЛАССО" ВЫГЛЯДИТ...

$$E(w) = (y - Xw)^T(y - Xw) + \alpha \sum_{i=1}^n |w_j|$$

...ТАК:

ПОСКОЛЬКУ w_0 - ПОСТОЯННОЕ СЛАГАЕМОЕ В УРАВНЕНИИ, ЕГО ВЕЛИЧИНА НЕ ПОВЛИЯЕТ НА ВЕЛИЧИНУ УРАВНЕНИЯ РЕГРЕССИИ, И ЕГО ОБЫЧНО НЕ РЕГУЛЯРИЗИРУЮТ. ЗДЕСЬ ЕСЛИ ВЕС ДОПОЛНИТЕЛЬНОГО ЧЛЕНА

УВЕЛИЧИВАЕТСЯ, РАСТЕТ ЧИСЛО ВЕЛИЧИН С ВЕСОМ, РАВНЫМ 0.

А КАК МОЖНО ОБЪЯСНИТЬ РЕГРЕССИЮ "ЛАССО"?

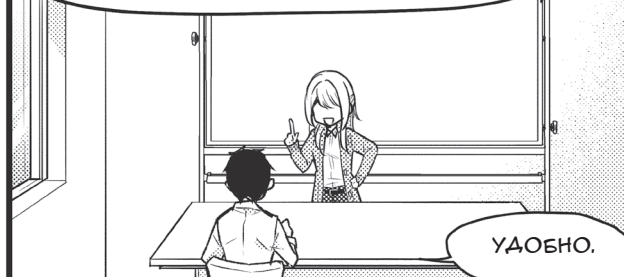
Функция
абсолютного
значения

Квадратичная
функция

ПОСКОЛЬКУ ТУДА ВХОДИТ АБСОЛЮТНАЯ ВЕЛИЧИНА, НЕДИФФЕРЕНЦИРУЕМАЯ В ТОЧКЕ НАЧАЛА КООРДИНАТ, НЕЛЬЗЯ НАЙТИ ЗНАЧЕНИЕ АНАЛИТИЧЕСКИ, ИСПОЛЬЗУЯ МЕТОД НАИМЕНЬШИХ КВАДРАТОВ, ПОЭТОМУ ВЕРХНИЙ

ПРЕДЕЛ/МАКСИМУМ ДОПОЛНИТЕЛЬНОГО ЧЛЕНА РЕГУЛЯРИЗАЦИИ ОГРАНИЧИВАЕТСЯ ДИФФЕРЕНЦИРУЕМОЙ КВАДРАТИЧНОЙ ФУНКЦИЕЙ. БЫЛ ПРЕДЛОЖЕН МЕТОД, ЧТОБЫ МНОГОКРАТНО ОБНОВЛЯТЬ ЕЕ ПАРАМЕТРЫ С ЦЕЛЬЮ УМЕНЬШЕНИЯ ОШИБКИ.

С ПОМОЩЬЮ МЕТОДА "ЛАССО" МОЖНО ПРОРЕДИТЬ НЕЗАВИСИМЫЕ ПЕРЕМЕННЫЕ С ВЕСАМИ, НЕ РАВНЫМИ НУЛЮ, И НАЙТИ ТЕ, КОТОРЫЕ ОКАЗЫВАЮТ ВЛИЯНИЕ.

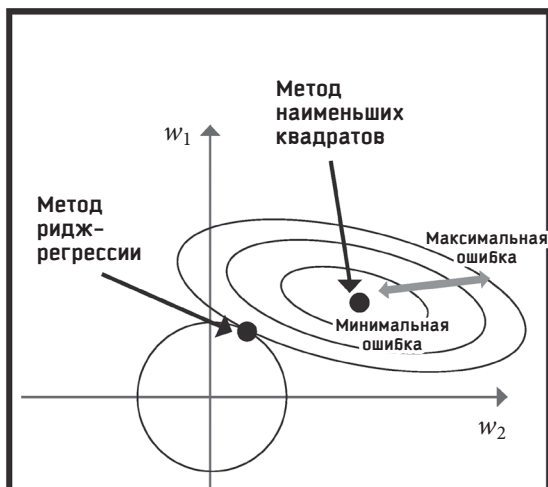


УДОБНО.

А ТЕПЕРЬ Я ОБЪЯСНЮ, ЧЕМ РИДЖ-РЕГРЕССИЯ ОТЛИЧАЕТСЯ ОТ РЕГРЕССИИ "ЛАССО".

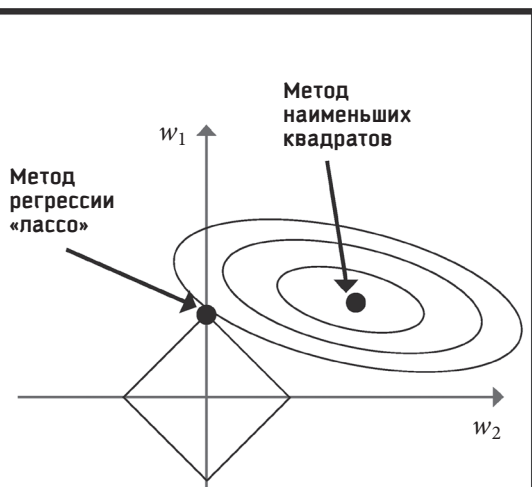


ПОЖАЛУЙСТА!



КАК ПОКАЗАНО НА РИСУНКЕ, ПРИ ИСПОЛЬЗОВАНИИ РИДЖ-РЕГРЕССИИ ОГРАНИЧЕНИЕ ДИАПАЗОНА ПАРАМЕТРОВ ОКРУЖНОСТЬЮ (ОБЩИЙ СЛУЧАЙ d -МЕРНОЙ ГИПЕРСФЕРЫ) НЕ ПОЗВОЛЯЕТ КАЖДОМУ ВЕСУ ПРИНИМАТЬ БОЛЬШОЕ ЗНАЧЕНИЕ. В ОБЩЕМ СЛУЧАЕ ТОЧКА КАСАНИЯ ИЗОЛИНИИ ФУНКЦИИ ОШИБКИ ЯВЛЯЕТСЯ ТОЧКОЙ НА ОКРУЖНОСТИ, КОТОРАЯ ЯВЛЯЕТСЯ ЗНАЧЕНИЕМ ВЕСА ПАРАМЕТРА.

ПОЭТОМУ ВЕЛИЧИНА ПАРАМЕТРА УМЕНЬШАЕТСЯ.



А В СЛУЧАЕ РЕГРЕССИИ "ЛАССО", ПРИ УСЛОВИИ ЧТО ОПРЕДЕЛЕНА СУММА ПАРАМЕТРОВ, ДИАПАЗОН ПАРАМЕТРОВ ОГРАНИЧЕН ОБЛАСТЬЮ (РОМБЕМ), УГЛЫ КОТОРОГО ЛЕЖАТ НА КАЖДОЙ ОСИ, КАК ПОКАЗАНО НА РИСУНКЕ.

И ОДИН ИЗ УГЛОВ РОМБА КАСАЕТСЯ ИЗОЛИНИИ ФУНКЦИИ ОШИБКИ.

КАЖЕТСЯ, В УГЛАХ БОЛЬШИНСТВО ПАРАМЕТРОВ СТАНОВЯТСЯ РАВНЫМИ 0.

ЭТО ВЛИЯНИЕ РЕГУЛЯРИЗАЦИИ РЕГРЕССИИ "ЛАССО".



А ТЕПЕРЬ ПОПРОБУЕМ ИСПОЛЬЗОВАТЬ ЯЗЫК ПРОГРАММИРОВАНИЯ PYTHON ДЛЯ РЕГРЕССИИ. В PYTHON МОЖНО ИСПОЛЬЗОВАТЬ БИБЛИОТЕКУ МАШИННОГО ОБУЧЕНИЯ SCIKIT-LEARN И ДЕЛАТЬ С ЕЕ ПОМОЩЬЮ ПРОГРАММЫ.

НА PYTHON
Я ПРОГРАММИРУЮ ПЛОХО...



Для начала загрузим библиотеку. В scikit-learn подготовлено несколько наборов данных, выберем методы из пакета datasets. Для регрессии это линейная регрессия, ридж-регрессия и регрессия «лассо».

```
from sklearn.datasets import load_boston
from sklearn.linear_model import LinearRegression, Ridge, Lasso
```



В качестве данных для анализа мы возьмем 13 параметров из стандартной выборки boston dataset, куда входят уровень преступности, количество комнат, географическое положение и прочие данные, связанные с недвижимостью.



Атрибут data экземпляра boston, который создан при помощи приведенного ниже кода, является матрицей, в которой признаковое описание объекта располагается в виде столбцов (для 13-мерного признака будет 506 векторов-строк), а атрибут target будет введен в качестве вектора-столбца – цены каждого свойства.



Можно показать детали данных boston с помощью атрибута descr функцией print (boston.DESCR).

```
boston = load_boston()
X = boston.data
y = boston.target
```



Код в scikit-learn удобен для обучения. А теперь используем учебный набор данных.

```
lr1 = LinearRegression()
```



В этом экземпляре можно вызвать метод `fit`, который выполняет обучение с набором признаков X и точной информацией y в качестве аргументов.

```
lr1.fit(X, y)
```



Когда получено уравнение линейной регрессии, можно узнать прогнозируемое значение, которое в качестве аргумента имеет 13-мерный вектор x , с помощью метода `predict`, который выведет прогнозируемое значение.



А теперь попробуем провести регуляризацию. Для начала используем сумму квадратов и коэффициенты формулы линейной регрессии, которые мы только что разобрали.

```
print("Linear Regression")
for f, w in zip(boston.feature_names, lr1.coef_) :
    print("{0:7s}: {1:6.2f}".format(f, w))
print("coef = {0:4.2f}".format(sum(lr1.coef_**2)))
```

Linear Regression

CRIM : -0.11

ZN : 0.05

INDUS : 0.02

CHAS : 2.69

NOX : -17.80

RM : 3.80

AGE : 0.00

DIS : -1.48

```
RAD      :    0.31
TAX      :   -0.01
PTRATIO:  -0.95
B        :    0.01
LSTAT   :  -0.53
coef = 341.86
```



Попробуем также провести ридж-регрессию. Поскольку в данные входят X и y , лучше всего начать с постройки экземпляра, по которому можно провести обучение.

Если есть параметр, который нужно указать, он задается в качестве аргумента экземпляра в формате «имя параметра = значение». Вес α дополнительного параметра регуляризации примем равным 10,0.

```
lr2 = Ridge(alpha=10.0)
lr2.fit(X, y)
print("Ridge")
for f, w in zip(boston.feature_names, lr2.coef_) :
    print("{0:7s}: {1:6.2f}".format(f, w))
print("coef = {0:4.2f}".format(sum(lr2.coef_**2)))
```

```
Ridge
CRIM   :  -0.10
ZN     :   0.05
INDUS  :  -0.04
CHAS   :   1.95
NOX    :  -2.37
RM     :   3.70
AGE    :  -0.01
DIS    :  -1.25
RAD    :   0.28
TAX    :  -0.01
PTRATIO:  -0.80
B      :   0.01
LSTAT  :  -0.56
coef = 25.73
```

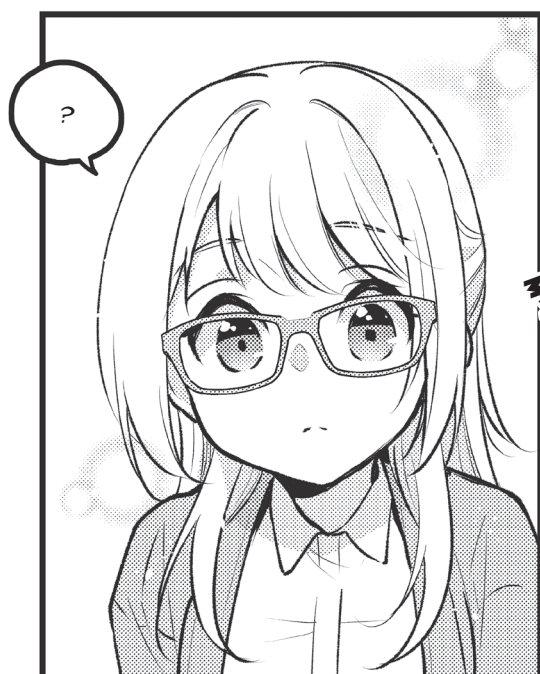
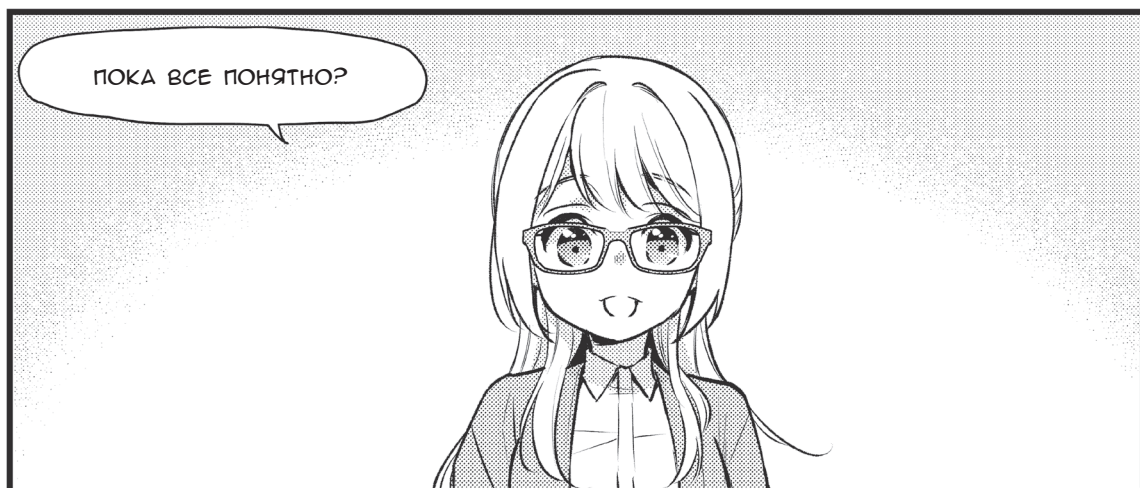


Мы видим, что сумма квадратов коэффициентов абсолютно мала. А теперь применим регрессию «лассо». Вес α дополнительного параметра регуляризации примем равным 2,0 и заметим, что несколько коэффициентов равны 0.

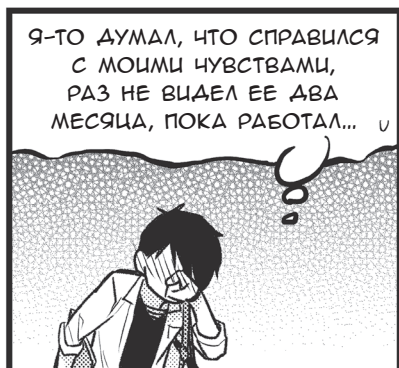
```
lr3 = Lasso(alpha=2.0)
lr3.fit(X, y)
print("Lasso")
for f, w in zip(boston.feature_names, lr3.coef_) :
    print("{0:7s}: {1:6.2f}".format(f, w))
print("coef = {0:4.2f}".format(sum(lr3.coef_**2)))
```

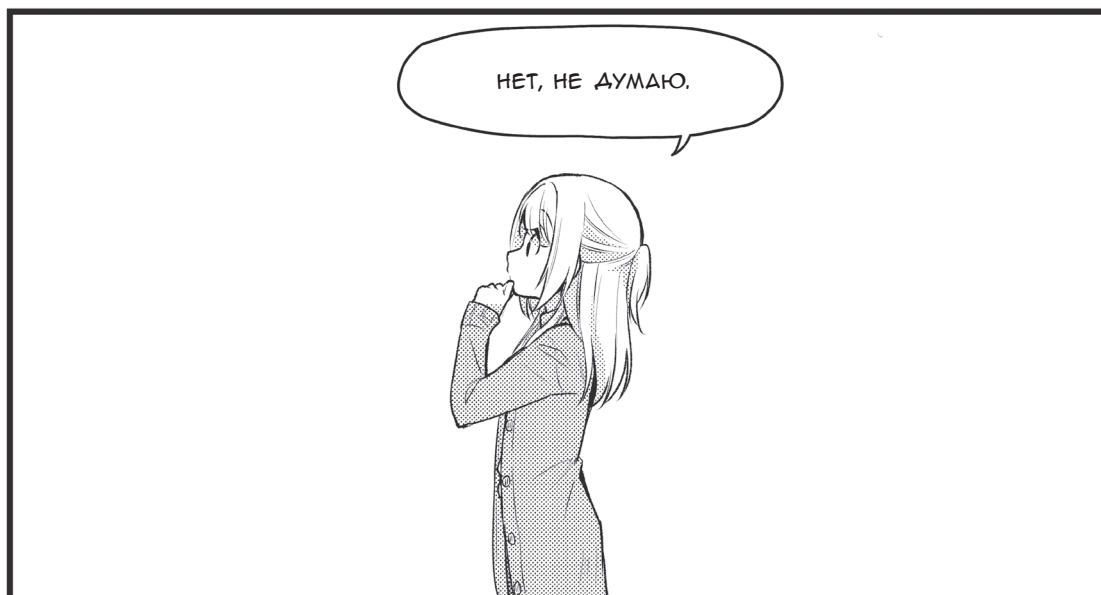
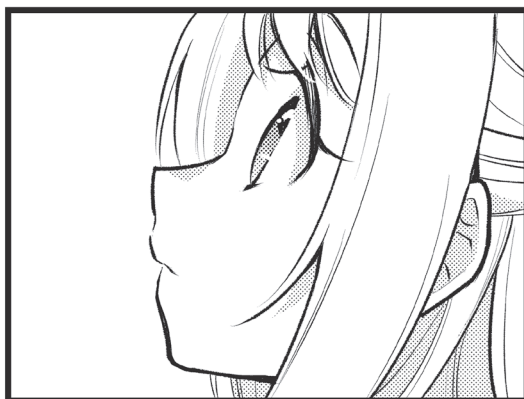
Lasso

```
CRIM   : -0.02
ZN     :  0.04
INDUS  : -0.00
CHAS   :  0.00
NOX    : -0.00
RM     :  0.00
AGE    :  0.04
DIS    : -0.07
RAD    :  0.17
TAX    : -0.01
PTRATIO : -0.56
B      :  0.01
LSTAT  : -0.82
coef = 1.02
```

1.4. РЕГУЛЯРИЗАЦИЯ РЕЗУЛЬТАТА





Математическое повторение (1)



Вот об этом мы говорили с Киёхара-куном. Ай-тян, ты все поняла?

Вы много говорили о векторах и матрицах. Вектор – это последовательность чисел, заключенная в скобки. Есть двухмерный вектор (a, b) , трехмерный (a, b, c) . А что такое d -мерный вектор?



Если d больше четырех, то мы не можем представить себе пространство, это трудно. Но необязательно представлять пространство, можно просто представить очень много чисел, выстроенных в ряд.

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

А что такое вектор-столбец – ряд чисел, выстроенных в ряд?



Нет, не совсем так. Просто когда есть несколько признаков, их обычно представляют вертикально. В машинном обучении часто складывают матрицы и векторы, и матрица обычно слева, а состав матрицы удобно выражать как произведение матриц.

А в школе матрицы не проходят.



Можно сказать, что матрицы – это числа, записанные в виде прямоугольника.



Напомним, что столбцы идут сверху вниз, а строки слева направо.

Aga!



Матрица, которая имеет 2 столбца и 2 строки, называется матрицей 2 на 2 (2×2). Сумма матриц – это матрица, элементы которой равны сумме соответствующих элементов слагаемых матриц, но умножение матриц – очень трудная операция.

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \cdot \begin{pmatrix} 5 & 6 \\ 7 & 8 \end{pmatrix} = \begin{pmatrix} 1 \times 5 + 2 \times 7 & 1 \times 6 + 2 \times 8 \\ 3 \times 5 + 4 \times 7 & 3 \times 6 + 4 \times 8 \end{pmatrix} = \begin{pmatrix} 19 & 22 \\ 43 & 50 \end{pmatrix}$$



Значение n -й строки и m -го столбца матрицы, которая является результатом умножения, определяется так: извлекаются n -я строка из первой матрицы и m -й столбец из второй и перемножаются, начиная с первого числа, а затем складываются.

А если число строк первой матрицы не совпадает с числом столбцов второй, то матрицы просто перемножить нельзя.



Да. Все так. При умножении матрицы ($i \times j$) на матрицу ($j \times k$) результатом будет матрица ($i \times k$), если говорить совсем просто.

$$\begin{pmatrix} \square & \square & \square \\ \square & \square & \square \end{pmatrix} \cdot \begin{pmatrix} \square & \square \\ \square & \square \\ \square & \square \end{pmatrix} = \begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix}$$

i строк, j столбцов

j строк, k столбцов

i строк, k столбцов



Теперь рассмотрим транспонированные и обратные матрицы. Транспонированная матрица X обозначается как X^T и получается при замене столбцов на строки и обратно.

Как мы транспонировали вектор w в формуле 1.1?



Векторы лучше представлять как особый случай. Например, d -мерный столбец-вектор может рассматриваться как матрица из d строк и 1 столбца.

А-а... Вот как. То есть если w – матрица из d строк и 1 столбца, то w^T – матрица из 1 строки и d столбцов. Произведение $(w^T x)$ матрицы из d строк и 1 столбца на транспонированную w^T даст матрицу из одной строки и одного столбца. Но это разве не обычное число?



Именно. Обычное число называется скаляром. $w^T x$ – скаляр, и w_0 – тоже скаляр. Их сумма y – тоже скаляр.

Угу.



А теперь посмотрим на более сложное уравнение 1.8. Матрицу, обратную матрице A , обозначим как A^{-1} . Ай-тян, какое число обратно 5?

Обратное число – это результат деления единицы на него.
Обратное число для 5 – $1/5$.



Да. В основном обратная матрица работает так же. В мире матриц число 1 называется единичной матрицей. Единичная матрица с одинаковым количеством столбцов и строк обозначается I . Внутри нее по диагонали расположены единицы, а остальные величины – нули.

А почему ее величина соответствует единице?



Попробуем перемножить ее с другой матрицей. Ничего не изменилось?

Ничего. Результат произведения такой же.



$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$



Например, матрица A^{-1} , обратная матрице A , в которой два столбца и две строки, будет вычисляться по представленной ниже формуле.

Обычно вычисления матриц, обратной матрице с d столбцами и d строками, доверяют компьютеру.

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$



Если матрицу A перемножить на обратную ей левую матрицу A^{-1} , то получится единичная матрица I .



Есть еще непонятные места?

Вот этот знак – Σ .



Это греческая буква сигма, используется для записи суммы. С использованием сигмы формула $\sum_{i=1}^d w_i x_i$ записывается просто так.

И еще, как дифференцировать функцию, заданную вектором.



Ну... это... Проще говоря, можно представить вектор как обычную переменную и дифференцировать.

Вот и все! А теперь надо искать вес по формуле 1.8.

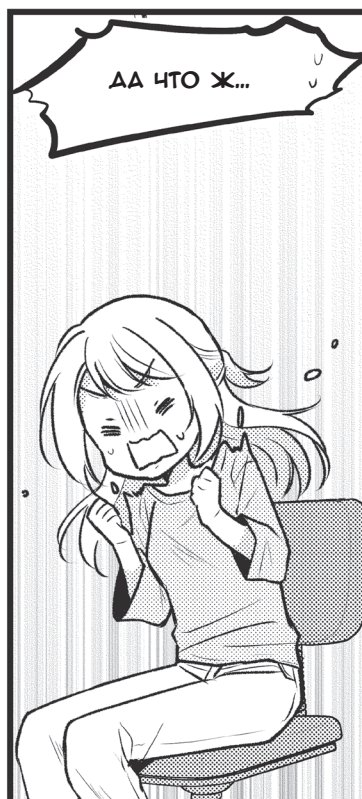


ГЛАВА 2

КАК ДЕЛАТЬ КЛАССИФИКАЦИЮ?

ГДЕ РАСТЕТ
РЕШАЮЩЕЕ ДЕРЕВО?









И ТУТ Я ЗАСТРЯЛ...

ЛАДНО,
ЧТО-НИБУДЬ
ПРИДУМАЮ.

ПОРА ЦАТИ!



СЕГОДНЯ Я РАБОТАЮ
АНЕМ, НАДО БЕЖАТЬ.
ЕЩЕ РАЗ СПАСИБО
ЗА ПОМОЩЬ.

ПОГОДИ-КА,
КИЁХАРА-КУН.

ВЗМАХ

ВЕДЬ МЫ ИЗУЧАЛИ
РЕГРЕССИЮ
НА ПРОШЛОМ ЗАНЯТИИ.

ХОЧЕШЬ, Я ПОМОГУ ТЕБЕ РАЗОБРАТЬСЯ
С КЛАССИФИКАЦИЕЙ?

УДИВЛЕННО

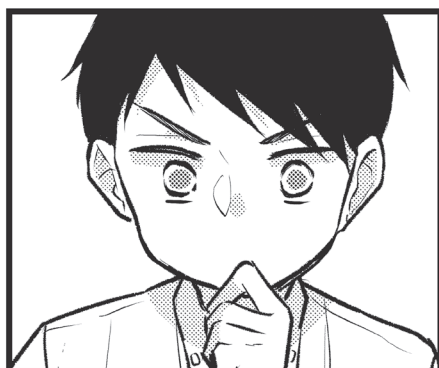
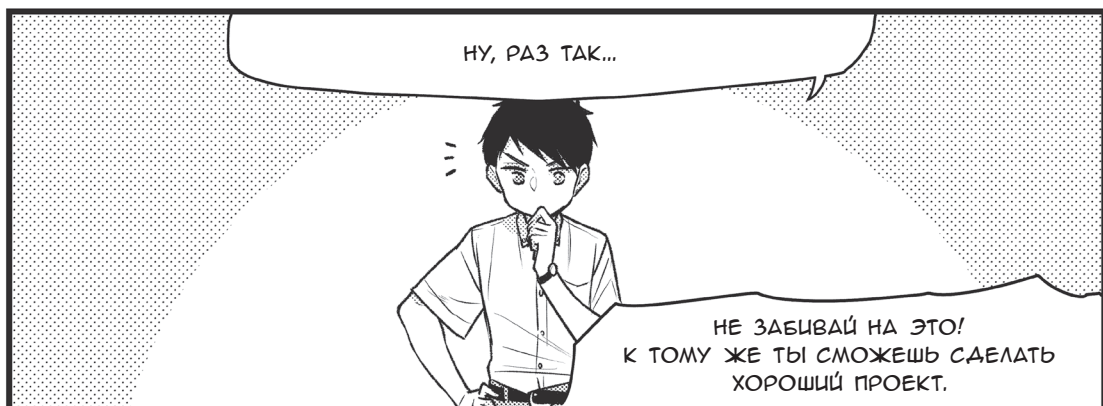
ЧТО?

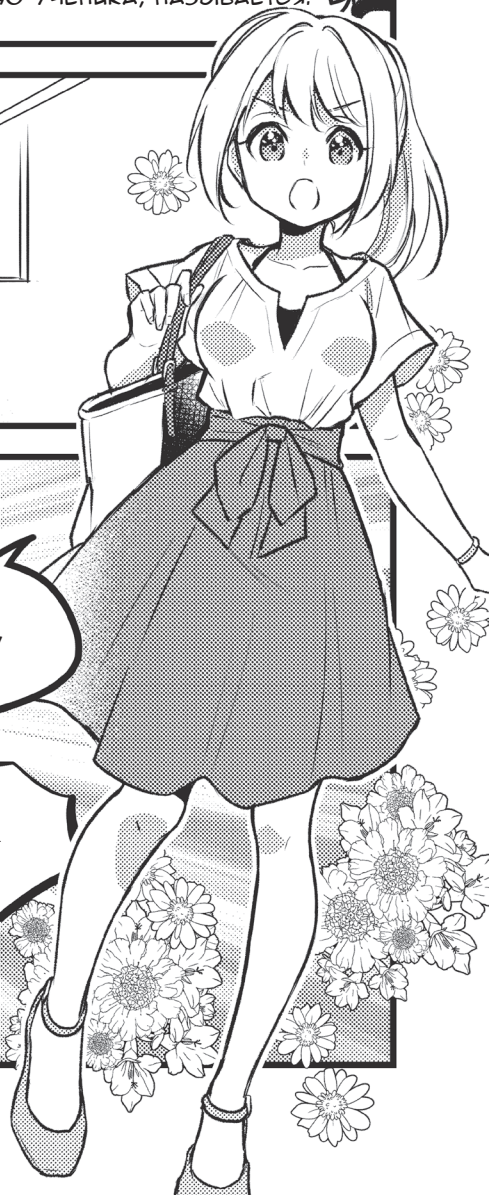
НЕТ-НЕТ, Я ЗАШЕЛ ТЕБЯ
ОТЕБЛАГОДАРЬТЬ!

НЕТ-НЕТ-НЕТ!
ХВАТИТ С МЕНЯ БЕЗОТВЕТНОЙ ЛЮБВИ!
Я БОЛЬШЕ ЭТОГО НЕ ВЫНЕСУ,
ОНА МНЕ ТАК НРАВИТСЯ!

КИЁХАРА-КУН,
ТЫ ЖЕ ХОЧЕШЬ, ЧТОБЫ ВСЕ ЖИТЕЛИ
БЫЛИ ЗДОРОВЫМИ?

ЧТО?







ИЗВИНИ, ЧТО
НАПРЯГАЮ ТЕБЯ
В ВЫХОДНОЙ.

РАЗВЕ?
ЭТО БЫЛА МОЯ ЦЕЛЯ.

ДАВАЙ-КА
ЧТО-НИБУДЬ ЗАКАЖЕМ.



МНЕ, ПОЖАЛУЙСТА,
ПЕРСИКОВЫЙ ПАРФЕ,
ШОКОЛАДНЫЙ ТОРТ,
ЖЕЛЕ АММИЦУ И КОФЕ
ИЗ ДРИНК-БАРА.

А МНЕ ЧИЗКЕЙК...



КАКАЯ ВЫ
СЛАДКОЕЖКА!

НЕ, ЖАРА ЖЕ,
САМОЕ ВРЕМЯ
ДЛЯ ДЕСЕРТОВ.

ЩЕЛК



ИТАК...

Надевает очки

НАЧИНАЕМ УРОК!

2.1. ПРИВОДИМ ДАННЫЕ В ПОРЯДОК

ОПРЕДЕЛИМСЯ, КАКИЕ ДАННЫЕ НУЖНЫ
ДЛЯ ПАЦИЕНТОВ В ГРУППЕ РИСКА
ЗАБОЛЕВАНИЯ ДИАБЕТОМ.

ЕСТЬ ДАННЫЕ ЗА 10 ЛЕТ,
НО ОНИ НЕ СОВСЕМ ПОЛНЫЕ.

Пол	Возраст	ИМТ	Уровень глюкозы	Давление	Диабет
Ж	65	22	180	135	Нет
М	60	28	200	140	Да
М	75	21		120	Нет
Ж	72	25		140	Нет
М	65	26	210		Да
М	80	19	175	112	Нет

Недостающие значения

ЕСЛИ У НАС ЕСТЬ НЕДОСТАЮЩИЕ
ЗНАЧЕНИЯ, ОБУЧЕНИЕ МОЖЕТ ПОЙТИ
ПЛОХО, ПОЭТОМУ НАДО ПРИВЕСТИ
ДАННЫЕ В ПОРЯДОК.

И КАК ЛУЧШЕ ЭТО СДЕЛАТЬ?

ЕСЛИ ДАННЫХ МНОГО,
ТО МОЖНО, КОНЕЧНО, ВЫБРОСИТЬ ВСЕ
С НЕДОСТАЮЩИМИ ЗНАЧЕНИЯМИ,
НО НАДО ЖЕ ИХ ЭФФЕКТИВНО
ИСПОЛЬЗОВАТЬ.

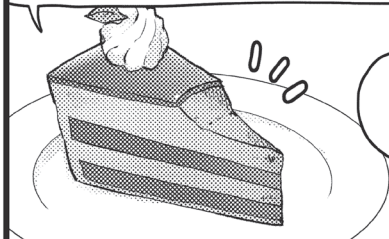


ВЫБРАСЫВАЕТ

Данные

ДА!

ПРОСТОЙ СПОСОБ - ЗАПОЛНИТЬ
НЕДОСТАЮЩИЕ ЗНАЧЕНИЯ
СРЕДНИМИ ДАННЫМИ.



ДА,
ТАК БУДЕТ ЛУЧШЕ.

ОДНАКО ЕСЛИ ДАННЫХ МАЛО,
ИЛИ ЖЕ СРЕДИ НИХ ЕСТЬ ВЫБРОСЫ,
КОТОРЫЕ НЕ ПОПАДАЮТ ПОД ОБЩЕЕ
РАСПРЕДЕЛЕНИЕ, ТО ЗАПОЛНЯТЬ
НЕДОСТАЮЩИЕ ДАННЫЕ СРЕДНИМИ
ЗНАЧЕНИЯМИ НЕ СТОИТ.



А ЧТО
ТОГАА ДЕЛАТЬ?

ЧТОБЫ ВЫБРОСЫ НЕ ОКАЗАЛИ ВЛИЯНИЕ,
ИСПОЛЬЗУЮТ НЕ СРЕДНЕЕ ЗНАЧЕНИЕ,
А МЕДИАННОЕ ИЛИ МОДУ.



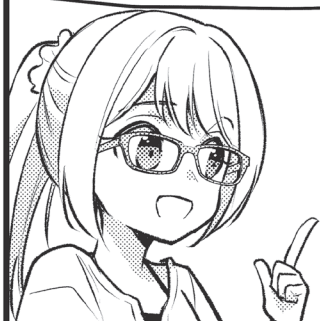
ТО ЕСТЬ НУЖНО ПРИВОДИТЬ
В ПОРЯДОК ДАННЫЕ, СОЗНАВАЯ,
ЧТО ЕСЛИ МЫ ДОБАВИМ
ОПРЕДЕЛЕННЫЕ ЗНАЧЕНИЯ,
ТО РАСПРЕДЕЛЕНИЕ ДАННЫХ
ИЗМЕНИТСЯ?



ДА!

2.2. ОПРЕДЕЛЯЕМ КЛАСС ДАННЫХ

ЗАДАЧА КЛАССИФИКАЦИИ -
ЭТО ЗАДАЧА РАСПРЕДЕЛЕНИЯ
ПО УЖЕ ЗАДАНЫМ КЛАССАМ.



ТИПИЧНЫЕ ЗАДАЧИ КЛАССИФИКАЦИИ
ВКЛЮЧАЮТ В СЕБЯ РАСПОЗНАВАНИЕ РЕЧИ
И ТЕКСТА, РН-КЛАССИФИКАЦИЮ РЕЦЕНЗИЙ,
ОПРЕДЕЛЕНИЕ НАЛИЧИЯ ИЛИ ОТСУТСТВИЯ
БОЛЕЗНИ.



РН-КЛАССИФИКАЦИЯ?

ПОЗИТИВНЫЙ
ИЛИ НЕГАТИВНЫЙ.

ХВАЛЯТ ИЛИ РУГАЮТ
ПРОДУКТ.

ПОЗИТИВНЫЙ

НЕГАТИВНЫЙ

САМАЯ ПРОСТАЯ
ИЗ ЗАДАЧ КЛАССИФИКАЦИИ -
ЭТО БИНАРНАЯ КЛАССИФИКАЦИЯ.

БИНАРНАЯ КЛАССИФИКАЦИЯ -
ЭТО РАЗДЕЛЕНИЕ НА ДВЕ ЧАСТИ

ИМЕННО. БОЛЕН ЧЕЛОВЕК ИЛИ НЕТ,
ЦАЕТ ПИСЬМО В СПАМ ИЛИ НЕТ -
ЭТО ЗАДАЧА НА РАЗДЕЛЕНИЕ
ПО ДВУМ КЛАССАМ.

Японские
сладости
или нет

ТОРТ
(НЕГАТИВНОЕ
ЗНАЧЕНИЕ)

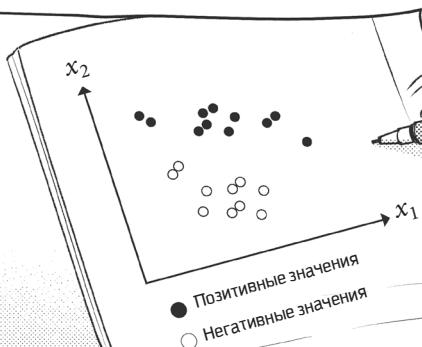
ХЛЕБ АММИНУ
(ПОЗИТИВНОЕ
ЗНАЧЕНИЕ)

Все остальное

Японская сладость

ПОКА ЧТО ДЛЯ ПРОСТОТЫ ПРЕДПОЛОЖИМ,
ЧТО ВВОДЯТСЯ ТОЛЬКО ВЕКТОРЫ
С ЧИСЛЕННЫМИ ЗНАЧЕНИЯМИ.

ЕСЛИ ВВОДЯТСЯ ДВУХМЕРНЫЕ ВЕКТОРЫ,
ТО ДАННЫЕ МОЖНО РАСПОЛОЖИТЬ
НА ПЛОСКОСТИ ВОТ ТАКИМ ОБРАЗОМ:



ПИШЕТ

В ЗАВИСИМОСТИ ОТ КЛАССА
ОБОЗНАЧИМ ИХ ЧЕРНЫМИ
ИЛИ БЕЛЫМИ ТОЧКАМИ.

НЕ ВИДИШЬ?
МОЖЕТ, ТЕБЕ ОТСЮДА
НЕ ПОНЯТНО.

ЧТО? НЕТ!
Я ВСЕ ВИЖУ,
И МНЕ ВСЕ
ПОНЯТНО!

МОЖЕТ МНЕ ПРОЩЕ ВСТАТЬ
И ТАК ОБЪЯСНЯТЬ...

НЕТ-НЕТ-НЕТ!
И ТАК ПОНИМАЮ!

ИЛИ Я СЯДУ РЯДОМ...

А-А-А-А-А. Я ПЕРЕСЯДУ!



2.3. ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ

ДРУГИМИ СЛОВАМИ, ЛОГИСТИЧЕСКУЮ РЕГРЕССИЮ МОЖНО ПРЕДСТАВИТЬ КАК РАСШИРЕНИЕ ЗАДАЧИ РЕГРЕССИИ. ПОГОВОРИМ О ЛОГИСТИЧЕСКОЙ РЕГРЕССИИ!

Логистическая регрессия – метод нахождения функции на основании ввода взвешенной суммы данных; если элемент данных принадлежит к положительному классу, то выход формулы регрессии близок к 1, а если к отрицательному, то он будет близок к 0.



В СЛУЧАЕ ЗАДАЧИ БИНАРНОЙ КЛАССИФИКАЦИИ, КАК СЕЙЧАС, ПРИЗНАКОВОМУ ОПИСАНИЮ ОБЪЕКТА $x = (x_1, \dots, x_d)^T$ СООТВЕТСТВУЕТ ВЗВЕШЕННАЯ СУММА ВСЕХ ПРИЗНАКОВ $w_1x_1 + \dots + w_dx_d$.

НЕОБХОДИМО СКОРРЕКТИРОВАТЬ ВЕС ТАК, ЧТОБЫ ДЛЯ ПОЛОЖИТЕЛЬНЫХ ПРИМЕРОВ ФУНКЦИЯ ПРИНИМАЛА ЗНАЧЕНИЯ, БЛИЗКИЕ К 1, А ДЛЯ ОТРИЦАТЕЛЬНЫХ – БЛИЗКИЕ К 0.

ДРУГИМИ СЛОВАМИ, ЕСЛИ НАСТРОИТЬ ВЕС ТАК, ЧТО РЕЗУЛЬТАТ ФОРМУЛЫ РЕГРЕССИИ БУДЕТ РАВЕН 1 ДЛЯ ПОЛОЖИТЕЛЬНЫХ ПРИМЕРОВ И 0 ДЛЯ ОТРИЦАТЕЛЬНЫХ, ТО НЕВОЗМОЖНО РЕШИТЬ ЭТО УРАВНЕНИЕ, ЕСЛИ $x = 0$; ПОЭТОМУ В КАЧЕСТВЕ ПАРАМЕТРА ДОБАВЛЯЕТСЯ ПОСТОЯННАЯ w_0 .

$$\hat{g}(x) = w_0 + w_1x_1 + \dots + w_dx_d = w_0 + w^Tx$$

скрип



- черный кружок \times вес = близко к 1
- белый кружок \times вес = близко к 0

Распределение весов

И ПОЛУЧАЕТСЯ ТАКАЯ ФОРМУЛА:

$$\hat{g}(x) = w_0 + w_1x_1 + \dots + w_dx_d = w_0 + \mathbf{w}^T \mathbf{x}$$



ЗАЕЩЬ $\mathbf{w}^T \mathbf{x}$ -
ВНУТРЕННЕЕ ПРОИЗВЕДЕНИЕ ВЕКТОРА \mathbf{w}
НА ВЕКТОР \mathbf{x} . ЭТО МОЖНО ОПРЕДЕЛИТЬ
КАК СУММУ ПРОИЗВЕДЕНИЙ
СООТВЕТСТВУЮЩИХ ЭЛЕМЕНТОВ.

ПОСКОЛЬКУ ДЛЯ ГИПЕРПЛОСКОСТИ $\hat{g}(x) = 0$,
 d -МЕРНУЮ ГИПЕРПЛОСКОСТЬ МОЖНО ВЫРАЗИТЬ
УРАВНЕНИЕМ. ЕСЛИ ПЛОСКОСТЬ ВЕДЕТ СЕБЯ,
КАК УКАЗАНО ВЫШЕ, ТО ЧТО БУДЕТ ПОСЛЕ ТОГО,
КАК ВСЕ ОПРЕДЕЛЕНО?

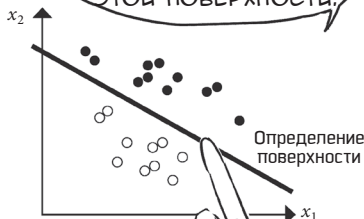
ЭМ...
НУ, ЕСЛИ ПЛОСКОСТЬ
РАВНА 0...

В ПОЛОЖИТЕЛЬНОЙ ПЛОСКОСТИ
ЗНАЧЕНИЕ БУДЕТ ПОЛОЖИТЕЛЬНОЕ,
В ОТРИЦАТЕЛЬНОЙ - ОТРИЦАТЕЛЬНОЕ?

ИМЕННО!

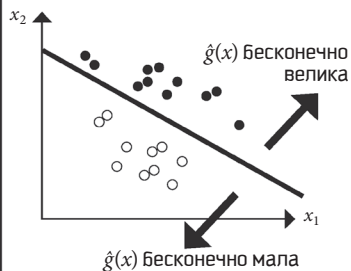
НЕ ЗАБЫВАЙ,
ЧТО НАДО УЗНАТЬ,
К КАКОМУ КЛАССУ
ОТНОСЯТСЯ ТОЧКИ
НА ПЛОСКОСТИ!

ПОВЕРХНОСТЬ, КОТОРАЯ
ДЕЛИТ ПРОСТРАНСТВО
ПРИЗНАКОВ НА КЛАССЫ,
НАЗЫВАЕТСЯ **РАЗДЕЛЯЮЩЕЙ**.
ТОЧНОСТЬ, С КОТОРОЙ
МОЖНО ОТНЕСТИ ТОЧКУ
К ТОМУ ИЛИ ИНОМУ
КЛАССУ, ОПРЕДЕЛЯЕТСЯ
ЕЕ РАССТОЯНИЕМ ОТ
ЭТОЙ ПОВЕРХНОСТИ.



А МОЖЕТ ТАК БЫТЬ, ЧТО
 $\hat{g}(x)$ В ЗАВИСИМОСТИ
ОТ ЗНАЧЕНИЯ x БУДЕТ
ТО УВЕЛИЧИВАТЬСЯ,
ТО УМЕНЬШАТЬСЯ?

КОНЕЧНО, ОТ МИНУС
БЕСКОНЕЧНОСТИ ДО
ПЛУС БЕСКОНЕЧНОСТИ.



ДА!

ПОСКОЛЬКУ ТРУДНО ОЦЕНИТЬ ВЕРОЯТНОСТЬ ПРИНАДЛЕЖНОСТИ x К ПОЛОЖИТЕЛЬНОМУ КЛАССУ, НАДО СДЕЛАТЬ ТАК, ЧТОБЫ РАЗБРОС РЕЗУЛЬТАТОВ НА ВЫХОДЕ ФУНКЦИИ $\hat{g}(x)$ КОЛЕБАЛСЯ ОТ 0 ДО 1; ЕСЛИ x ПРИНАДЛЕЖИТ К ПОЛОЖИТЕЛЬНОМУ КЛАССУ, ТО РЕЗУЛЬТАТ ДОЛЖЕН БЫТЬ БЛИЗОК К 1, А ЕСЛИ К ОТРИЦАТЕЛЬНОМУ, ТО БЛИЗОК К 0.

И КАК ЭТО МОЖНО СДЕЛАТЬ?

ЕСЛИ МЫ ВОЗЬМЕМ ПРЕОБРАЗОВАННУЮ ФУНКЦИЮ $p(+|x) = \frac{1}{1 + e^{-(w_0 + w^T x)}}$, ТО ПО ФОРМУЛЕ НИЖЕ МОЖНО ПОЛУЧИТЬ ВЕРОЯТНОСТЬ ПРИНАДЛЕЖНОСТИ x К ПОЛОЖИТЕЛЬНОМУ КЛАССУ.

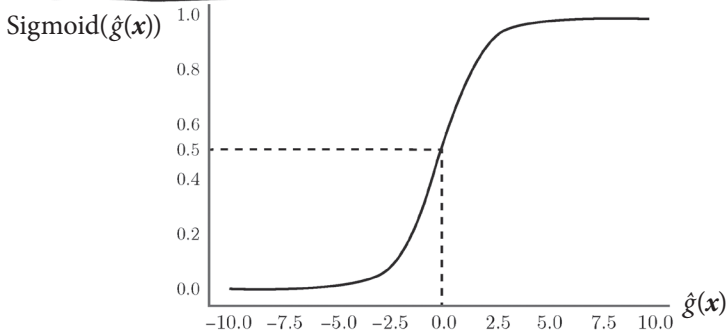
$$p(+|x) = \frac{1}{1 + e^{-(w_0 + w^T x)}}$$

А КАКОВА ВЕРОЯТНОСТЬ ТОГО, ЧТО x ПРИНАДЛЕЖИТ К ОТРИЦАТЕЛЬНОМУ КЛАССУ?

ЕСЛИ ВЫЧЕСТЬ ВЕРОЯТНОСТЬ ПОЛОЖИТЕЛЬНОГО КЛАССА ИЗ 1, ТО $p(-|x) = 1 - p(+|x)$?

ИМЕННО!

ВОТ ГРАФИК ЭТОЙ ФУНКЦИИ.



ЕСЛИ ПОСМОТРЕТЬ НА ГРАФИК, ЯСНО, ЧТО КАКОЕ БЫ ЗНАЧЕНИЕ НИ ПРИНЯЛА ФУНКЦИЯ $\hat{g}(x) = w_0 + w^T x$, РЕЗУЛЬТАТ БУДЕТ НАХОДИТЬСЯ В ДИАПАЗОНЕ ОТ 0 ДО 1. ЭТО - СИГМОИДАЛЬНАЯ ФУНКЦИЯ.

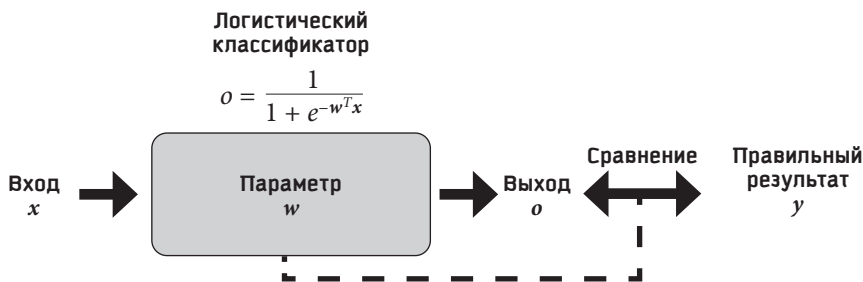
ЕСЛИ $\hat{g}(x) = 0$, ТО ФУНКЦИЯ ПРИМЕТ ЗНАЧЕНИЕ 0,5.

ЭТО ФУНКЦИЯ, ПОКАЗЫВАЮЩАЯ ЗНАЧЕНИЕ ВЕРОЯТНОСТИ.



Далее поговорим о том, как провести обучение логистической классификации. Логистический классификатор можно рассматривать как вероятностную модель с весовым параметром \mathbf{w} .

Здесь и далее для простоты объяснения \mathbf{w} будет включать в себя w_0 .



Пусть в данных для обучения D этой модели при входе \mathbf{x}_i , а на выходе o_i . Желаемый выход обозначим положительным исходом y_i . Предположим, что у нас задача бинарной классификации, $y_i = 1$, если данные принадлежат к положительному классу, а если к отрицательному, то $y_i = 0$.

Чтобы правильно провести обучение созданной модели и оценить значение, необходимо уравнение правдоподобия, которое приведено ниже. Π обозначает произведение.

$$P(D | \mathbf{w}) = \prod_{\mathbf{x}_i \in D} o_i^{y_i} (1 - o_i)^{(1-y_i)}.$$

$o_i^{y_i} (1 - o_i)^{(1-y_i)}$ принимает значение o_i , если элемент данных принадлежит к положительному классу, и тогда $(y_i = 1)$, и значение $(1 - o_i)$, если к отрицательному $(y_i = 0)$. Иными словами, если настроить веса \mathbf{w} так, что при положительном результате выходное значение o_i будет близко к 1, а при отрицательном результате o_i будет близко к 0, то во всех данных величина произведения $P(D | \mathbf{w})$ будет увеличиваться в зависимости от данных.



При расчете значения максимального правдоподобия для простоты расчета используется логарифмическая функция правдоподобия.

$$L(D) = \log P(D | \mathbf{w}) = \sum_{\mathbf{x}_i \in D} \{y_i \log o_i + (1 - y_i) \log (1 - o_i)\}.$$



Чтобы представить ясней задачу оптимизации, в дальнейшем будем рассматривать задачу минимизации функции ошибки $E(\mathbf{w})$, которая может быть определена как логарифмическая функция со знаком минус.

$$E(\mathbf{w}) = -\log P(D|\mathbf{w}).$$



Продифференцировав, найдем предельное значение \mathbf{w} . Поскольку модель – логистический классификатор, выход o_i будет представлять собой сигмоидную функцию.

$$S(z) = \frac{1}{1 + e^{-z}}.$$

Производная сигмоидной функции будет выглядеть так:

$$S'(z) = S(z) \cdot (1 - S(z)).$$



Поскольку на выходе модели есть функция веса \mathbf{w} , при ее изменении меняется величина ошибки. Решение таких задач можно найти методом градиентного спуска. Метод градиентного спуска – это метод сходимости к оптимальному решению путем многократного постепенного уменьшения параметров в направлении градиента минимизируемой функции.

В этом случае мы немного меняем параметр \mathbf{w} , чтобы найти направление наискорейшего спуска функции ошибки $E(\mathbf{w})$. Это «немного» обозначается коэффициентом обучаемости η . Используя метод градиентного спуска, можно вывести новую формулу веса.

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}}.$$



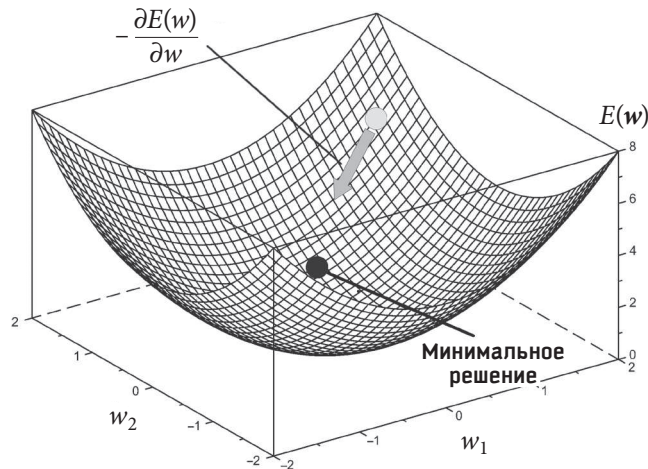
Затем вычислим направление градиента функции ошибки $E(\mathbf{w})$ по формуле ниже:

$$\frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = - \sum_{\mathbf{x}_i \in D} \left(\frac{y_i}{o_i} - \frac{1 - y_i}{1 - o_i} \right) o_i (1 - o_i) \mathbf{x}_i = - \sum_{\mathbf{x}_i \in D} (y_i - o_i) \mathbf{x}_i.$$

Следовательно, новая формула веса будет выглядеть так:

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \sum_{\mathbf{x}_i \in D} (y_i - o_i) \mathbf{x}_i.$$

Когда новое значение веса будет ниже заранее определенного значения, метод градиентного спуска закончен.



Метод, при котором градиент вычисляют на основании всех данных для обучения D , называется **пакетным** (batch method). Если из D выбирают данные определенной величины и вычисляют отдельный градиент по ним, то это называется **мини-пакетный метод** (mini batch), а метод, когда из D случайно выбирают элемент данных и вычисляют градиент для него, называется **методом стохастического градиентного спуска**.

2.4. КЛАССИФИКАЦИЯ ПО РЕШАЮЩЕМУ ДЕРЕВУ

ДАЛЕЕ – КЛАССИФИКАЦИЯ ПУТЕМ РЕШАЮЩЕГО ДЕРЕВА... ВОЗМОЖНО, ЭТОТ СПОСОБ ПОДХОДИТ ДЛЯ ОПРЕДЕЛЕНИЯ ТЕХ, КТО РИСКУЕТ ЗАБОЛЕТЬ ДИАБЕТОМ.

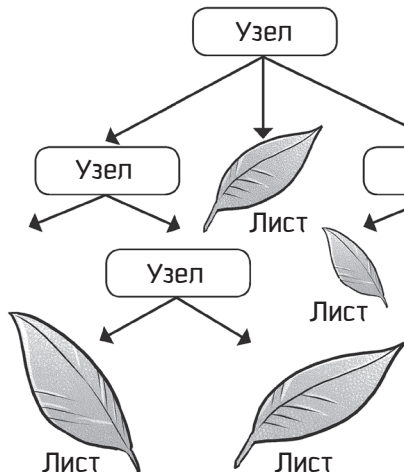


ДЕРЕВО?
КОТОРОЕ РАСТЕТ?



ДА, ТОЛЬКО
КОРНИ У НЕГО БУДУТ СВЕРХУ,
А ЛИСТЬЯ – СНИЗУ.

Структура решающего дерева



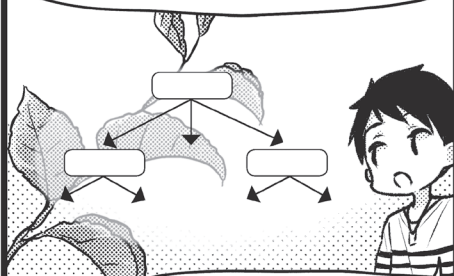
Решающее дерево состоит из узлов (веток), в которых данные классифицируются, и листьев, в которых выводится результат классификации.



МОЖНО ПРЕОБРАЗОВАТЬ ДЕРЕВО В ЭКВИВАЛЕНТНОЕ ЛОГИЧЕСКОЕ ВЫРАЖЕНИЕ, КОМБИНИРУЯ ЗНАЧЕНИЯ ВЕТВЕЙ УЗЛОВ, КОТОРЫЕ ВЕДУТ К ПОЛОЖИТЕЛЬНОМУ РЕЗУЛЬТАТУ С УСЛОВИЕМ И, А ВСЕ ОСТАЛЬНОЕ – С УСЛОВИЕМ ИЛИ.

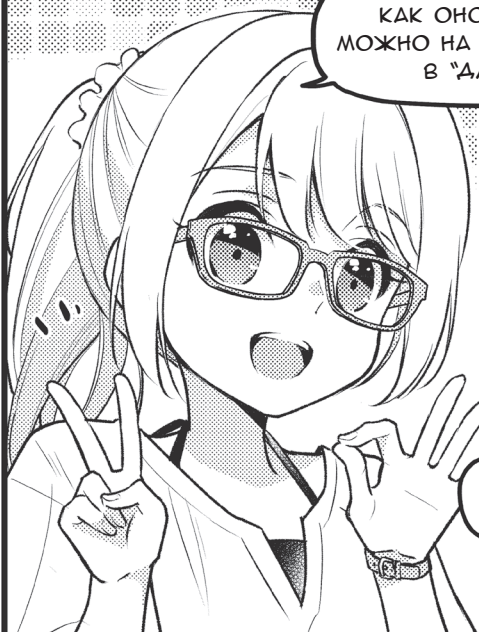
ОГО!

ОДНАКО ПОСКОЛЬКУ ПО СТРУКТУРЕ ДЕРЕВА ЛЕГЧЕ НАГЛЯДНО ПРЕДСТАВИТЬ РЕЗУЛЬТАТ ОБУЧЕНИЯ, ТО ЭТОТ МЕТОД ПРЕДПОЧТИТЕЛЕН.



РЕШАЮЩЕЕ ДЕРЕВО ПРОЩЕ ПОНЯТЬ, ЧЕМ ЛОГИЧЕСКИЕ СХЕМЫ.

ОБЪЯСНИТЬ, КАК ОНО РАБОТАЕТ, МОЖНО НА ПРИМЕРЕ ИГРЫ В "ДАНЕТКИ".



ДАНЕТКИ?

ЭТО ИГРА, В КОТОРОЙ НАДО ОТГАДАТЬ ТО,
ЧТО ЗАГАДАЛ ВЕДУЩИЙ, ЗАДАВ НЕ БОЛЕЕ
20 ВОПРОСОВ, НА КОТОРЫЕ МОЖНО ОТВЕТИТЬ
ТОЛЬКО "ДА" ИЛИ "НЕТ".

МОЖЕТ БЫТЬ,
ИГРАЛ В ДЕТСТВЕ В НЕЕ?



Животное?

Летает?

Лает?

СЕКРЕТ УСПЕХА В ТОМ,
ЧТОБЫ СНАЧАЛА НЕ ЗАДАВАТЬ
СЛИШКОМ КОНКРЕТНЫЕ ВОПРОСЫ.

ЕСЛИ ПЕРВЫЕ ВОПРОСЫ БУДУТ
КОНКРЕТНЫМИ, МОЖНО СУЗИТЬ
СПИСОК КАНДИДАТОВ,
НО В БОЛЬШИНСТВЕ СЛУЧАЕВ
НИКАКОЙ КОНКРЕТНОЙ
ИНФОРМАЦИИ УЗНАТЬ
НЕ ПОЛУЧИТСЯ.

Оно ядовитое?



СНАЧАЛА НАДО
ЗАДАВАТЬ ОБЩИЕ ВОПРОСЫ,
ЧТОБЫ ПОНЯТЬ, ЧТО ЭТО?

ИМЕННО!

РЕШАЮЩЕЕ ДЕРЕВО СТРОИТСЯ
ИМЕННО НА ОСНОВЕ ЭТОГО СЕКРЕТА -
ВСЕ ВОПРОСЫ, НА КОТОРЫЕ ДАЕТСЯ ОТВЕТ,
ЯВЛЯЮТСЯ УЗЛАМИ ДАННОГО ДЕРЕВА.

В ОСНОВНОМ ДЛЯ ПОСТРОЙКИ
ТАКОГО ДЕРЕВА ИСПОЛЬЗУЕТСЯ
АЛГОРИТМ ID3.

	Погода	Температура	Влажность	Ветер	Играем?
1	Ясно	Высокая	Высокая	Нет	Нет
2	Ясно	Высокая	Высокая	Да	Нет
3	Облачно	Высокая	Высокая	Нет	Да
4	Дождь	Средняя	Высокая	Нет	Да
5	Дождь	Низкая	Стандартная	Нет	Да
6	Дождь	Низкая	Стандартная	Да	Нет
7	Облачно	Низкая	Стандартная	Да	Да
8	Ясно	Средняя	Высокая	Нет	Нет
9	Ясно	Низкая	Стандартная	Нет	Да
10	Дождь	Средняя	Стандартная	Нет	Да
11	Ясно	Средняя	Стандартная	Да	Да
12	Облачно	Средняя	Высокая	Да	Да
13	Облачно	Высокая	Стандартная	Нет	Да
14	Дождь	Средняя	Высокая	Да	Нет

ЭТО ДАННЫЕ ПО ИГРЕ В ГОЛЬФ,
НА ПРИМЕРЕ КОТОРЫХ
МЫ РАССМОТРИМ ID3-АЛГОРИТМ.



ОНИ ПОКАЗЫВАЮТ,
ИГРАЛИ ЛЮДИ В ГОЛЬФ
НА ПРОТЯЖЕНИИ ДВУХ НЕДЕЛЬ
ИЛИ НЕТ?



ЗАДАЧА СОСТОИТ В ТОМ,
ЧТОБЫ ОПРЕДЕЛИТЬ, МОЖНО ИГРАТЬ
В ГОЛЬФ ИЛИ НЕТ ПРИ СООТВЕТСТВУЮЩИХ
ПОГОДНЫХ УСЛОВИЯХ. НАДО ВЫДЕЛИТЬ
ОДНУ КОНКРЕТНУЮ ВЕЛИЧИНУ И УЗНАТЬ
ОТВЕТ ПРИ ПОМОЩИ НАИМЕНЬШЕГО
КОЛИЧЕСТВА ВОПРОСОВ.



С КАКОГО ВОПРОСА
ТОГДА ЛУЧШЕ НАЧАТЬ?



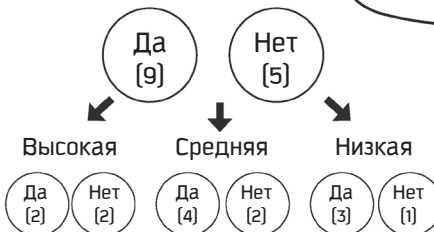
НЕ МОГУ СРАЗУ СКАЗАТЬ...

ДОПУСТИМ, МЫ ЗАДАДИМ ВОПРОС ПРО ВЛАЖНОСТЬ.

В ЗАВИСИМОСТИ ОТ ОТВЕТА ДАННЫЕ БУДУТ ДЕЛИТЬСЯ ТАК, КАК ПОКАЗАНО НИЖЕ:

ТАК... КАЖЕТСЯ, ТУТ НЕЛЬЗЯ НИЧЕГО ВЫДЕЛИТЬ.

Данные D



Если мы сначала спросим про влажность, то какой бы ответ ни был выбран, условия для задания второго вопроса не будут отличаться от предыдущих.

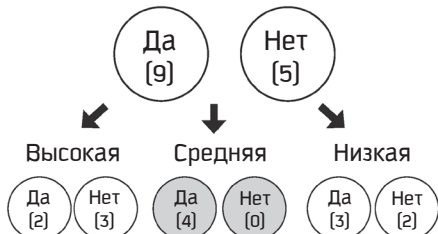
АГА. ПОТОМУ, ЧТО УСЛОВИЯ СОВСЕМ НЕ ОТЛИЧАЮТСЯ.

А ЕСЛИ МЫ ВЫБЕРЕМ ПОГОДУ?

В СЛУЧАЕ ОБЛАЧНОЙ ПОГОДЫ ВСЕ РЕЗУЛЬТАТЫ ПОЛОЖИТЕЛЬНЫЕ.

ТУТ МОЖНО КОЕ-ЧТО ВЫДЕЛИТЬ.

Данные D



Если мы сначала спросим про погоду, то в случае облачной погоды (основываясь на данных) все результаты будут положительными.

ДА, ОДНАКО НАДО ПОНЯТЬ, КАКОЕ ИМЕННО УСЛОВИЕ ВЛИЯЕТ НА РЕЗУЛЬТАТ.

ЕСЛИ ВЫБЕРЕМ ТОЛЬКО ЯСНУЮ ПОГОДУ, ДАННЫЕ БУДУТ ТАКИМИ:

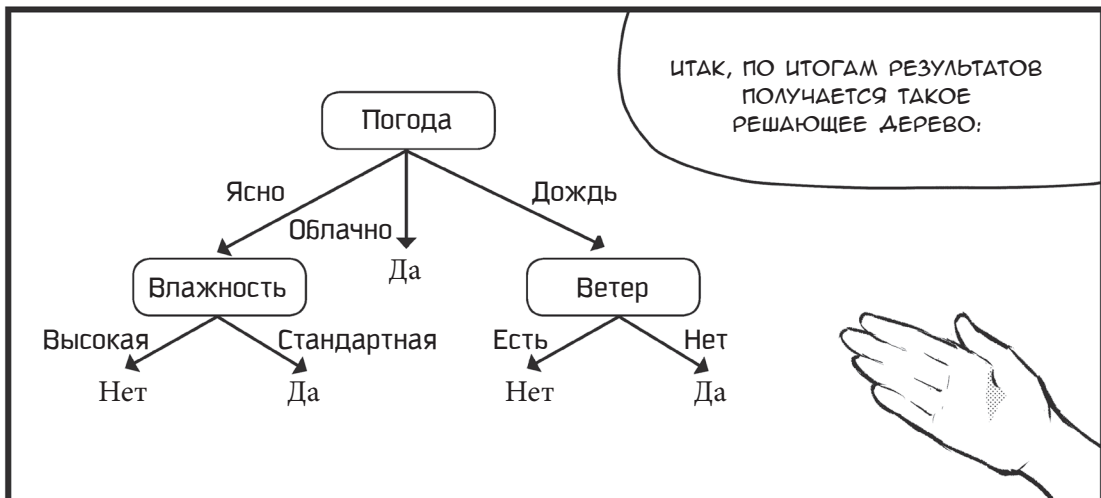
А ЕСЛИ МЫ ВЫБЕРЕМ ДОЖДЬ, ПОСМОТРИМ, ЧТО ПОЛУЧИТСЯ...

	Погода	Температура	Влажность	Ветер	Играем?
1	Ясно	Высокая	Высокая	Нет	Нет
2	Ясно	Высокая	Высокая	Да	Нет
8	Ясно	Средняя	Высокая	Нет	Нет
9	Ясно	Низкая	Стандартная	Нет	Да
11	Ясно	Средняя	Стандартная	Да	Да

	Погода	Температура	Влажность	Ветер	Играем?
4	Дождь	Средняя	Высокая	Нет	Да
5	Дождь	Низкая	Стандартная	Нет	Да
6	Дождь	Низкая	Стандартная	Да	Нет
10	Дождь	Средняя	Стандартная	Нет	Да
14	Дождь	Средняя	Высокая	Да	Нет

А ЕСЛИ ТЕПЕРЬ СПРОСИТЬ ПРО ВЛАЖНОСТЬ, ТО ВСЕ ОТВЕТЫ БУДУТ "ДА", ЕСЛИ ВЛАЖНОСТЬ СТАНДАРТНАЯ.

ЗНАЧИТ, ЕСЛИ ДОЖДЬ, НАДО СМОТРЕТЬ, ЕСТЬ ЛИ ВЕТЕР.



ИТАК...
ПОДУМАЙ, КАК ЛУЧШЕ ВСЕГО
НАЙТИ САМЫЙ ЭФФЕКТИВНЫЙ
ВОПРОС, КАК ПОГОДА,
В ЭТОМ СЛУЧАЕ?

☀ Погода?

°C Температура?

% Влажность?

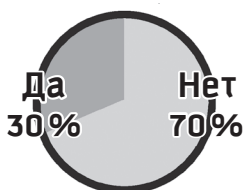
↘ Ветер?



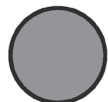
НАЧНЕМ
С НЕПРЕДСКАЗУЕМОСТИ.

Непредсказуемость означает, что трудно определить заранее, какие элементы из набора данных дадут на выходе ответы «да» или «нет».

Набор данных



Данные {1}



Это да?
Или нет?

СЛОЖНОСТЬ?

САМЫЙ ТРУДНЫЙ СЛУЧАЙ –
ЭТО КОГДА ПОЛОВИНА ОТВЕТОВ
“ДА”, А ПОЛОВИНА – “НЕТ”.

А КАКОЙ
ТОГДА САМЫЙ
ЛЕГКИЙ СЛУЧАЙ?



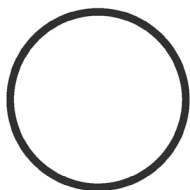
КОГДА ВСЕ ОТВЕТЫ “ДА”
ИЛИ ВСЕ ОТВЕТЫ “НЕТ”.

ДА... ИНЫМИ СЛОВАМИ,
ВСЕ ДАННЫЕ МОЖНО
ОТНЕСТИ К ОДНОМУ
И ТОМУ ЖЕ КЛАССУ.

ЗАТЕМ ЦАЕТ
ИНФОРМАЦИОННАЯ ЭНТРОПИЯ.

Информационная энтропия – это количество информации, определяемое вероятностью получения определенного результата (ДА или НЕТ) из набора данных.

Набор данных



Данные {1}

Да

По крайней мере,
не все – НЕТ!

ВЕРОЯТНОСТЬЮ
ПОЛУЧЕНИЯ
РЕЗУЛЬТАТА?

НАПРИМЕР, ЕСЛИ МЫ ПОЛУЧИЛИ ОТВЕТ "ДА"
ИЗ НАБОРА ДАННЫХ, ГДЕ ВСЕ ДАННЫЕ -
"ДА", МЫ НЕ ПОЛУЧИЛИ НИКАКОЙ
ИНФОРМАЦИИ.

Да!

Да!

Да!



НАВЕРНОЕ, ТАК.

А ЕСЛИ ИЗ 14 ЭЛЕМЕНТОВ ДАННЫХ
13 БУДУТ "ДА", А ОДИН - "НЕТ"?

Да!

Да!

Нет!



У НАС БУДЕТ
ИНФОРМАЦИЯ О ТОМ,
ЧТО СЛУЧИЛОСЬ ЧТО-ТО
НЕОБЫЧНОЕ.

МОЖНО СКАЗАТЬ, ЧТО ИНФОРМАЦИОННАЯ ЭНТРОПИЯ
ВЕЛИКА, ЕСЛИ ВЕРОЯТНОСТЬ КАКОГО-НИБУДЬ СОБЫТИЯ
НИЗКАЯ, И НИЗКА, ЕСЛИ ВЕРОЯТНОСТЬ ЕГО ПОЯВЛЕНИЯ
ВЫСОКАЯ.

ЕЕ МОЖНО СЧИТАТЬ
ВЕЛИЧИНОЙ, ОБРАТНОЙ
ВЕРОЯТНОСТИ.

Если вероятность высокая ↗, информационная энтропия низкая ↘.

Если вероятность низкая ↘, информационная энтропия высокая ↗.

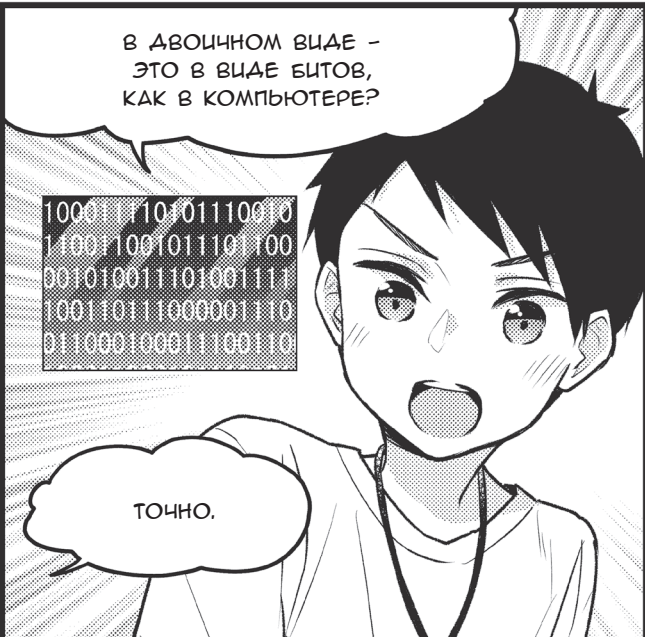
АГА! ТО ЕСТЬ ЕСЛИ ПОЯВЛЯЕТСЯ КАКОЕ-ТО РЕДКОЕ СОБЫТИЕ,
ТО ВЕРОЯТНОСТЬ НИЗКАЯ, А ИНФОРМАЦИОННАЯ ЭНТРОПИЯ СТАНОВИТСЯ БОЛЬШЕ!

ДА... ЕСЛИ ВЫЧИСЛИТЬ ЛОГАРИФМ
ПО ОСНОВАНИЮ 2 ДЛЯ ВЕРОЯТНОСТИ,
Т. Е. ОБРАТНОЙ ВЕЛИЧИНЫ ЭНТРОПИИ,
МОЖНО УЗНАТЬ КОЛИЧЕСТВО ЦИФР,
НЕОБХОДИМОЕ ДЛЯ ПРЕДСТАВЛЕНИЯ
ЭТОЙ ИНФОРМАЦИИ В ДВОИЧНОМ
ВИДЕ.



В ДВОИЧНОМ ВИДЕ -
ЭТО В ВИДЕ БИТОВ,
КАК В КОМПЬЮТЕРЕ?

```
10001110101110010
110011001011101100
001010011101001111
100110111000001110
011000100011100110
```



ТОЧНО.

А ТЕПЕРЬ РАССМОТРИМ ФОРМУЛУ
НЕПРЕДСКАЗУЕМОСТИ ДАННЫХ.

ПОЖАЛУЙСТА.

ЕЕ МОЖНО НАЙТИ,
ЕСЛИ МЫ СЛОЖИМ ПРОИЗВЕДЕНИЯ
ИНФОРМАЦИОННОЙ ЭНТРОПИИ
КАЖДОГО КЛАССА И СООТНОШЕНИЕ
КАЖДОГО КЛАССА С ОБЩИМИ
ДААННЫМИ, И ЗАПИСАТЬ СЛЕДУЮЩЕЙ
ФОРМУЛОЙ:

$$E(D) = -P_{\text{Да}} \log_2 P_{\text{Да}} - P_{\text{Нет}} \log_2 P_{\text{Нет}}$$



КОГДА МЫ ЗАДАЕМ ВОПРОС,
ТО ДАННЫЕ ВЕДЬ МОЖНО РАЗДЕЛИТЬ
В ЗАВИСИМОСТИ ОТ ОТВЕТА? ЗАТЕМ
МОЖНО НАЙТИ НЕПРЕДСКАЗУЕМОСТЬ
ПО ФОРМУЛЕ, КОТОРАЯ
ПРИВЕДЕНА ВЫШЕ.

МОЖНО ОПРЕДЕЛИТЬ
ЗНАЧЕНИЕ СНИЖЕНИЯ НЕПРЕДСКАЗУЕМОСТИ
КАК ЗНАЧЕНИЕ ИНФОРМАЦИОННОГО ВЫИГРЫША,
И ТАМ, ГДЕ ЭТО ЗНАЧЕНИЕ ВЫШЕ ВСЕГО,
НАХОДИТСЯ ВОПРОС, КОТОРЫЙ ПОМОЖЕТ
УМЕНЬШИТЬ ВЕРОЯТНОСТЬ
РАЗНЫХ ОТВЕТОВ.



А ТЕПЕРЬ РАССЧИТАЕМ
НАШИ ДАННЫЕ ПО ГОЛЬФУ!





Следуя шагам 1–5, рассчитаем непредсказуемость и соотношение прироста информации в данных для гольфа.

Шаг 1

Так как в наборе данных D ответов «да» 9, а ответов «нет» – 5, то рассчитаем непредсказуемость по формуле ниже.

$$E(D) = -\frac{9}{14}\log_2\frac{9}{14} - \frac{5}{14}\log_2\frac{5}{14} = -0.643 \times (-0.637) - 0.357 \times (-1.495) = 0.94.$$

Шаг 2

Найдем непредсказуемость данных соответственно для ясной погоды, облачной и дождя.

$$E(\text{ясно}) = -\frac{2}{5}\log_2\frac{2}{5} - \frac{3}{5}\log_2\frac{3}{5} = -0.4 \times (-1.32) - 0.6 \times (-0.74) = 0.971.$$

$$E(\text{облачно}) = -\frac{4}{4}\log_2\frac{4}{4} - \frac{0}{4}\log_2\frac{0}{4} = 0 - 0 = 0.$$

$$E(\text{дождь}) = -\frac{3}{5}\log_2\frac{3}{5} - \frac{2}{5}\log_2\frac{2}{5} = -0.6 \times (-0.74) - 0.4 \times (-1.32) = 0.971.$$

Шаг 3

Возьмем эти величины в качестве весов к данным и рассчитаем непредсказуемость после разделения.

$$\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{4}{14} \times 0.971 = 0.694.$$

Шаг 4

Вычтя из непредсказуемости изначальных данных величину непредсказуемости данных после разделения, найдем отношение информационного выигрыша к внутренней информации, или же Gain.

$$\text{Gain}(D, \text{погода}) = 0.94 - 0.694 = 0.246.$$

Шаг 5

Таким же методом рассчитаем информационный выигрыш для других данных.

$$\text{Gain}(D, \text{температура}) = 0.029.$$

$$\text{Gain}(D, \text{влажность}) = 0.151.$$

$$\text{Gain}(D, \text{ветер}) = 0.048.$$

Следовательно, если первым вопросом для деления данных должен стать вопрос о погоде, непредсказуемость будет самой большой, а потом будет уменьшаться. После разделения данных можно использовать оставшиеся признаки в том же порядке.



В качестве метода расчета непредсказуемости данных вместо вышеуказанного способа можно воспользоваться коэффициентом Джини.

$$\text{Gini}(D) = 1 - P_{\text{Да}}^2 - P_{\text{Нет}}^2$$

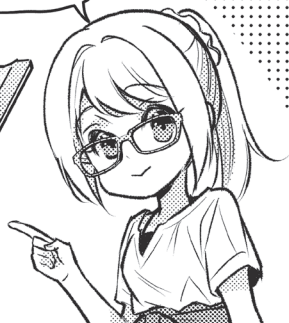
ТАК МОЖНО
ЧИСЛЕННО СРАВНИТЬ
ВЕЛИЧИНЫ ЭФФЕКТИВНЫХ
ВОПРОСОВ?



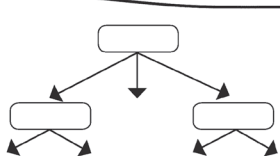
В ОСНОВЕ МЕТОДА,
КОТОРЫМ МЫ СЕГОДНЯ ПОЛЬЗОВАЛИСЬ,
ЛЕЖИТ "БРИТВА ОККАМА", КОТОРАЯ ГЛАСИТ:
"ВЫБИРАЙ САМУЮ ПРОСТУЮ ГИПОТЕЗУ
ДЛЯ ПРИМЕНЕНИЯ К ДАННЫМ".

**БРИТВА
ОККАМА**

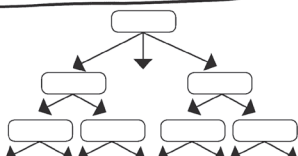
САМУЮ ПРОСТУЮ?



ЕСЛИ ГИПОТЕЗА СЛИШКОМ СЛОЖНАЯ, ТО ПОЛУЧЕННЫЕ
ДАННЫЕ МОЖНО ОБЪЯСНИТЬ СЛУЧАЙНОСТЬЮ,
НЕ ТАК ЛИ? А ЕСЛИ ОНА ПРОСТАЯ, ТО ВЕРОЯТНОСТЬ
СЛУЧАЙНОГО ОБЪЯСНЕНИЯ ДАННЫХ ПАДАЕТ.



Простая (короткая) гипотеза



Сложная (длинная) гипотеза

ОДНАКО ЕСЛИ,
СЛЕДУЯ ЭТОМУ МЕТОДУ,
ПОСТРОИТЬ РЕШАЮЩЕЕ
ДЕРЕВО, В КОТОРОМ
НЕ БУДЕТ ОШИБОК,
ТО МОЖНО ДОБИТЬСЯ
ПЕРЕОБУЧЕНИЯ -

КОГДА ДЕРЕВО
БУДЕТ СЛИШКОМ
ХОРОШО ПОДХОДИТЬ
К ДАННЫМ
ДЛЯ ОБУЧЕНИЯ.

ТЫ ХОЧЕШЬ СКАЗАТЬ,
ЧТО ЕСЛИ МЫ ПРИ ОБУЧЕНИИ ПОЛУЧИЛИ
НЕБОЛЬШОЕ РЕШАЮЩЕЕ ДЕРЕВО, ТО ОНО
И БУДЕТ ИСПОЛЬЗОВАТЬСЯ ДАЛЬШЕ?



Может
использоваться!

ЕСЛИ МЫ СМОГЛИ ПРОВЕСТИ ОБУЧЕНИЕ
НЕБОЛЬШОГО ДЕРЕВА, ПОСТРОЕННОГО С УЧЕТОМ
ID3-АЛГОРИТМА, ТО ВЕРОЯТНОСТЬ СЛУЧАЙНОСТЕЙ
БУДЕТ СНИЖЕНА. ДРУГИМИ СЛОВАМИ, СЛУЧАЙНОСТИ
НЕ БУДЕТ, БУДЕТ ЗАКОНОМЕРНОСТЬ.

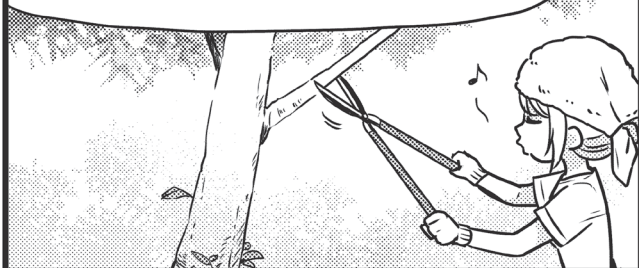
ПЕРЕ-
ОБУЧЕНИЯ?



ПОСКОЛЬКУ МОДЕЛЬ СЛИШКОМ ХОРОШО ОБЪЯСНЯЕТ ДАННЫЕ ДЛЯ ОБУЧЕНИЯ, ПРИ ВВОДЕ НОВЫХ ДАННЫХ ПРАВИЛЬНЫЕ ЗНАЧЕНИЯ НЕ ПОЛУЧАТСЯ.



ЧТОБЫ СПРАВИТЬСЯ С ПЕРЕОБУЧЕНИЕМ, МОЖНО ЛИБО ИЗНАЧАЛЬНО ОГРАНИЧИТЬ ТОЛЩИНУ ДЕРЕВА, ЛИБО ЖЕ ПОСЛЕ ОБУЧЕНИЯ ОБРЕЗАТЬ ВЕТВИ.



ДО ЭТОГО МЫ СТРОИЛИ РЕШАЮЩЕЕ ДЕРЕВО, КЛАССИФИЦИРУЯ ПО КАТЕГОРИЯМ, А В СЛУЧАЕ ЧИСЛЕННЫХ ПРИЗНАКОВ ПОПРОБУЕМ ОБУЧИТЬ ДЕРЕВО ПО МОДЕЛИ **ДИСКРЕТИЗАЦИИ**, КОТОРАЯ РАЗДЕЛЯЕТ РЯДЫ ЧИСЛОВЫХ ЗНАЧЕНИЙ НА НЕСКОЛЬКО ГРУПП.

МОЖНО РАЗДЕЛЯТЬ И ЧИСЛЕННЫЕ ПРИЗНАКИ.

Классификация по категориям

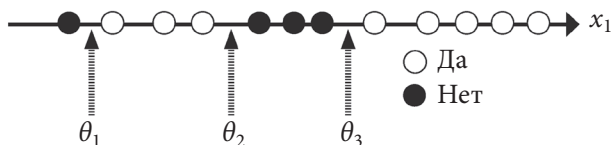
	Погода	Температура	Влажность	Ветер	Играем?
1	Ясно	Высокая	Высокая	Нет	Нет
2	Ясно	Высокая	Высокая	Да	Нет
3	Облачно	Высокая	Высокая	Нет	Да
4	Дождь	Средняя	Высокая	Нет	Да
5	Дождь	Низкая	Стандартная	Нет	Да

Классификация по числовым признакам

ID	Радиус	Текстура	Окружность	Опухоль
44	13.17	21.81	85.42	Злокачественная
45	18.65	17.60	123.7	Злокачественная
46	8.20	16.84	51.71	Доброкачественная
47	13.17	18.66	85.98	Злокачественная
48	12.02	14.63	78.04	Доброкачественная

ПОСКОЛЬКУ МЫ ХОТИМ НАЙТИ МЕСТА, ГДЕ НЕПРЕДСКАЗУЕМОСТЬ НИЗКАЯ, МЫ НЕ РАЗДЕЛЯЕМ ОДИНАКОВЫЕ КЛАССЫ.

КОГДА МЫ ИЩЕМ ГРАНИЦУ КЛАССА, ОНА БУДЕТ ПОКАЗАНА ВЕРТИКАЛЬНОЙ СТРЕЛКОЙ НА РИСУНКЕ. ЭТО ПОГРАНИЧНОЕ ЗНАЧЕНИЕ БУДЕТ СРЕДНИМ ОТ ЗНАЧЕНИЙ ДО И ПОСЛЕ НЕГО.



ВЫБЕРЕМ МЕСТО, ГДЕ НАИБОЛЕЕ ВЫСОК ИНФОРМАЦИОННЫЙ ВЫИГРЫШ. ВЫПОЛНЯЯ ТЕ ЖЕ ВЫЧИСЛЕНИЯ, ЧТО И В СЛУЧАЕ КАТЕГОРИЙ, МЫ ОБНАРУЖИМ, ЧТО ПРИ РАЗДЕЛЕНИИ С ПОРОГОМ θ_3 ИНФОРМАЦИОННЫЙ ВЫИГРЫШ НАИБОЛЕЕ ВЫСОК.



Попробуем построить модель логистической классификации и решающее дерево.

```
from sklearn.datasets import load_breast_cancer
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
```



В качестве данных возьмем обучающую выборку `breast_cancer`, которая показывает, доброкачественная или злокачественная опухоль.

```
breast_cancer = load_breast_cancer()
X = breast_cancer.data
y = breast_cancer.target
```



В `scikit-learn` и для регрессии, и для классификации используются в основном экземпляры класса, и обучение ведется путем метода `fit`. Сначала – логистическая классификация.

```
clf1 = LogisticRegression()
clf1.fit(X, y)
```



Находим те же коэффициенты, что и при регрессии.

```
for f, w in zip(breast_cancer.feature_names, clf1.coef_[0]) :
    print("{0:<23}: {1:6.2f}".format(f, w))
```

```

mean radius      : 2.10
mean texture     : 0.12
mean perimeter   : -0.06
...
worst concave points : -0.65
worst symmetry    : -0.69
worst fractal dimension: -0.11

```



Некоторые коэффициенты, имеющие большие положительные значения, влияют на положительный результат. Большие отрицательные значения влияют на отрицательный результат. Строим решающее дерево тем же методом, что и раньше.

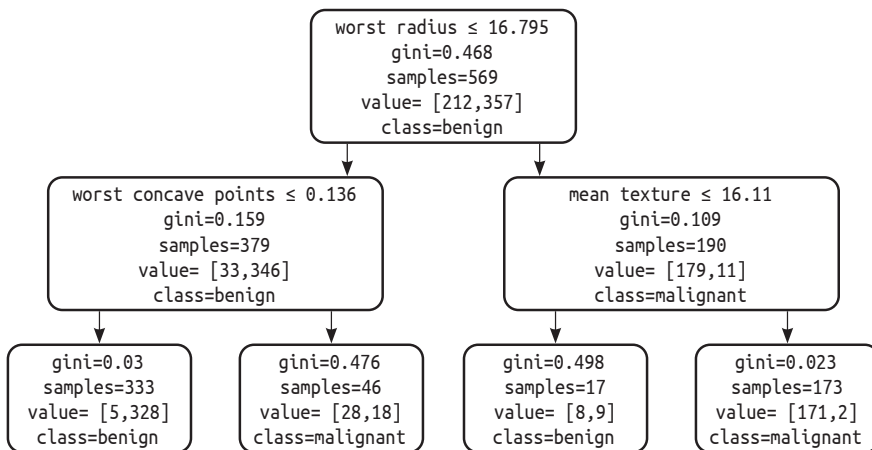
```

clf2 = DecisionTreeClassifier(max_depth=2)
clf2.fit(X, y)

```



В результате получится вот такое дерево. На вершине окажется параметр, показывающий средний радиус опухоли. Здесь происходит деление данных по этому параметру, радиус больше или меньше 16.795. Если он меньше, то дальше деление происходит по параметру «вмятины», и если он меньше 0.136, то опухоль доброкачественная, а если больше, то злокачественная. С другой стороны, если средний радиус опухоли больше, чем 16.795, то далее разделение происходит по параметру «текстура», и если он меньше 16.11, то опухоль доброкачественная, а если больше – то злокачественная.





Понедельник

Отдел здравоохранения
и благосостояния

ЧТО?

РАЗРАБОТЧИК ИГР?

НУ, НАШ СИСАДМИН КУАЗЭ МАКОТО. ОН ПРОГРАММИСТ
И ЕЩЕ В ШКОЛЕ СДЕЛАЛ ОЧЕНЬ ПОПУЛЯРНУЮ ИГРУ.

ого...

ОН НЕМНОГО ЛЕГКОМЫСЛЕННО
ОТНОСИТСЯ К РАБОЧИМ
ОБЯЗАННОСТЯМ, И ГОВОРЯТ,
ЧТО ОН РАЗРАБАТЫВАЕТ ИГРЫ
ПРЯМО НА РАБОТЕ...

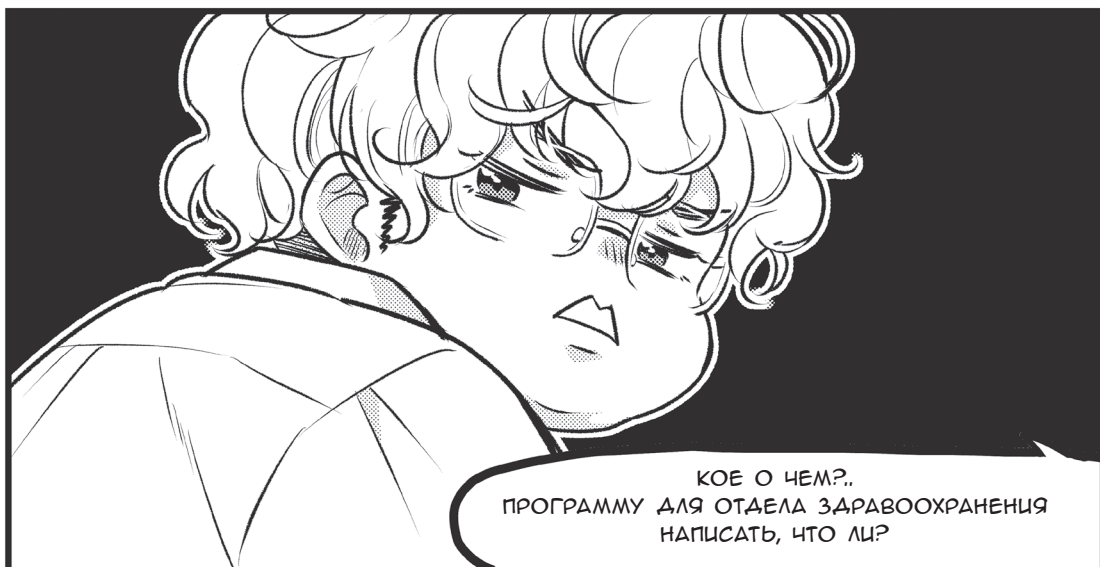
НИЧЕГО СЕБЕ ЧЕЛОВЕК.

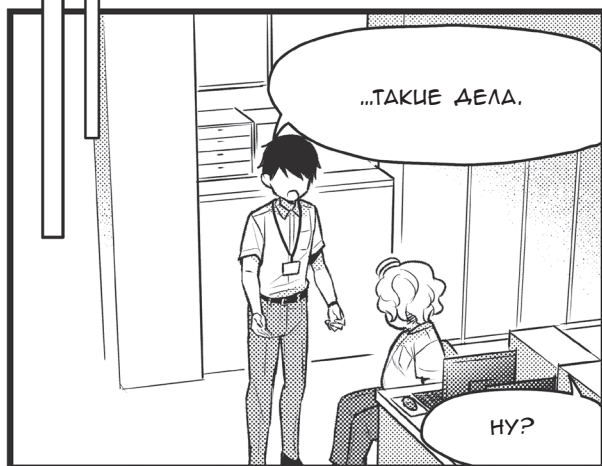
СПАСИБО
ЗА ИНФОРМАЦИЮ!

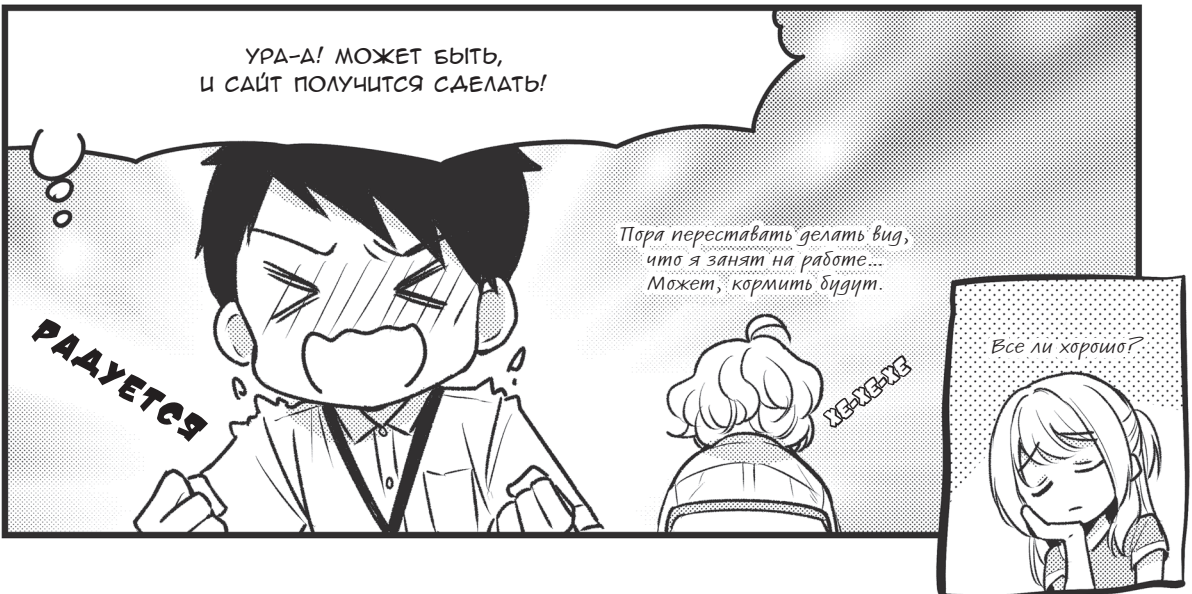
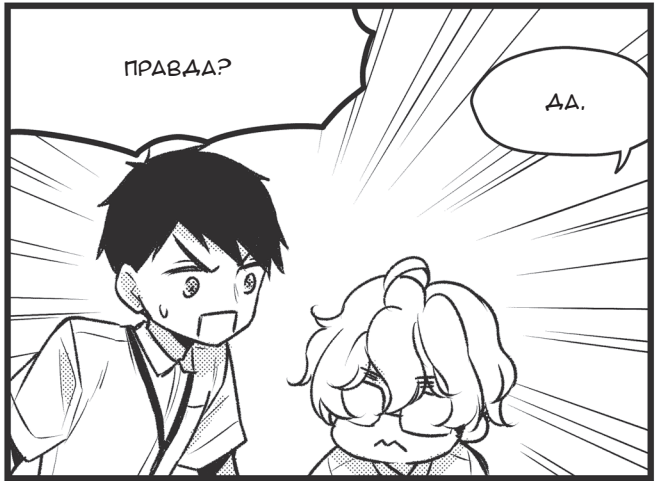
урачи!

Отдел поддержки
оборудования

ОН НАВЕРНЯКА ТУТ.







В кабинете у Саяка (3)

Математическое повторение (2)

То есть тебя угостил Киёхара-кун? Что-то как-то странно.



Именно! Он меня, видимо, за сэмпая не держит!

Наверное, хм... (Бедный Киёхара-кун! Не заметил этого!)



О чем сегодня поговорим? Было что-то непонятное?

А что такое e на стр. 51?



е? А, e в сигмоидной функции? Это число Непера, которое используется для основания натуральных логарифмов. Это бесконечная дробь, которая равняется 2.71828.

Логарифм – это показатель степени, в которую надо возвести основание, чтобы получилось число. Если основание равно двум, то понятно, почему в решающем дереве встречаются двоичные разряды, но что за странное основание e , и почему они натуральные?



При дифференцировании функции e^x получается e^x , а при дифференцировании $\log e^x$ получается $1/x$, поэтому это очень удобная штука. На самом деле должно быть наоборот, e – это число, которое обладает таким свойством.

Далее – производная вектора на стр. 53.



У... Тебя не проведешь! Давай расскажу!



Если в функции ошибки $E(w)$ изменить величину веса модели w , значение тоже изменится. Так как существует несколько весов, то функция ошибки становится функцией нескольких переменных. Если мы выразим сумму весов в виде вектора, то в функции ошибки появляется аргумент в виде вектора.

Ага, пока понимаю.





А теперь найдем частную производную по вектору. ∂ (закругленная d) – это обозначение дифференцирования по одной из переменных. Например, в формуле $\frac{\partial E}{\partial w_0}$ ∂ означает, что в формуле E нужно найти производную по переменной w_0 .

$$\nabla E = \frac{\partial E}{\partial \mathbf{w}} = \left(\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_d} \right)^T.$$



Тогда нам нужно дифференцировать вектор!



Именно. Это так называемый вектор градиента.



Понятно. А если производная функции одной переменной определяется углом наклона касательной, то производная функции нескольких переменных – это ее градиент?



Да. Как на картинке на стр. 54. Текущий вес обозначается точкой на склоне, и если спускаться в обратном направлении (вниз по склону), то можно немного приблизиться к минимуму функции ошибки.



Ага. Тогда все! А Киёхара-кун собирается делать сайт?

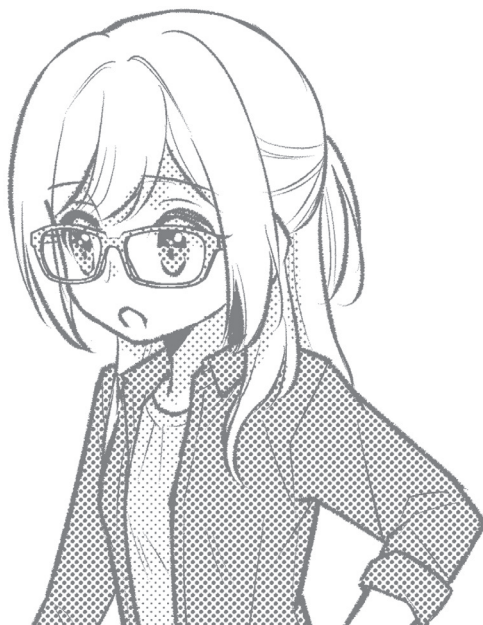


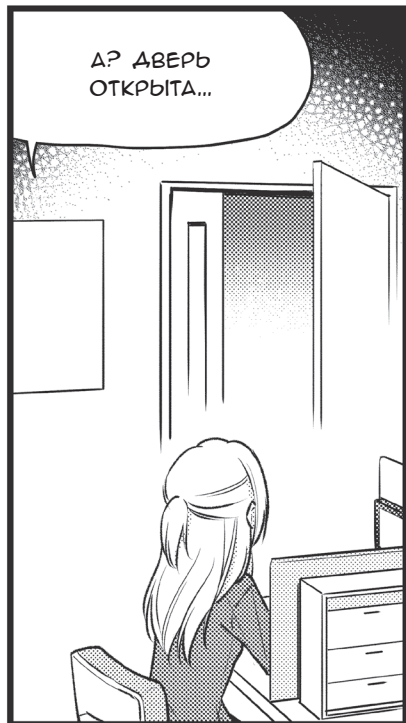
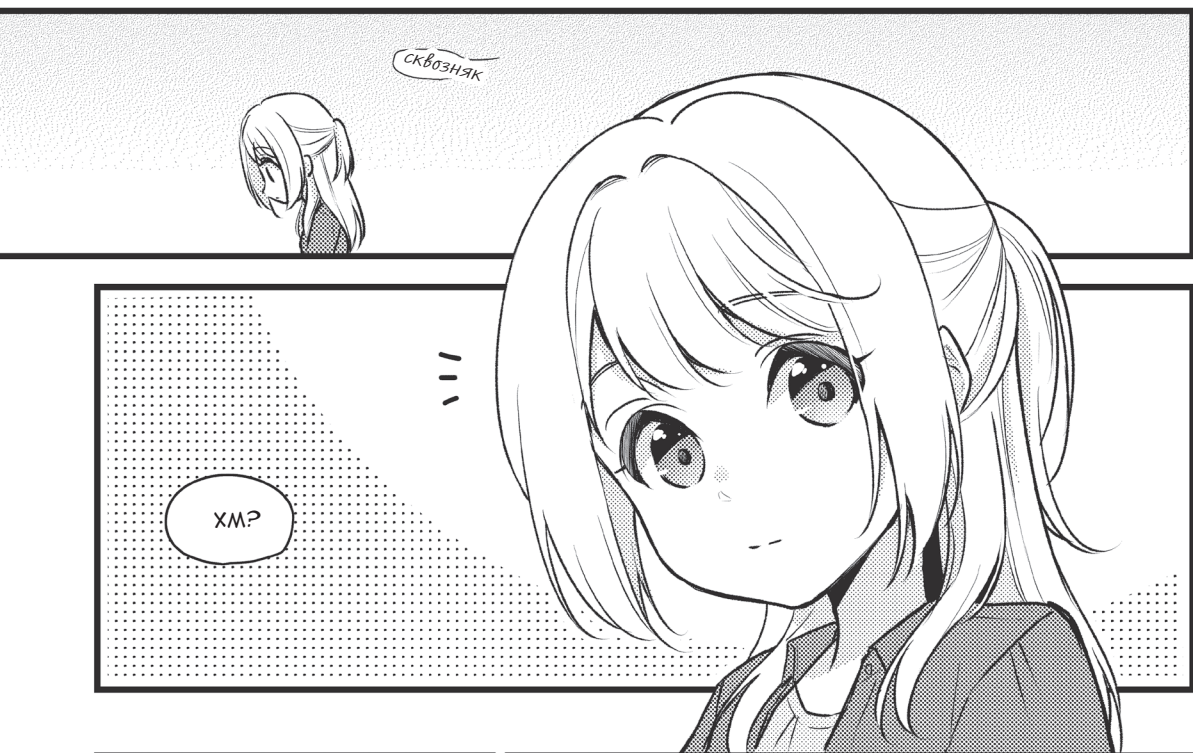
Я волнуюсь за него. Он же сбежал на полпути...

ГЛАВА 3

ОЦЕНКА РЕЗУЛЬТАТОВ

ОЦЕНКА РЕЗУЛЬТАТОВ
ОСОБЕННО ВАЖНА!





ДА ТЫ МЕНЯ ОПЯТЬ НАПУГАЛ!

ИЗВИНИТЕ...

СКОЛЬКО МО-О-О-ОЖНО-О-О-О!

Я ЖЕ ПРОСИЛА ПОТИШЕ,
НО ЭТО СЛИШКОМ ТИХО!

ИЗВИНИТЕ...

ЧТО СЛУЧИЛОСЬ?
Я ТРИ МЕСЯЦА ПОСЛЕ РЕСТОРАНА
ТЕБЯ НЕ ВИДЕЛА, ДУМАЛА,
ВСЕ ХОРОШО.

Я НАШЕЛ ПРОГРАММИСТА,
И ВМЕСТЕ МЫ СДЕЛАЛИ
САЙТ С ИСПОЛЬЗОВАНИЕМ
РЕШАЮЩЕГО
ДЕРЕВА, КОТОРЫЙ
С ВЕРОЯТНОСТЬЮ 100 %
ОПРЕДЕЛЯЕТ СТЕПЕНЬ РИСКА
ЗАБОЛЕВАНИЯ ДИАБЕТОМ.

Возможен ли
у вас диабет?

ДА

или

НЕТ

Выберите
нужный
ответ

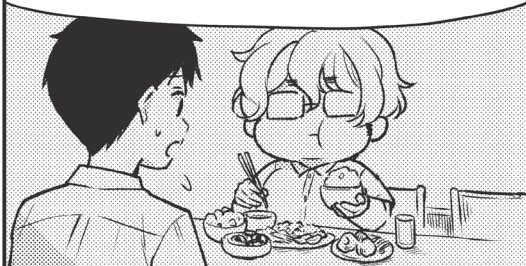
НАЧАТЬ

НО НА САЙТ МНОГО ЖАЛУЮТСЯ.

УЖЕ ТРОЕ ЛЮДЕЙ
НЕ ИЗ ГРУППЫ РИСКА
ПРОДОЛЖИЛИ ВЕСТИ
ОБЫЧНЫЙ ОБРАЗ ЖИЗНИ
И ОКАЗАЛИСЬ С ДИАБЕТОМ.

А ТЕ, КТО БЫЛ В ГРУППЕ
РИСКА, ИДУТ В БОЛЬНИЦУ,
И ИМ ГОВОРЯТ, ЧТО ВСЕ
НОРМАЛЬНО.

ПОЧЕМУ ВЕРОЯТНОСТЬ
ПРАВИЛЬНОГО ОТВЕТА СТОПРОЦЕНТНАЯ,
А КЛАССИФИКАЦИЯ ВСЕ РАВНО ВЫДАЕТ
ОШИБКУ? Я ПОСОВЕТОВАЛСЯ
С ПРОГРАММИСТОМ КУДЗЁ-САН...



ЧТО? ЖАЛОБЫ?
Я ТОЛЬКО НАСТРОИЛ ДАННЫЕ,
КАК ТЫ И ПРОСИЛ, И САМ НИЧЕГО
НЕ ПОНИМАЮ...



...ВОТ ЧТО ОН СКАЗАЛ.

СОВСЕМ НЕ МОГУ ПОНЯТЬ,
В ЧЕМ ДЕЛО И ГДЕ Я ОШИБСЯ.



КЦЁХАРА-КУН, А ТЫ ПРОВОДИЛ
ОЦЕНКУ ТЕСТОВЫХ ДАННЫХ
В КЛАССИФИКАТОРЕ?



ОЦЕНКУ ТЕСТОВЫХ
ДАННЫХ?



ВОТ И ПРИЧИНА...

НЕ ИМЕЕТ СМЫСЛА ПОЛЬЗОВАТЬСЯ
МАШИНЫМ ОБУЧЕНИЕМ,
ЕСЛИ ТЕСТОВЫЕ ДАННЫЕ
НЕ ПРОВЕРЯЮТСЯ
НА ДРУГИХ УСЛОВИЯХ.



НАДО БЫЛО ТЕБЕ РАССКАЗАТЬ
ОБ ОЦЕНКЕ ТЕСТОВЫХ ДАННЫХ...
ЭТО МОЯ ВИНА...

НЕТ-НЕТ, ЭТО НЕ ВАША ВИНА,
А МОЯ.

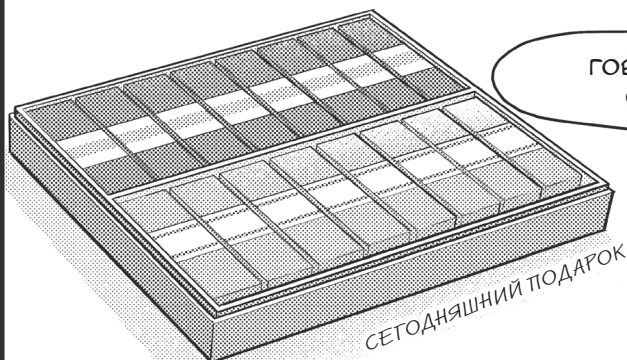
У ТЕБЯ ЕСТЬ ВРЕМЯ,
КИЁХАРА-КУН?

Я РАССКАЖУ
ТЕБЕ ОБ ЭТОМ.

ДА.

ПОЖАЛУЙСТА.

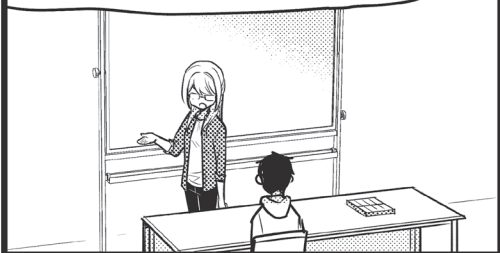
3.1. БЕЗ ПРОВЕРКИ ТЕСТОВЫХ ДАННЫХ НИКАК НЕЛЬЗЯ



ГОВОРИШЬ, КЛАССИФИКАТОР
С ТОЧНОСТЬЮ 100 %?

АА, ЕСЛИ ВЕРИТЬ ДАННЫМ
ДЛЯ ОБУЧЕНИЯ.

НО ЕСЛИ В ДАННЫХ ДЛЯ ОБУЧЕНИЯ
ВЫСОКАЯ ВЕРОЯТНОСТЬ
ПРАВИЛЬНОГО ОТВЕТА,
ТО В ЭТОМ НЕТ СМЫСЛА.



МОЖНО ПОСТРОИТЬ ДЕРЕВО
С ТОЧНОСТЬЮ РЕШЕНИЙ В 100 %
В СЛУЧАЕ, ЕСЛИ НЕТ ОГРАНИЧЕНИЙ
НА РАЗМЕР ДЕРЕВА И ЕСЛИ НЕТ
ПРОТИВОРЕЧИЯ В ДАННЫХ,
ТО ЕСТЬ КОГДА ОДИНАКОВЫЕ
ПРИЗНАКОВЫЕ ОПИСАНИЯ
ОТНЕСЕНЫ К РАЗНЫМ КЛАССАМ.

НО ЕСЛИ ПРИ ИСПОЛЬЗОВАНИИ ДАННЫХ
ИЗ ОБУЧАЮЩЕЙ ВЫБОРКИ ТОЧНОСТЬ РЕШЕНИЯ
100 %, ТО ПОЧЕМУ ЭТО ПЛОХО?



В ТАКИХ СИСТЕМАХ,
ЕСЛИ В ДАННЫХ ДЛЯ ОБУЧЕНИЯ
ВСЕ ЧЕТКО НАСТРОЕНО, ВЕЛИКА
ВЕРОЯТНОСТЬ ТОГО, ЧТО НОВЫЕ ДАННЫЕ
БУДУТ ИНТЕРПРЕТИРОВАТЬСЯ
НЕПРАВИЛЬНО?

КОРОЧЕ, ЧТО ЛУЧШЕ ДЕЛАТЬ?

3.2. ОБУЧАЮЩАЯ, ТЕСТОВАЯ И КОНТРОЛЬНАЯ ВЫБОРКИ



ЕСЛИ НЕЛЬЗЯ ПОЛУЧИТЬ
ДОСТАТОЧНО НЕИЗВЕСТНЫХ
ДАННЫХ, ТО И ТЕСТИРОВАНИЕ
НЕВОЗМОЖНО?

ЕСЛИ ИХ НЕ ХВАТАЕТ...
Я ОБ ЭТОМ.

ТАК... ТОЛЬКО В СЛУЧАЕ,
ЕСЛИ ЕСТЬ ДОСТАТОЧНО ДАННЫХ
И ДЛЯ ОБУЧЕНИЯ, И ДЛЯ ТЕСТОВОЙ ВЫБОРКИ,
МОЖНО ИСПОЛЬЗОВАТЬ ЭТОТ МЕТОД.

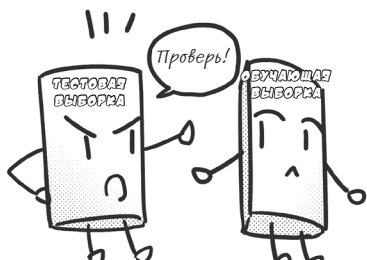
НО ДАЖЕ ЕСЛИ ДАННЫХ
ДОСТАТОЧНО, ТО НЕ ВСЕГДА
ЭТО ХОРОШИЙ МЕТОД.

ПОЧЕМУ?

У НАС ЕСТЬ ГИПЕРПАРАМЕТР,
КОТОРЫЙ ВЛИЯЕТ НА РЕЗУЛЬТАТЫ ОБУЧЕНИЯ?
ВЕС ДОПОЛНИТЕЛЬНОГО ЧЛЕНА ПРИ ЛИНЕЙНОЙ
РЕГРЕССИИ ИЛИ ЖЕ ТОЛЩИНА ДЕРЕВА
ПРИ ПОСТРОЕНИИ РЕШАЮЩЕГО ДЕРЕВА?

НУЖНА ПРОВЕРКА АДЕКВАТНОСТИ
ВЕЛИЧИНЫ ГИПЕРПАРАМЕТРА, КОТОРЫЙ
ИСПОЛЬЗУЕТСЯ ПРИ ОБУЧЕНИИ.

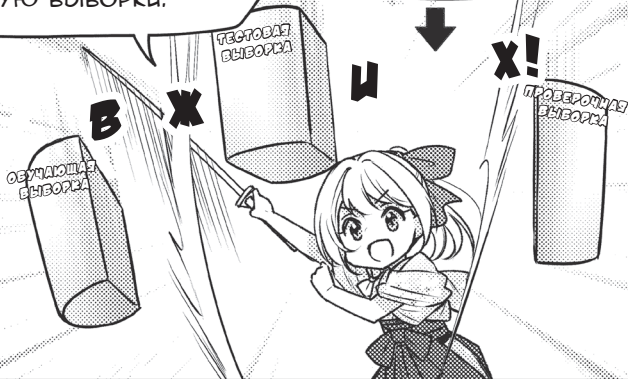
ЕСЛИ ИСПОЛЬЗОВАТЬ ТЕСТОВУЮ ВЫБОРКУ
ДЛЯ ОЦЕНКИ, ТО ДАННЫЕ ИЗ ТЕСТОВОЙ
ВЫБОРКИ НЕ СМОГУТ РАССМАТРИВАТЬСЯ
КАК НЕИЗВЕСТНЫЕ ДАННЫЕ.



ДЕЙСТВИТЕЛЬНО.

МЕТОД ПРОВЕРКИ НА ЗАРЕЗЕРВИРОВАННЫХ ДАННЫХ
ДЛЯ ОЦЕНКИ ЭФФЕКТИВНОСТИ ЗАКЛЮЧАЕТСЯ В ТОМ,
ЧТО ДАННЫЕ ДЕЛЯТ НА ТРИ ГРУППЫ - ОБУЧАЮЩУЮ,
ТЕСТОВУЮ И КОНТРОЛЬНУЮ ВЫБОРКИ.

ВСЕ ДАННЫЕ

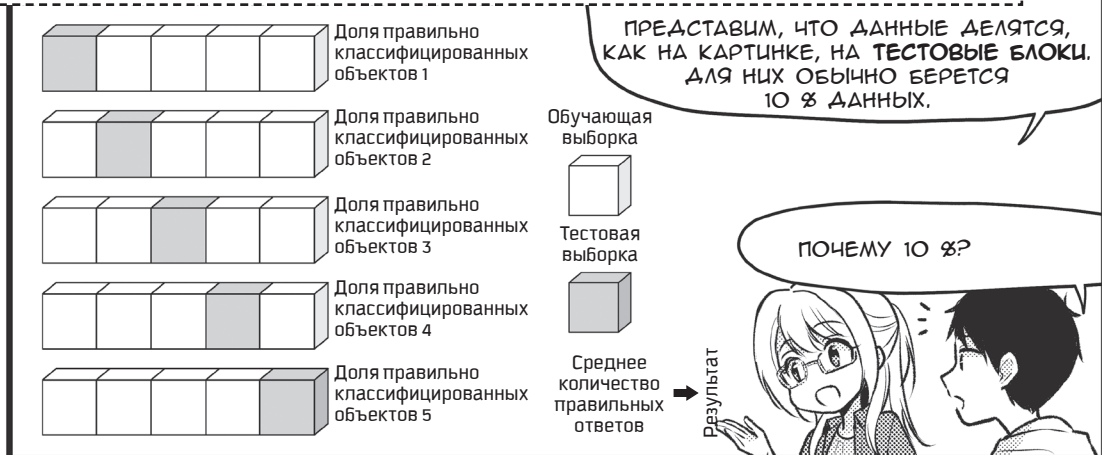


Ничего себе

АПЛОДИСМЕНТЫ



3.3. МЕТОД ПЕРЕКРЕСТНОЙ ПРОВЕРКИ (КРОСС-ВАЛИДАЦИИ)



МЕТОД, ПРИ КОТОРОМ ДЛЯ ОЦЕНКИ
ИСПОЛЬЗУЕТСЯ ОДИН БЛОК ДАННЫХ,
НАЗЫВАЕТСЯ **КОНТРОЛЕМ ПО ОТДЕЛЬНЫМ
ОБЪЕКТАМ**.

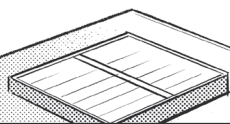
А РАЗВЕ У КРОСС-ВАЛИДАЦИИ
НЕТ НЕДОСТАТКОВ?



ПРИ ИСПОЛЬЗОВАНИИ ЭТОГО МЕТОДА
ОЦЕНКА ЭФФЕКТИВНОСТИ МОЖЕТ
ЗАНЯТЬ НЕКОТОРОЕ ВРЕМЯ,
ЕСЛИ КОЛИЧЕСТВО ТЕСТОВЫХ БЛОКОВ
 m ВЕЛИКО.

m

А ЕСЛИ ЭТО ЧИСЛО m
СЛИШКОМ НИЗКОЕ, ТО СИСТЕМА
НЕ СМОЖЕТ ОБУЧИТЬСЯ.



ОДНАКО ПОСКОЛЬКУ МЕТОД ПЕРЕКРЕСТНОЙ
ПРОВЕРКИ ИСПОЛЬЗУЕТСЯ, КОГДА ДАННЫХ
НЕ ТАК УЖ И МНОГО И КОГДА ХВАТАЕТ ОДНОГО
РАУНДА ДЛЯ ОЦЕНКИ, ТО ЕМУ МОЖНО ДОВЕРЯТЬ.
КОНТРОЛЬ ПО ОТДЕЛЬНЫМ ЭЛЕМЕНТАМ ТОЖЕ
МОЖНО ИСПОЛЬЗОВАТЬ, ЕСЛИ ТАСОВАТЬ
ФРАГМЕНТЫ И ПРОДЕЛАТЬ ПРОВЕРКУ
НЕСКОЛЬКО РАЗ.



А ТЕПЕРЬ РАССМОТРИМ
ЭТИ МЕТОДЫ ОЦЕНКИ
В ДЕТАЛЯХ, С ЧИСЛОВЫМИ
ВЕЛИЧИНАМИ.



ПОЖАЛУЙСТА!

3.4. ДОЛЯ ПРАВИЛЬНО ПРЕДСКАЗАННЫХ ОБЪЕКТОВ, ТОЧНОСТЬ, ПОЛНОТА И F-МЕРА

ДО ЭТОГО
МЫ ИСПОЛЬЗОВАЛИ ПОКАЗАТЕЛЬ
КОЛИЧЕСТВА ПРАВИЛЬНО ОЦЕНЕННЫХ
ОБЪЕКТОВ ДЛЯ ПРИМЕРНОЙ ОЦЕНКИ
ТОЧНОСТИ, НО ОН ОСНОВЫВАЕТСЯ НА
КОЛИЧЕСТВЕ ПРАВИЛЬНО РАСПРЕДЕ-
ЛЕННЫХ ПО КЛАССАМ ДАННЫХ
В ТЕСТОВОЙ ВЫБОРКЕ.

А ТЕПЕРЬ РАССМОТРИМ
ЧУТЬ-ЧУТЬ ПОДРОБНЕЕ
МЕТОДЫ ОЦЕНКИ
ЭФФЕКТИВНОСТИ САМОЙ
КЛАССИФИКАЦИИ.



ПРЕЖДЕ ВСЕГО...
ДЛЯ ПРОСТОТЫ
РАССМОТРИМ
МЕТОД ОЦЕНКИ,

ПРИ КОТОРОМ ДАННЫЕ РАЗДЕЛЕНЫ
НА ДВА КЛАССА. ДОПУСТИМ, У НАС ВОПРОСЫ
ТАКОГО ПЛАНА: ЕСТЬ БОЛЕЗНЬ ИЛИ НЕТ,
СПАМ ПИСЬМО ИЛИ НЕ СПАМ.
ДАННЫЕ, КОТОРЫЕ ПОДХОДЯТ,
БУДУТ НАЗЫВАТЬСЯ **ИСТИННЫМИ ПРИМЕРАМИ**,
А КОТОРЫЕ НЕ ПОДХОДЯТ – **ЛОЖНЫМИ**.

ИСТИННЫЕ



ЛОЖНЫЕ



ТО ЕСТЬ СПАМ У НАС – ИСТИННЫЙ?



НЕМНОГО СТРАННО, ДА,
НО БУДЕМ ОСНОВЫВАТЬСЯ НА ТОМ,
ПОДХОДИТ НАМ РЕЗУЛЬТАТ ИЛИ НЕТ, ПОЭТОМУ
РАЗДЕЛИМ ИХ НА ИСТИННЫЕ И ЛОЖНЫЕ.

ЕСЛИ МЫ ОБЪЕДИНИМ ИСТИННЫЙ
КЛАСС, РАЗДЕЛЕННЫЙ НА ИСТИННОЕ
«ДА» И ИСТИННОЕ «НЕТ», С ПРЕДСКА-
ЗАННЫМ КЛАССОМ, КОТОРЫЙ ТОЖЕ
РАЗДЕЛЕН НА «ДА» И «НЕТ»,
ТО ПОЛУЧИМ ЧЕТЫРЕ ГРУППЫ.
ИСТИННОЕ «ДА» И ИСТИННОЕ «НЕТ» –
ЭТО РЕЗУЛЬТАТЫ ОТВЕТОВ «ДА»
ИЛИ «НЕТ», А ПРОГНОЗЫ ПО КЛАССИ-
ФИКАТОРУ БУДУТ ДЕЛИТЬСЯ
НА ПРЕДСКАЗАННОЕ «ДА»
И ПРЕДСКАЗАННОЕ «НЕТ».

	Предсказанное «да»	Предсказанное «нет»
Истинное «да»	30	20
Истинное «нет»	10	40

ЭТО МАТРИЦА НЕТОЧНОСТЕЙ.
ЦИФРЫ У НАС ТОЛЬКО ДЛЯ ПРИМЕРА,
НО ЧТО ТЫ МОЖЕШЬ ПОНЯТЬ ОТСЮДА,
КИЁХАРА-КУН?

НУ...

	Предсказанное «да»	Предсказанное «нет»
Истинное «да»	30	20
Истинное «нет»	10	40

	Предсказанное «да»	Предсказанное «нет»
Истинное «да»	30	20
Истинное «нет»	10	40

Точность: $30 + 40 = 70$

Ошибки: $20 + 10 = 30$

ЕСЛИ СЛОЖИТЬ РЕЗУЛЬТАТЫ
ПО ЭТОЙ ДИАГОНАЛИ,
ТО УЗНАЕМ ТОЧНОСТЬ, А ЕСЛИ
ПО ДРУГОЙ, ТО ПОЛУЧИМ
КОЛИЧЕСТВО ОШИБОК?

ИМЕННО!

ЕСЛИ МЫ СЛОЖИМ РЕЗУЛЬТАТЫ «ДА»,
ТО ПОЛУЧИМ 50, ИЗ НИХ 30 ПРЕДСКАЗАНЫ
ЛОГИЧЕСКИМ КЛАССИФИКАТОРОМ,
А 20 ОШИБОЧНЫ.

ПРОСТЕЙШИЙ ПОКАЗАТЕЛЬ, КОТОРЫЙ
МОЖНО УЗНАТЬ ИЗ ЭТОЙ ТАБЛИЦЫ,
РАССЧИТЫВАЕТСЯ ПУТЕМ ДЕЛЕНИЯ
КОЛИЧЕСТВА ПРАВИЛЬНО ПРЕДСКАЗАН-
НЫХ КЛАССИФИКАТОРОМ ОТВЕТОВ НА
ОБЩЕЕ КОЛИЧЕСТВО ДАННЫХ.



	Предсказанное «да»	Предсказанное «нет»
Истинное «да»	30	20
Истинное «нет»	10	40

В СЛУЧАЕ ЭТОЙ ТАБЛИЦЫ
ПОЛУЧИТСЯ 0.7 - ЭТО ДОЛЯ ПРАВИЛЬНО
ПРЕДСКАЗАННЫХ ОБЪЕКТОВ.

С ПЕРВОГО ВЗГЛЯДА КАЖЕТСЯ,
ЧТО ЭТОГО ДОСТАТОЧНО,
НО ДЛЯ ОЦЕНКИ МАШИННОГО
ОБУЧЕНИЯ ЭТО ЕЩЕ НЕ ВСЕ.

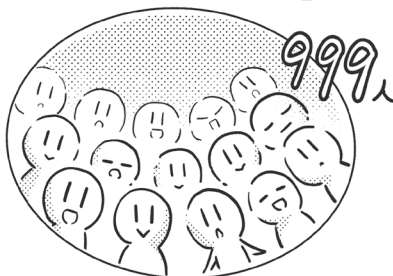
ДАВАЙ ПРЕДСТАВИМ СИТУАЦИЮ,
КОГДА ЧИСЛО ЛЮДЕЙ, КОТОРЫЕ
НЕ БОЛЕЮТ, НАМНОГО БОЛЬШЕ
ЧИСЛА ЗАБОЛЕВШИХ.

ТО ЕСТЬ СКЛАДЫВАТЬ
ПОЛОЖИТЕЛЬНЫЕ
И ОТРИЦАТЕЛЬНЫЕ РЕЗУЛЬТАТЫ
НЕ ПОМОЖЕТ?

ЧТО?
ДОЛЕЙ ПРАВИЛЬНО ПРЕДСКАЗАН-
НЫХ ОТВЕТОВ НЕ ОБОЙТИСЬ?

ИМЕННО!
ДОПУСТИМ, ЧТО БОЛЕЕТ
ОДИН ЧЕЛОВЕК ИЗ ТЫСЯЧИ,
ВСЕ НЕГАТИВНЫЕ РЕЗУЛЬТАТЫ
ПОВЛИЯЮТ НА ЛОГИЧЕСКИЙ
КЛАССИФИКАТОР И НА ДОЛЮ
ПРАВИЛЬНЫХ ОТВЕТОВ.

ЧТОБЫ РАСПОЗНАТЬ
ТАКИЕ СИТУАЦИИ, РЕЗУЛЬТАТЫ
МАШИННОГО ОБУЧЕНИЯ
НЕОБХОДИМО ПЕРЕПРОВЕРЯТЬ.



А, ОНА БУДЕТ 0.999.

А НУ РАССКАЖИТЕ МНЕ!

ЭТО СНОВА ТА ЖЕ САМАЯ
МАТРИЦА НЕТОЧНОСТЕЙ.

	Предсказанное «да»	Предсказанное «нет»
Истинное «да»	Истинное «за» (ИЗ)	Ложное «против» (ЛП)
Истинное «нет»	Ложное «за» (ЛЗ)	Истинное «против» (ИП)



НАПРИМЕР, ЛЕВЫЙ ВЕРХНИЙ ЭЛЕМЕНТ НАЗЫВАЕТСЯ ИСТИННЫМ «ЗА», ПОТОМУ ЧТО КЛАССИФИКАТОР ПРЕДСКАЗАЛ «ДА» И ЭТО СОВПАЛО С ИСТИННЫМ РЕЗУЛЬТАТОМ. ДЛЯ СОКРАЩЕНИЯ ИСПОЛЬЗУЕТСЯ АББРЕВИАТУРА ИЗ.

А «ЛОЖНОЕ ПРОТИВ» – ЭТО КОГДА БЫЛО ПРЕДСКАЗАНО «НЕТ», НО РЕЗУЛЬТАТ НА САМОМ ДЕЛЕ «ДА».



ТАКИМ ОБРАЗОМ, ДОЛЯ ПРАВИЛЬНО ПРЕДСКАЗАННЫХ ДАННЫХ БУДЕТ ВЫЧИСЛЯТЬСЯ ПО ФОРМУЛЕ НИЖЕ:

$$\text{Доля правильно предсказанных данных} = \frac{\text{ИЗ} + \text{ИП}}{\text{ИЗ} + \text{ЛП} + \text{ЛЗ} + \text{ИП}}$$



НУЖНО РАЗДЕЛИТЬ КОЛИЧЕСТВО ПРАВИЛЬНЫХ ОТВЕТОВ НА ОБЩЕЕ КОЛИЧЕСТВО ДАННЫХ.

ПОГОВОРИМ О **ТОЧНОСТИ**. ОНА ПОКАЗЫВАЕТ, МОЖНО ЛИ ДОВЕРЯТЬ КЛАССИФИКАТОРУ ПРИ ДЕЛЕНИИ ДАННЫХ.

ТОЧНОСТЬ ВЫЧИСЛЯЕТСЯ ПО ПРЕДСТАВЛЕННОЙ ФОРМУЛЕ.



ДОПУСТИМ, МЫ ХОТИМ УЗНАТЬ, НАСКОЛЬКО ТОЧНО МЫ ОПРЕДЕЛИЛИ БОЛЕЗНЬ.

$$\text{Точность} = \frac{\text{ИЗ}}{\text{ИЗ} + \text{ЛЗ}}$$

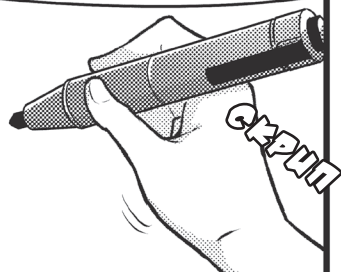


МЫ ДЕЛИМ КОЛИЧЕСТВО ИСТИННЫХ «ЗА» НА СУММУ «ДА», ПРЕДСКАЗАННЫХ КЛАССИФИКАТОРОМ.

ПОСЛЕ ТОЧНОСТИ ЦАЕТ **ПОЛНОТА** ДАННЫХ. ОНА УКАЗЫВАЕТ, НАСКОЛЬКО АДЕКВАТНО ОЦЕНИВАЮТСЯ ИСТИННЫЕ "ДА". НАПРИМЕР, МЫ МОЖЕМ УЗНАТЬ, НАСКОЛЬКО АДЕКВАТНА ВЫБОРКА НА ОСНОВАНИИ КОЛИЧЕСТВА НА САМОМ ДЕЛЕ ЗАБОЛЕВШИХ ЛЮДЕЙ.

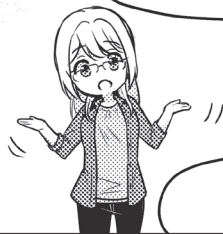
ЧТОБЫ НАЙТИ ПОЛНОТУ, НАДО РАЗДЕЛИТЬ КОЛИЧЕСТВО ИСТИННЫХ "ДА", ПРАВИЛЬНО ПРЕДСКАЗАННЫХ КЛАССИФИКАТОРОМ, НА ОБЩУЮ СУММУ ИСТИННЫХ "ДА".

$$\text{Полнота} = \frac{\text{ИЗ}}{\text{ИЗ} + \text{ЛП}}$$



ИМЕННО!

ТОЧНОСТЬ И ПОЛНОТА НАХОДЯТСЯ В ТАКИХ ОТНОШЕНИЯХ: ЕСЛИ ИСПОЛЬЗОВАТЬ ОДИН ИЗ ЭТИХ ПАРАМЕТРОВ, ДРУГОЙ ТЕРЯЕТ ЗНАЧИМОСТЬ, И НАОБОРОТ.



ЭТО КАК?

ДОПУСТИМ, ТОЧНОСТЬ КЛАССИФИКАТОРА, КОТОРЫЙ ВЫДАЕТ ПОЛОЖИТЕЛЬНЫЙ РЕЗУЛЬТАТ, БУДЕТ ВЫШЕ, ЕСЛИ МЫ ЗАРАНЕЕ ЗНАЕМ, ЧТО В НАШИХ ДАННЫХ ЕСТЬ БОЛЬШОЕ КОЛИЧЕСТВО ЛЮДЕЙ, ЗАБОЛЕВШИХ ТОЙ ИЛИ ИНОЙ БОЛЕЗНЬЮ.



ХМ. А ЕСЛИ У ЧЕЛОВЕКА ЛИШЬ НЕЗНАЧИТЕЛЬНЫЙ СИМПТОМ И МЫ ТОЧНО НЕ ЗНАЕМ, БОЛЕН ОН ИЛИ НЕТ, ТО БОЛЕЗНЬ И ПРОГЛЯДЕТЬ МОЖНО?



ДА! ПОТОМУ ЧТО ПОЛНОТА ДАННЫХ В ЭТОМ СЛУЧАЕ НИЗКАЯ.

И НАОБОРОТ, ЕСЛИ В ПРИОРИТЕТЕ ВЫСОКАЯ ПОЛНОТА ДАННЫХ, ТО ПРИ МАЛЕЙШЕМ СОМНЕНИИ КЛАССИФИКАТОР БУДЕТ ВЫДАВАТЬ РЕЗУЛЬТАТЫ "ДА".



ЧТО?

ВЕРОЯТНОСТЬ ПРОПУСТИТЬ ЗАБОЛЕВШЕГО НИЗКАЯ, ОДНАКО ПРИДЕТСЯ ПРОВОДИТЬ БОЛЬШОЕ КОЛИЧЕСТВО МЕДОСМОТРОВ ЛЮДЕЙ, КОТОРЫЕ НЕ ЗАБОЛЕЛИ.



ПОНЯТНО.

ТОЧНОСТЬ И ПОЛНОТУ ДАННЫХ ОБЪЕДИНЯЕТ ТАК НАЗЫВАЕМАЯ F-МЕРА.

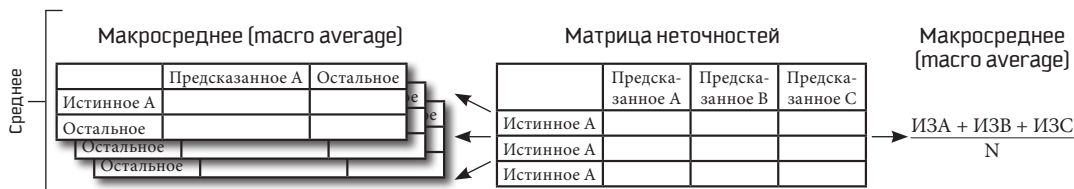
$$F\text{-мера} = 2 \times \frac{\text{Точность} \times \text{Полнота}}{\text{Точность} + \text{Полнота}}$$

ЭТО СРЕДНЕЕ ГАРМОНИЧЕСКОЕ?



А ЧТО, ЕСЛИ У НАС ТРИ КЛАССИФИКАТОРА?

МАТРИЦА НЕТОЧНОСТЕЙ БУДЕТ 3x3.

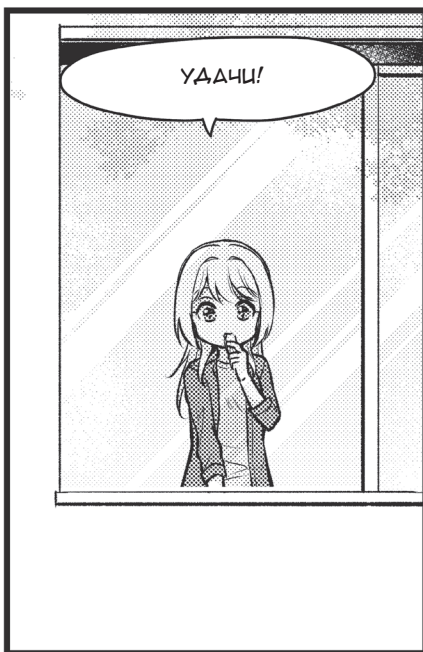
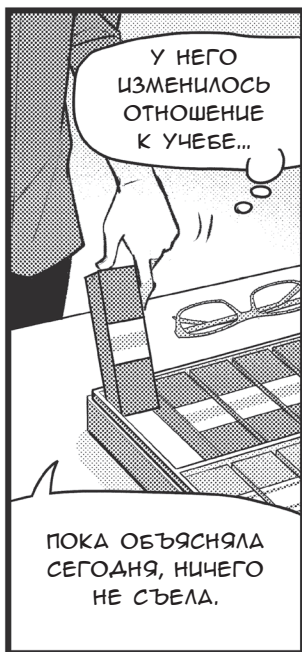
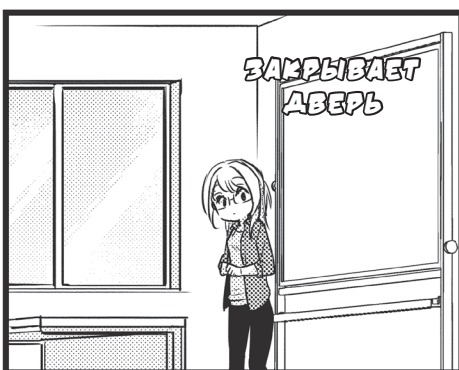
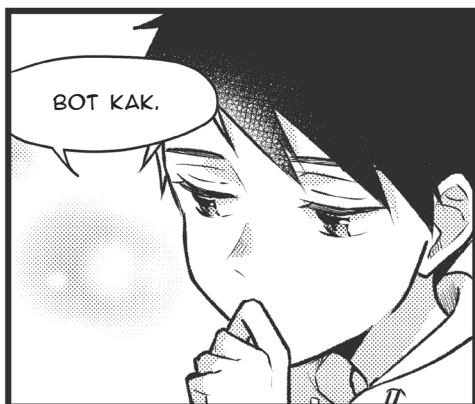


У НАС ЕСТЬ MACRO AVERAGE - ЭТО НАХОЖДЕНИЕ СРЕДНЕГО ИЗ ЭФФЕКТИВНОСТИ КАЖДОГО КЛАССА В МАТРИЦЕ НЕТОЧНОСТЕЙ, И MICRO AVERAGE, КОГДА СКЛАДЫВАЮТСЯ ПОКАЗАТЕЛИ ИЗ, ЛЗ, ЛП, ЛП КАЖДОГО КЛАССА И ДЕЛЯТСЯ НА ОБЩЕЕ КОЛИЧЕСТВО ДАННЫХ. MICRO AVERAGE ОТРАЖАЕТ СООТНОШЕНИЕ ДАННЫХ С ВЕЛИЧИНОЙ ОЦЕНКИ.

ТО ЕСТЬ ВСЕ ЭТИ ЗНАЧЕНИЯ НАДО ИСПОЛЬЗОВАТЬ?

В ЗАВИСИМОСТИ ОТ ЗАДАЧИ МОЖНО СОСРЕДОТОЧИТЬСЯ НА ТОЧНОСТИ, А ИНОГДА НА ПОЛНОТЕ. ЕСЛИ ЖЕ НАМ НУЖНЫ ОБА ПОКАЗАТЕЛЯ, ТО ЛУЧШЕ ВСЕГО ИСПОЛЬЗОВАТЬ F-МЕРУ.





Отдел здравоохранения
и благосостояния

ТАК, ПЕРЕКРЕСТНАЯ
ПРОВЕРКА.

ЧИСЛО F... 0,60

ТЯЖКО ВЗДЫХАЕТ

САМО СОБОЙ, СКОЛЬКО ТАМ ОШИБОК...
САЙТ НАДО ЗАКРЫТЬ.

ИЗ-ЗА СВОЕЙ ГЛУПОСТИ
Я ДОСТАВИЛ ВСЕМ ЖИТЕЛЯМ
СТОЛЬКО БЕСПОКОЙСТВ!

И ВСЕ ИЗ-ЗА ТОГО,
ЧТО Я СБЕЖАЛ ОТ НЕЕ!

ДА ЧТО Ж
Я ТАКОЕ СДЕЛАЛ...

ЧЕК

УБЕГАЕТ

В кабинете у Саяка (4)

Математическое повторение (3)



Сегодняшняя беседа была довольно проста с точки зрения математики, однако она важна для машинного обучения. Можно скопипастить полученные данные, но нехорошо, если возникает ситуация «все сделал, но не понимаю, что получилось».

Я немного не понимаю число F . Это всего лишь среднее от точности и аккуратности?



Обычное среднее – это арифметическое среднее. Если у нас есть числа a и b , то оно определяется по формуле $(a + b)/2$. Так определяют средний балл теста, температуру и т. п., но для прогнозов используется другой метод.



Что касается точности или полноты, то среднее находится другим способом.

Почему не используется среднее арифметическое?



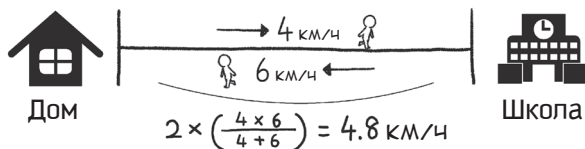
Ну, давай рассмотрим пример из повседневной жизни – скорость. Как она определяется?

Расстояние, поделенное на время!





Да. Возьмем такой пример.



Вот когда ты идешь в школу, скорость 4 км/ч.



А когда ты возвращаешься из школы, то 6 км/ч. Расстояние одинаковое, верно? Но какая будет средняя скорость?

Расстояние не указано, поэтому проведем вычисления так. В одну сторону расстояние будет равно x , тогда туда и обратно – $2x$. Время туда будет равно $x/4$, а обратно – $x/6$, и средняя скорость будет рассчитана по формуле:



$$\frac{2x}{\frac{x}{4} + \frac{x}{6}} = 2 \times \frac{x}{\frac{x(4+6)}{4 \times 6}} = 2 \times \frac{4 \times 6}{4 + 6} = \frac{48}{10} = 4.8$$

Ага, x исчез. То есть не 5 км/ч, а 4,8 км/ч.



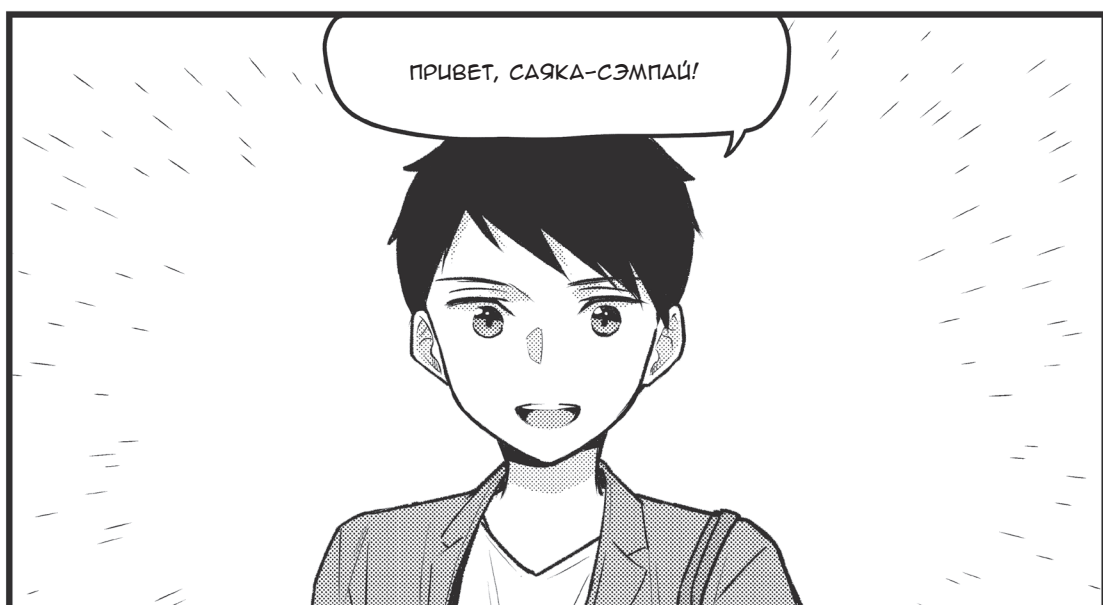
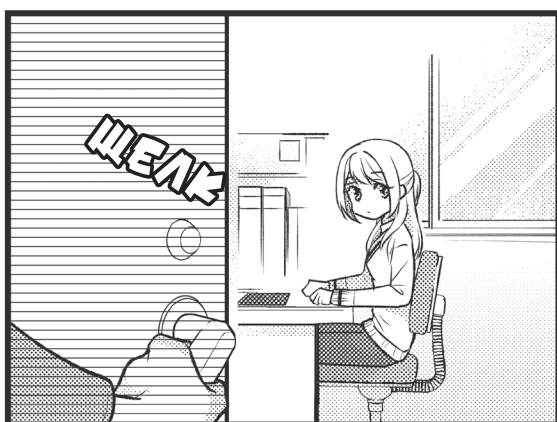
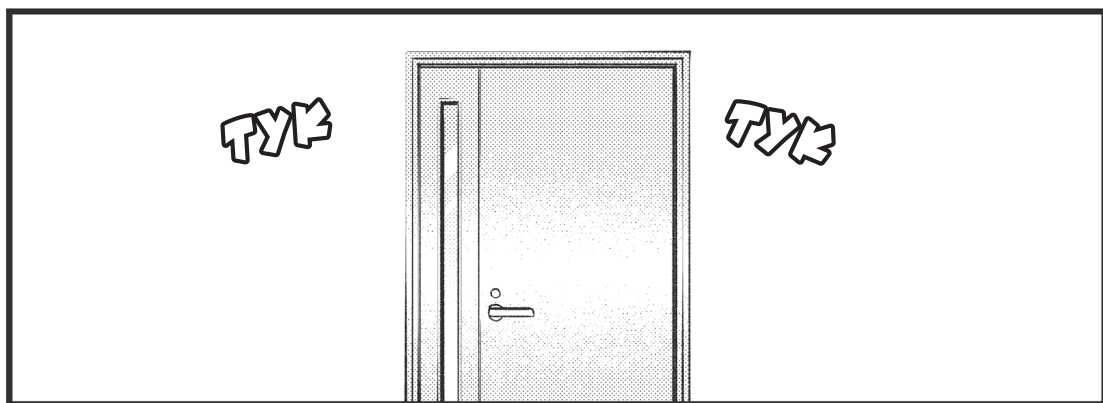
Да. Эта формула $2 \times \left(\frac{4 \times 6}{4 + 6} \right)$ используется и для нахождения F -меры. Это среднее гармоническое значение.

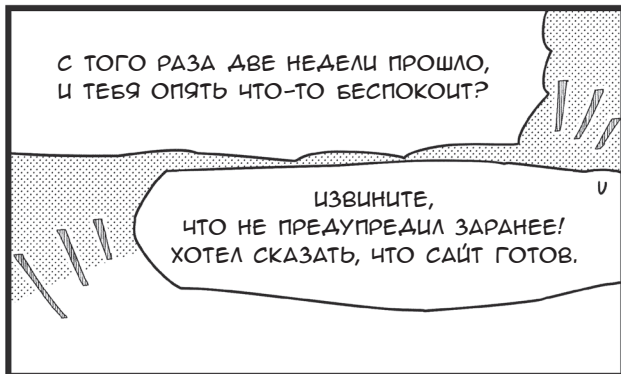
ГЛАВА 4

ГЛУБОКОЕ ОБУЧЕНИЕ

ГЛУБОКОЕ ОБУЧЕНИЕ
ДЛЯ РАСПОЗНАВАНИЯ
КАРТИНОК!

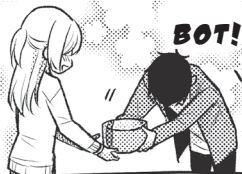






ГЛУБОКОЕ
ОБУЧЕНИЕ?

Н-ДА... НУ, СНАЧАЛА
РАССКАЖИ, ЧТО И ЗАЧЕМ
ТЫ СЕГОДНЯ ПРИНЕС.



КОНЕЧНО...
ВОТ СЛАДОСТИ!

НЕДАВНО, КОГДА Я СИДЕЛ,
РАССТРОЕННЫЙ ИЗ-ЗА САЙТА,
ПРИШЕЛ ЧЕЛОВЕК ИЗ ДЕПАРТАМЕНТА
СЕЛЬСКОГО ХОЗЯЙСТВА И ПОПРОСИЛ
МЕНЯ КОЕ О ЧЕМ.

МЫ ХОТЕЛИ БЫ,
ЧТОБЫ КТО-НИБУДЬ СДЕЛАЛ
АВТОМАТИЧЕСКУЮ СИСТЕМУ, КОТОРАЯ
СМОГЛА БЫ СОРТИРОВАТЬ ВИНОГРАД.

МЕСТНАЯ КОМПАНИЯ ПО ИЗГОТОВЛЕНИЮ ЭЛЕКТРОНИКИ
СДЕЛАЛА ПРОГРАММУ ДЛЯ АВТОМАТИЧЕСКОЙ УПАКОВКИ
ВИНОГРАДА НА КОНВЕЙЕРЕ, НО СОРТИРОВАТЬ ВИНОГРАД
ОНА НЕ УМЕЕТ.

ВОТ ОНО ЧТО...

МОИ ПРОШЛЫЕ НЕУДАЧИ ПРИВЕЛИ МЕНЯ К ТОМУ,
ЧТО Я РЕШИЛ ПОСОВЕТОВАТЬСЯ С КУДЗЁ-САН.

Столовая

☑ Меню обеда

М-М-М, ТУТ В ЖУРНАЛЕ
НАПИСАНО, ЧТО ПРИ ПОМОЩИ
ГЛУБОКОГО ОБУЧЕНИЯ
МОЖНО РАЗДЕЛЯТЬ
ПРЕДМЕТЫ РАЗНЫХ
РАЗМЕРОВ.



ДАЙ-КА ПОЧИТАТЬ!



ГЛУБОКОЕ ОБУЧЕНИЕ
НА ИЗОБРАЖЕНИЯХ...



ЭТО БЫЛО КАК РАЗ ТО,
ЧТО МНЕ НУЖНО, И Я ПОБЛАГОДАРИЛ
КУАЗЁ-САН ЗА ЭТУ СТАТЬЮ.

НА САМОМ ДЕЛЕ Я ИГРАЛ ТУТ
С КАКИМ-ТО СПЕЦИОМ ПО МАШИННОМУ
ОБУЧЕНИЮ, ПОГОВОРИЛ С НИМ, И ОН
ПОСОВЕТОВАЛ МНЕ ЭТУ КНИГУ.

НО ЭТО... У ТЕБЯ НИЧЕГО
НЕ ПОЛУЧАЕТСЯ, А ТЫ
ВСЕ СТАРАЕШЬСЯ?

АА! Я ЖЕ ХОЧУ ПОМОЧЬ ФЕРМЕРАМ
СОРТИРОВАТЬ ВИНОГРАД.

ДОБРАЯ ДУША.

ХЛЮП

ТЫ ЧТО-ТО
СКАЗАЛ?

НЕ, ЕСЛИ ПОМОЩЬ
С ПРОГРАММОЙ ПОНАДОБИТСЯ,
ЗОВИ!

СПАСИБО!

Я ЧИТАЛА ЭТУ ЗАМЕТКУ.
ХОРОШИЙ ЧЕЛОВЕК
ЭТОТ КУАЗЁ-САН.

Я СНАЧАЛА ПОДУМАЛ, ЧТО
С НИМ ТРУДНО ПОЛАДИТЬ,
НО ОН ВСЕГДА ВЫРУЧИТ.

ХРУМ

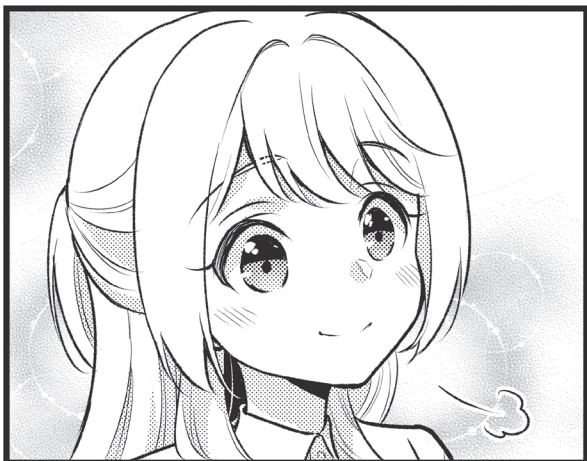
ХРУМ

А ТЫ ХОЧЕШЬ, ЧТОБЫ Я ТЕБЕ
ПОМОГАЛА? НАМИГОЭ-СЭНСЭЙ
УЖЕ ВЕРНУЛСЯ.

НУ, Я ДУМАЛ, САЯКА-СЭМПАЙ,
ВЫ МНЕ ПОМОЖЕТЕ.

Я НЕ ХОЧУ СНОВА
ОБЛАЖАТЬСЯ!

ОТЧАЯННО



ПУСТЬ В ПРОШЛЫЙ РАЗ У ТЕБЯ
НЕ ПОЛУЧИЛОСЬ, НО В ЭТОТ ПОЛУЧИТСЯ!
Я ТЕБЕ ВСЕ РАССКАЖУ!

ДА! СПАСИБО ВАМ!

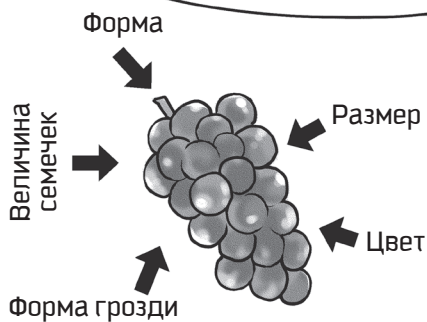


ИТАК, РАЗБЕРЕМ СЕГОДНЯ,
КАК СОЗДАТЬ АВТОМАТИЗИРОВАННУЮ
СИСТЕМУ ПО СОРТИРОВКЕ
ВИНОГРАДА.



КАК МЫ МОЖЕМ ЕГО СОРТИРОВАТЬ?

ПО РАЗМЕРУ, ФОРМЕ,
ЦВЕТУ, ВЕЛИЧИНЕ СЕМЕЧЕК
И ФОРМЕ ГРОЗДИ.



АА, С НАСКОКА С ЭТОЙ ЗАДАЧЕЙ НЕ СПРАВИТЬСЯ. НУЖЕН ОПЫТ!



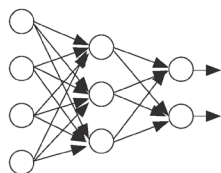
ПО ВСЕЙ ВЕЩНОСТИ, ДЛЯ СОРТИРОВКИ НАМ ПОНАДОБИТСЯ БОЛЬШОЕ КОЛИЧЕСТВО ИЗОБРАЖЕНИЙ И СВЕРТОЧНАЯ НЕЙРОННАЯ СЕТЬ, КОТОРАЯ БУДЕТ ИХ РАСПОЗНАВАТЬ И ЗАТЕМ КЛАССИФИЦИРОВАТЬ.

ЭТО И ЕСТЬ НАШ МЕТОД?

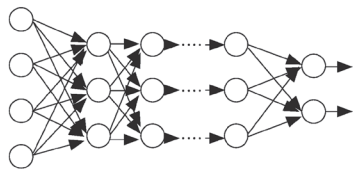


БЕЗ ПАНИКИ! СНАЧАЛА ПОГОВОРИМ О НЕЙРОННЫХ СЕТЯХ, ЗАТЕМ О ГЛУБОКОМ ОБУЧЕНИИ, ИНЫМИ СЛОВАМИ, О ТОМ, ЧТО ЭТО ТАКОЕ.

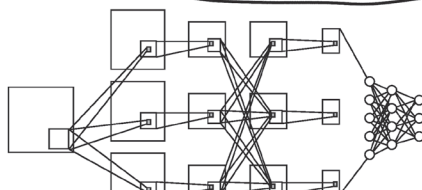
А ПОТОМ УЖЕ ПЕРЕИДЕМ К ОБЪЕКТУ ГЛУБОКОГО ОБУЧЕНИЯ - К СВЕРТОЧНОЙ НЕЙРОННОЙ СЕТИ!



① Основы нейронной сети



② Многослойная нейронная сеть

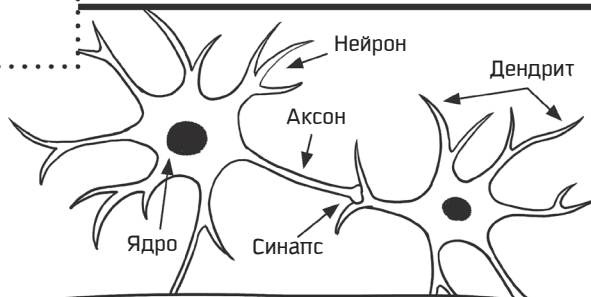


③ Сверточная нейронная сеть

ХОРОШО!

4.1. НЕЙРОННАЯ СЕТЬ

НЕЙРОННАЯ СЕТЬ - ЭТО ОСНОВНОЙ РАСЧЕТНЫЙ МЕХАНИЗМ, СОЗДАННЫЙ ПО МОДЕЛИ НЕРВНЫХ КЛЕТОК ЖИВЫХ СУЩЕСТВ.



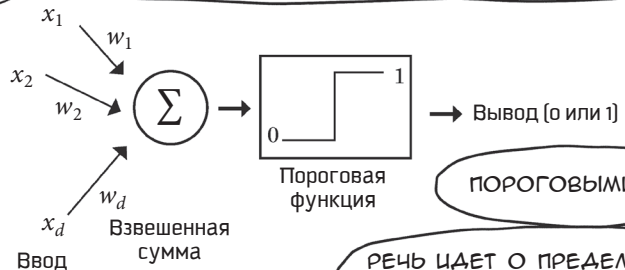
КАК ПОКАЗАНО НА РИСУНКЕ, КЛЕТКИ (НЕЙРОНЫ) СОЕДИНЕНЫ МЕЖАУ СОБОЙ ЧЕРЕЗ СИНАПСЫ - СОЕДИНИТЕЛЬНЫЕ УЧАСТКИ И ОБРАЗУЮТ СЛОЖНУЮ СЕТЬ.

А КАК УСТРОЕНЫ НЕРВНЫЕ КЛЕТКИ ЖИВЫХ СУЩЕСТВ?

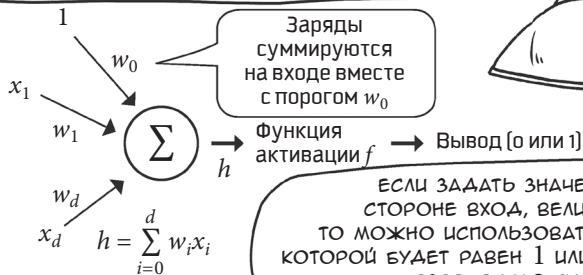
КАЖДЫЙ НЕЙРОН ПОЛУЧАЕТ ПОЛОЖИТЕЛЬНО ИЛИ ОТРИЦАТЕЛЬНО ЗАРЯЖЕННЫЙ СИГНАЛ ОТ СВЯЗАННЫХ С НИМ НЕЙРОНОВ. КОГДА СУММА ЗАРЯДОВ ДОСТИГАЕТ ОПРЕДЕЛЕННОГО ЗНАЧЕНИЯ, ОН НАЧИНАЕТ ГЕНЕРИРОВАТЬ ЭЛЕКТРИЧЕСКИЙ СИГНАЛ.



ТАКИЕ НЕЙРОНЫ, КАК НА КАРТИНКЕ, НАЗЫВАЮТСЯ ПОРОГОВЫМИ ЛОГИЧЕСКИМИ ЭЛЕМЕНТАМИ.



ФУНКЦИЯ $f(h)$ ОТ СУММЫ ЗАРЯДОВ НА ВХОДЕ h НАЗЫВАЕТСЯ ФУНКЦИЕЙ АКТИВАЦИИ.

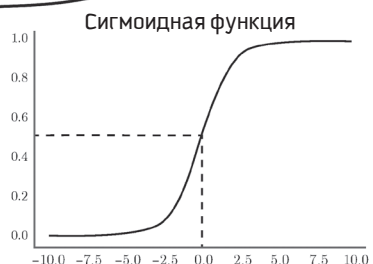
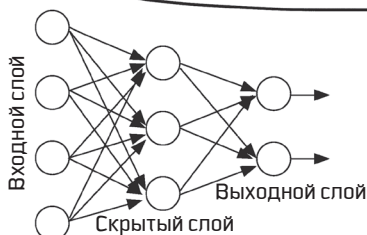


ЕСЛИ ЗАДАТЬ ЗНАЧЕНИЕ ПОРОГА w_0 И ДОБАВИТЬ НА ВХОДНОЙ СТОРОНЕ ВХОД, ВЕЛИЧИНА КОТОРОГО ВСЕГДА БУДЕТ РАВНА w_0 , ТО МОЖНО ИСПОЛЬЗОВАТЬ ПРОСТУЮ ФУНКЦИЮ АКТИВАЦИИ, ВЫХОД КОТОРОЙ БУДЕТ РАВЕН 1 ИЛИ 0 В ЗАВИСИМОСТИ ОТ ТОГО, ЯВЛЯЕТСЯ ЛИ ВЗВЕШЕННАЯ СУММА ПОЛОЖИТЕЛЬНОЙ ИЛИ ОТРИЦАТЕЛЬНОЙ.

СОЕДИНЕНИЕ МНОЖЕСТВА СЛОЕВ ТАКИХ ЭЛЕМЕНТОВ НАЗЫВАЕТСЯ НЕЙРОННОЙ СЕТЬЮ. КОГДА МНОЖЕСТВО СЛОЕВ ЭЛЕМЕНТОВ ОБРАЗУЕТ НЕЙРОННУЮ СЕТЬ, ТО ПОРОГОВАЯ ФУНКЦИЯ, ФУНКЦИЯ АКТИВАЦИИ, ВЕДЕТ СЕБЯ КАК СИГМОИДНАЯ ФУНКЦИЯ. МЫ ЕЕ ВСТРЕЧАЛИ, КОГДА ГОВОРИЛИ О ЛОГИСТИЧЕСКОЙ РЕГРЕССИИ.

А ПОЧЕМУ ИСПОЛЬЗУЕТСЯ СИГМОИДНАЯ ФУНКЦИЯ?

ПРИ ОБУЧЕНИИ ФУНКЦИЮ АКТИВАЦИИ НЕОБХОДИМО ДИФФЕРЕНЦИРОВАТЬ, А ПОРОГОВУЮ ФУНКЦИЮ ДИФФЕРЕНЦИРОВАТЬ НЕЛЬЗЯ.



ЭТО НЕЙРОННАЯ СЕТЬ
С ПРЯМЫМ РАСПРОСТРАНЕНИЕМ.

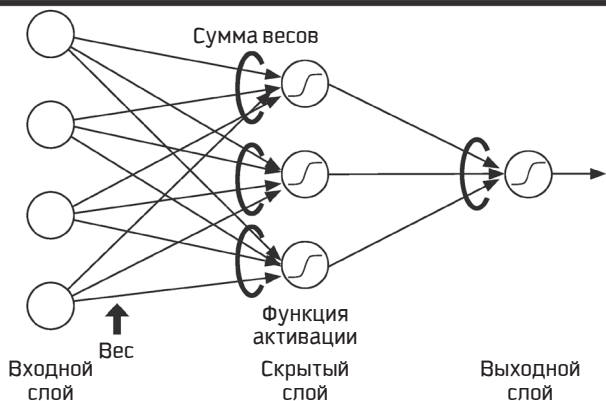
МОДЕЛЬ, СИГНАЛ
В КОТОРОЙ МОЖЕТ
ТОЛЬКО ИДТИ ВПЕРЕД.

ДА.
ВСЕ ЭЛЕМЕНТЫ СОЕДИНЕНЫ
С СОСЕДНИМИ СЛОЯМИ,
НО БЕЗ ОБРАТНОЙ СВЯЗИ.
СИГНАЛ ИДЕТ ОТ ВХОДА
К ВЫХОДУ В ПРЯМОМ
НАПРАВЛЕНИИ.

ЧЕРЕЗ ТРИ СЛОЯ?

ПРИ ЧИСЛЕННЫХ РАСЧЕТАХ
МОЖНО ИСПОЛЬЗОВАТЬ
ТОЛЬКО ДВА, СКРЫТЫЙ СЛОЙ
И ВЫХОДНОЙ.

ЕСЛИ ПРОВЕСТИ АНАЛОГИЮ С ЖИВЫМИ
ОРГАНИЗМАМИ, ТО **ВХОДНОЙ СЛОЙ** -
ЭТО КЛЕТКИ, КОТОРЫЕ РЕАГИРУЮТ
НА ВНЕШНИЕ РАЗДРАЖИТЕЛИ, **СКРЫТЫЙ СЛОЙ** -
КЛЕТКИ, КОТОРЫЕ ПЕРЕДАЮТ СИГНАЛЫ,
А **ВЫХОДНОЙ СЛОЙ** - КЛЕТКИ МОЗГА,
КОТОРЫЕ КЛАССИФИЦИРУЮТ СИГНАЛЫ.



ВХОДНОЙ СИГНАЛ ВЫХОДИТ ИЗ ВХОДНОГО СЛОЯ
БЕЗ ИЗМЕНЕНИЙ, ЕГО ЗАРЯДЫ СУММИРУЮТСЯ
И ПОСТУПАЮТ НА СКРЫТЫЙ СЛОЙ. НА СКРЫТОМ
СЛОЕ СУММА ЗАРЯДОВ, ПОСТУПИВШИХ
СО ВСЕХ ВХОДНЫХ СЛОЕВ, ОБРАБАТЫВАЕТСЯ
С ПОМОЩЬЮ ФУНКЦИИ АКТИВАЦИИ. В СЛУЧАЕ
ЕСЛИ В НЕЙРОННОЙ СЕТИ С ПРЯМЫМ
РАСПРОСТРАНЕНИЕМ РЕШАЕТСЯ ЗАДАЧА
БИНАРНОЙ КЛАССИФИКАЦИИ, НА ВЫХОДНОМ
СЛОЕ ОДИН ЭЛЕМЕНТ.

КАК И В СЛУЧАЕ ЛОГИСТИЧЕСКОЙ РЕГРЕССИИ,
ВЫХОДНОЕ ЗНАЧЕНИЕ ЭТОГО СЛОЯ МОЖНО
РАССМАТРИВАТЬ КАК ВЕРОЯТНОСТЬ ТОГО,
ЧТО ВХОДНОЙ СИГНАЛ ОТНОСИТСЯ
К ПОЛОЖИТЕЛЬНОМУ КЛАССУ.

ЕСЛИ ЖЕ КЛАССИФИКАЦИЯ ВЕДЕТСЯ ПО НЕСКОЛЬКИМ КАТЕГОРИЯМ, ТО КОЛИЧЕСТВО ВЫХОДНЫХ ЭЛЕМЕНТОВ СОВПАДАЕТ С КОЛИЧЕСТВОМ КЛАССОВ. В ТАКОМ СЛУЧАЕ ВЕРОЯТНО, ЧТО БОЛЬШОЕ КОЛИЧЕСТВО ЭЛЕМЕНТОВ НА ВЫХОДНОМ СЛОЕ ВЫДАСТ ЗНАЧЕНИЕ, БЛИЗКОЕ К 1.

ТОГДА В КАЧЕСТВЕ ФУНКЦИИ АКТИВАЦИИ $f(h)$ ИСПОЛЬЗУЕТСЯ НЕ СИГМОИДА, А ФУНКЦИЯ SOFTMAX.

$$g_k = \frac{\exp(h_k)}{\sum_{j=1}^c \exp(h_j)}$$

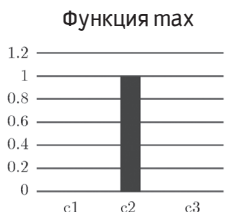


ЗДЕСЬ h_k - СУММА ЗАРЯДОВ НА ВЫХОДЕ ИЗ СКРЫТОГО СЛОЯ, КОТОРАЯ СООТВЕТСТВУЕТ КОЛИЧЕСТВУ ЭЛЕМЕНТОВ ВЫХОДНОГО СЛОЯ И КОЛИЧЕСТВУ КЛАССОВ k .

ЕСЛИ В КАЧЕСТВЕ ФУНКЦИИ АКТИВАЦИИ ИСПОЛЬЗУЕТСЯ SOFTMAX, ТО ПРИ СЛОЖЕНИИ ВЫХОДОВ ВСЕХ ЭЛЕМЕНТОВ ВЫХОДНОГО СЛОЯ g_k ПОЛУЧИТСЯ 1, И ПОЭТОМУ МОЖНО ГОВОРИТЬ О ВЕРОЯТНОСТИ.

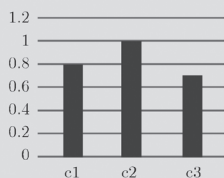
КАК ВЫ ДУМАЕТЕ, МОЖНО ПРИМЕНИТЬ СИГМОИДНУЮ ФУНКЦИЮ К ВЗВЕШЕННОЙ СУММЕ ЗАРЯДОВ, КОТОРАЯ ПРИНИМАЕТ ЗНАЧЕНИЯ ОТ 0 ДО 1, А ЗАТЕМ ВЗЯТЬ МАКСИМАЛЬНОЕ ЗНАЧЕНИЕ?

И ПОНЯТЬ ВЕРОЯТНОСТЬ ИСПОЛЬЗОВАНИЯ КАЖДОГО КЛАССА.

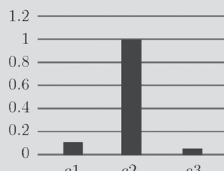


Ясно только, какой класс используется чаще всего

Величина взвешенной суммы зарядов на входе выходного слоя нейронной сети

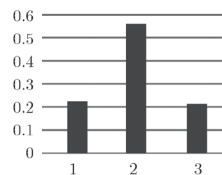
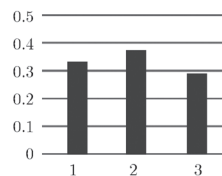


Случай 1: классы более-менее равны



Случай 2: один класс выделяется

Функция softmax



Ясна вероятность появления каждого класса

4.2. ОБУЧЕНИЕ МЕТОДОМ ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШИБОК

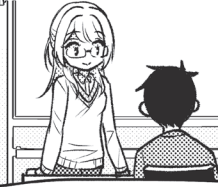
КОГДА МЫ ГОВОРИМ ОБ "ОБУЧЕНИИ" В НЕЙРОННЫХ СЕТЯХ С ПРЯМЫМ РАСПРОСТРАНЕНИЕМ, МЫ ГОВОРИМ

О "ВЗВЕШЕННОЙ СУММЕ ЗАРЯДОВ ВХОДНОГО СИГНАЛА".

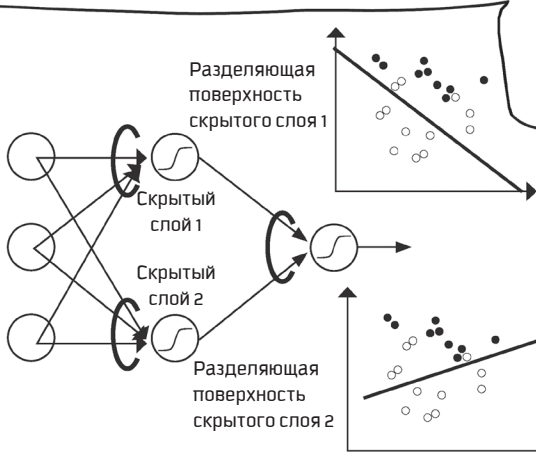


ПО ПОЛУЧЕННЫМ ДАННЫМ МЫ ОПРЕДЕЛЯЕМ И РЕГУЛИРУЕМ СУММУ ЗАРЯДОВ КАЖДОГО ЭЛЕМЕНТА.

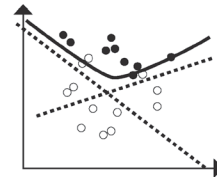
КАК ЭТО ПРОИСХОДИТ?



КАЖДАЯ ЭЛЕМЕНТ ВЫПОЛНЯЕТ НЕЛИНЕЙНОЕ ПРЕОБРАЗОВАНИЕ С ПОМОЩЬЮ СИГМОИДНОЙ ФУНКЦИИ. ЕСЛИ СЛОЖИТЬ СУММЫ ЗАРЯДОВ (ВЕСОВ), ТО ПОЛУЧИТСЯ НЕЛИНЕЙНЫЙ КЛАССИФИКАТОР ПРИЗНАКОВОГО ПРОСТРАНСТВА.



ЧТОБЫ УМЕНЬШИТЬ КОЛИЧЕСТВО ОШИБОК КЛАССИФИКАЦИИ В ЭТОМ НЕЛИНЕЙНОМ КЛАССИФИКАТОРЕ, КАК И В СЛУЧАЕ КЛАССИФИКАЦИИ, КОТОРУЮ МЫ ДЕЛАЛИ ПЕРЕД НАСТРОЙКОЙ ВСЕГО ВЕСА, НАДО НАЙТИ РАЗДЕЛЯЮЩУЮ ПОВЕРХНОСТЬ И МИНИМИЗИРОВАТЬ ОШЕЧКУ.



Разделяющая поверхность, представляющая собой взвешенную сумму двух разделяющих поверхностей

НАСТРОЙКА ВЕСА ОТ СКРЫТОГО СЛОЯ К ВЫХОДНОМУ ПРОИЗВОДИТСЯ С ПОМОЩЬЮ ВЫХОДА НЕЙРОННОЙ СЕТИ И **ОБУЧАЮЩЕГО СИГНАЛА**. ИНЫМИ СЛОВАМИ, ОБУЧЕНИЕ ВЕДЕТСЯ ПУТЕМ СРАВНЕНИЯ С ПРАВИЛЬНЫМ РЕЗУЛЬТАТОМ И ОЦЕНКИ ОШЕЧКИ.

А ЕСЛИ ОБУЧАЮЩЕГО СИГНАЛА НА СКРЫТОМ СЛОЕ НЕТ, ТО И ОШЕЧКИ, ПОЛУЧАЕТСЯ, СРАВНИВАТЬ НЕ С ЧЕМ.

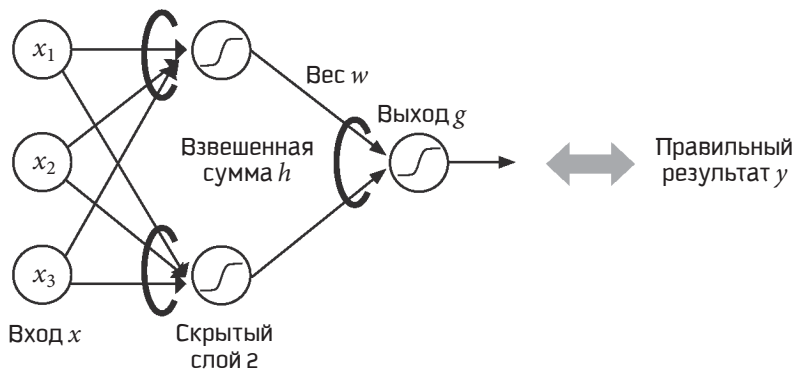


ИМЕННО! В ТАКОМ СЛУЧАЕ В МНОГОСЛОЙНОЙ СЕТИ ЦЕЛЕТ ОБУЧЕНИЕ ПУТЕМ **МЕТОДА ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШЕБОК**.





Поговорим о методе обратного распространения ошибок. Пусть обучение ведется в нейронной сети с архитектурой, как на рисунке ниже:



В качестве данных для обучения возьмем пары: признаковое описание объекта x и результат y . Обозначим набор данных D , и в нем есть i -я пара (x_i, y_i) . Функция ошибки может быть определена по-разному, но в данном случае будем минимизировать ее, вычисляя квадрат ошибки.

$$E(\mathbf{w}) \equiv \frac{1}{2} \sum_{x_i \in D} (g_i - y_i)^2. \quad (4.1)$$

Здесь \mathbf{w} обозначает всю сумму весов нейронной сети. Используя описанный в главе 2 метод градиентного спуска, найдем один вес w из весов \mathbf{w} и, регулируя его величину, снизим ошибку.

$$w \leftarrow w - \eta \frac{\partial E(\mathbf{w})}{\partial w}. \quad (4.2)$$

Далее найдем частную производную по w функции ошибки $E(\mathbf{w})$, используя метод градиентного спуска. В этом случае изменение веса w изменит и выход g функции активации. Используя формулу дифференцирования сложной функции, получим следующее выражение:

$$\frac{\partial E(\mathbf{w})}{\partial w} = \frac{\partial E(\mathbf{w})}{\partial h} \frac{\partial h}{\partial w}. \quad (4.3)$$

После вычисления второго множителя в правой части формулы (4.3) из определения суммы весов h можно получить значение выхода предыдущего слоя в сочетании с весом w . Первый множитель можно записать, как показано ниже, с использованием производной. После вычислений у нас получится величина ошибки ε .

$$\varepsilon = \frac{\partial E(w)}{\partial h} = \frac{\partial E(w)}{\partial g} \frac{\partial g}{\partial h}. \quad (4.4)$$

В формуле 4.4 второй множитель справа – это производная функции активации. Поскольку в качестве функции активации используется сигмоида, то это $g(1 - g)$. Значение первого множителя справа зависит от того, между какими слоями вес определяется. Если w – это вес от скрытого слоя к выходному, то первый множитель – это производная функции ошибки.

$$\frac{\partial E(w)}{\partial g} = g - y. \quad (4.5)$$

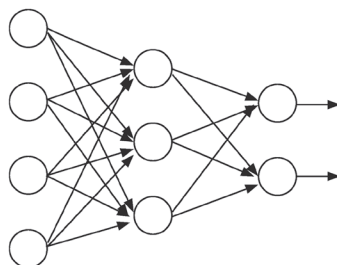
С другой стороны, если w – это вес от входного слоя к скрытому, то g в первом множителе – это выход скрытого слоя, и его величина будет влиять на выход выходного слоя h . Обычно выходных слоев много, и их обозначают h_j , где j – номер слоя, и первый множитель уравнения (4.4) примет такой вид:

$$\frac{\partial E(w)}{\partial g} = \sum_j \frac{\partial E(w)}{\partial h_j} \frac{\partial h_j}{\partial g} = \sum_j \varepsilon_j w_j. \quad (4.5)$$

Мы используем ε из формулы (4.4), ε_j – это ошибка от скрытого слоя j , суть метода в том, что эта величина используется для корректировки суммы весов от входного слоя к скрытому. После обобщения ε можно вычислить и по формулам, приведенным ниже:

$$\varepsilon = \begin{cases} (g - y)g(1 - g) & \text{в случае от скрытого слоя к выходному} \\ \sum_j \varepsilon_j w_j g(1 - g) & \text{в случае от начального слоя к скрытому} \end{cases}$$

Расчет во время обучения
(обратное направление)



Входной слой Скрытый слой Выходной слой

Вычисление веса нейронов
(прямое направление)

СТРЕЛКА ПОД КАРТИНКОЙ ПОКАЗЫВАЕТ ПРЯМОЕ РАСПРОСТРАНЕНИЕ СИГНАЛА ПО СЕТИ, ВЫЧИСЛЕНИЕ ВЕСА НЕЙРОНОВ. НАПРАВЛЕНИЕ ВО ВРЕМЯ ОБУЧЕНИЯ УКАЗАНО ВЕРХНЕЙ СТРЕЛКОЙ.

ТО ЕСТЬ ОБУЧЕНИЕ ВЕДЕТСЯ ОТ ВЫХОДНОГО СЛОЯ?

МОЖНО ПРЕДСТАВИТЬ, ЧТО ВЫХОДНОЙ СЛОЙ ОТПРАВЛЯЕТ С ВЫХОДА НАЗАД СИГНАЛ ОБУЧЕНИЯ, И ЕСЛИ ОШИБКА ВЕЛИКА, ТО ОН ЗАЛТИСЯ.



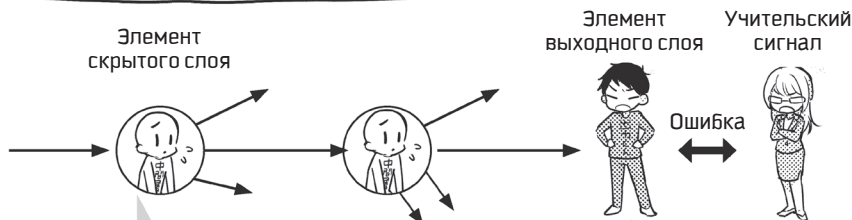
РАЗ ОДНИ ОШИБКИ, ТО ОН, КОНЕЧНО, БУДЕТ СЕРДИТ.

ЗАТЕМ РАЗОЗЛЕННЫЕ ВЫХОДНЫМ СЛОЕМ СИГНАЛЫ КОРРЕКТИРУЮТ ВЕЛИЧИНУ ВЕСОВ НА СРЕДНЕМ СЛОЕ И РУГАЮТ ЕГО.



КАК В ФИРМЕ....

МОЖНО ТАК ПРЕДСТАВИТЬ СЕБЕ ОБУЧЕНИЕ МЕТОДОМ ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШИБКИ:



Вес разозленного элемента \times величина веса, направленного к элементу

Сравнение ошибок злит учительский сигнал

4.3. ВЫЗОВЫ ГЛУБОКОГО ОБУЧЕНИЯ

ТАК, А ТЕПЕРЬ ПОГОВОРИМ О ГЛУБОКОМ ОБУЧЕНИИ (DEEP LEARNING) КАК О МНОГОСЛОЙНОЙ НЕЙРОННОЙ СЕТИ.



ДА, УЧИТЕЛЬ!
МЫ БУДЕМ ГОВОРИТЬ О ТОМ, ЧТО ТАИТСЯ
В ГЛУБИНЕ НЕЙРОННЫХ СЕТЕЙ?

НУ... ЦЕЛЬЮ
ОБУЧЕНИЯ ЯВЛЯЕТСЯ ВЫДЕЛЕНИЕ
ПРИЗНАКОВ. ДЛЯ ЭТОГО НАДО ПРОВОДИТЬ
ВЫДЕЛЕНИЕ ПРИЗНАКОВ НА ОСНОВАНИИ
СЛОЖНЫХ ПРОЦЕДУР РАСПОЗНАВАНИЯ ЗВУКА
ИЛИ ИЗОБРАЖЕНИЯ. ОДНАКО ГЛУБОКИЕ
НЕЙРОННЫЕ СЕТИ МОГУТ ОБРАБАТЫВАТЬ...

...ПОСТУПИВШИЕ ДАННЫЕ ИЗОБРАЖЕНИЙ
ИЛИ ЗВУКОВЫЕ СИГНАЛЫ ПРИ ПОМОЩИ
ПРОСТЫХ ВЫРАЖЕНИЙ ПРИЗНАКОВ,
КОТОРЫЕ ОБРАЗУЮТ СЛОЖНЫЕ
ПРОЦЕДУРЫ. ЭТО ОЧЕНЬ
ВЫСОКОЭФФЕКТИВНЫЙ МЕТОД.

С чем легко справляются глубокие
нейронные сети:

- распознавание звука;
- распознавание изображений;
- распознавание естественной речи.



ОГО...

НА САМОМ ДЕЛЕ ОБУЧЕНИЕ
МЕТОДОМ ОБРАТНОГО
РАСПРОСТРАНЕНИЯ ОШИБКИ ВОШЛО
В МОДУ ТОЛЬКО ВО ВТОРОЙ ПОЛОВИНЕ
1980-Х, ДО ЭТОГО ЭФФЕКТИВНОСТЬ
НЕЙРОННЫХ СЕТЕЙ
ПЫТАЛИСЬ УЛУЧШИТЬ, НО
НИЧЕГО НЕ ПОЛУЧАЛОСЬ.



ПОЧЕМУ?

ПОГОВОРИМ ОБ ЭТОМ. СНАЧАЛА
РАССКАЖУ О ПРОБЛЕМЕ ГЛУБОКИХ
НЕЙРОННЫХ СЕТЕЙ. А ПОТОМ О ДВУХ
ЕЕ РЕШЕНИЯХ:

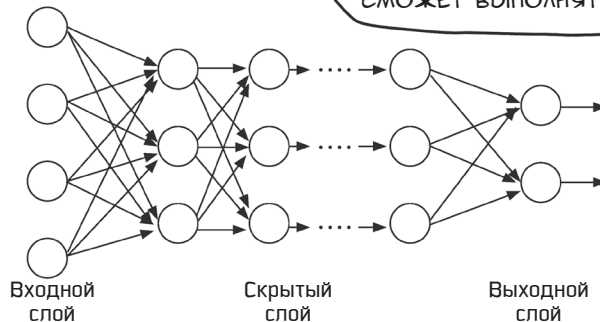
- 1) ХИТРОСТЯХ, КОТОРЫЕ ИСПОЛЬЗУЮТ
В МНОГОУРОВНЕВОМ ОБУЧЕНИИ;
- 2) ИСПОЛЬЗОВАНИИ СПЕЦИАЛИЗИ-
РОВАННЫХ НЕЙРОННЫХ СЕТЕЙ.

ХОРОШО!



4.3.1. Проблема глубокой нейронной сети

ТАК, ДЛЯ НАЧАЛА - ПРОБЛЕМА...



КАЗАЛОСЬ БЫ, ЧТО ЕСЛИ УВЕЛИЧИТЬ КОЛИЧЕСТВО СКРЫТЫХ СЛОЕВ В НЕЙРОННОЙ СЕТИ С ПРЯМЫМ ОБУЧЕНИЕМ, КАК НА КАРТИНКЕ, ТО ЕЕ ЭФФЕКТИВНОСТЬ ТОЖЕ ВЫРАСТЕТ, И ОНА СМОЖЕТ ВЫПОЛНЯТЬ РАЗНЫЕ ЗАДАЧИ.

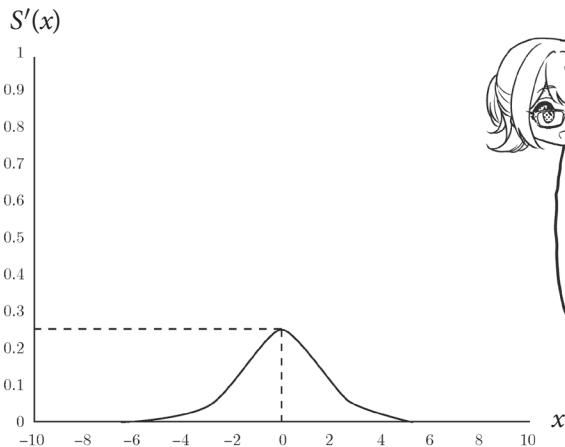
ОДНАКО В ОБУЧЕНИИ МНОГОСЛОЙНОЙ НЕЙРОННОЙ СЕТИ МЕТОДОМ ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШИБОК ПРИ ВОЗВРАЩЕНИИ К КОРРЕКТИРОВКЕ ВЕЛИЧИНЫ ВЕСА ВОЗНИКАЕТ

ПРОБЛЕМА ИСЧЕЗАЮЩЕГО ГРАДИЕНТА...

ПРОБЛЕМА ИСЧЕЗАЮЩЕГО ГРАДИЕНТА?

ФОРМУЛА ДЛЯ ВЫЧИСЛЕНИЯ ИЗМЕНЕНИЯ ВЕЛИЧИНЫ ОШИБКИ НА КАЖДОМ ЭТАПЕ ПРЕДСТАВЛЯЕТ СОБОЙ ПРОИЗВОДНУЮ СИГМОИДНОЙ ФУНКЦИИ:
$$S(x) = S(x)(1 - S(x)).$$

ПОСМОТРИ НА ГРАФИК, КАК ОНА СЕБЯ ВЕДЕТ В РЕАЛЬНОЙ ЖИЗНИ.



МАКСИМАЛЬНЫЙ ГРАДИЕНТ - 0,25.

ДА.

ЭТО СРАВНИТЕЛЬНО ВЫСОКОЕ ЗНАЧЕНИЕ ПОЯВЛЯЕТСЯ, ТОЛЬКО КОГДА ВХОДНЫЕ СИГНАЛЫ РАВНЫ 0 ИЛИ БЛИЗКИ К НЕМУ. А В ОСТАЛЬНЫХ СЛУЧАЯХ ГРАДИЕНТ БЛИЗИТСЯ К 0.

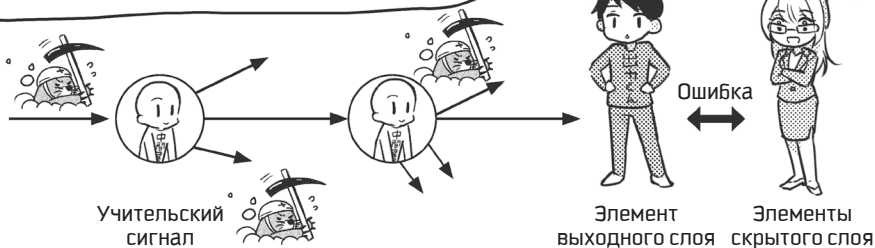
ОГО, ТОГДА ОБУЧЕНИЕ НЕ ПРОАВНЕТСЯ СЛИШКОМ ДАЛЕКО ОТ ВХОДНОГО СЛОЯ.



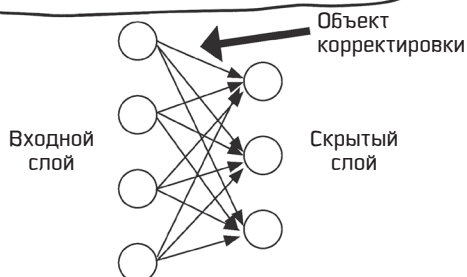
4.3.2. Хитрости многоступенчатого обучения

1. МЕТОД ПРЕДВАРИТЕЛЬНОГО ОБУЧЕНИЯ

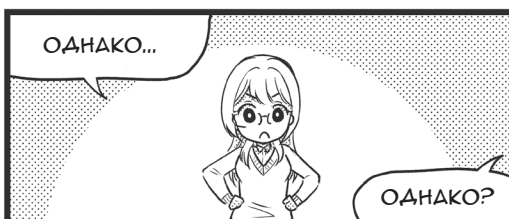
ПРИ ИСПОЛЬЗОВАНИИ МЕТОДА ПРЕДВАРИТЕЛЬНОГО ОБУЧЕНИЯ МЫ ОПРЕДЕЛЯЕМ НУЖНЫЕ ПАРАМЕТРЫ ВЕСОВ В САМОМ НАЧАЛЕ, ПЕРЕД ОБУЧЕНИЕМ, МЕТОДОМ ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШИБКИ.



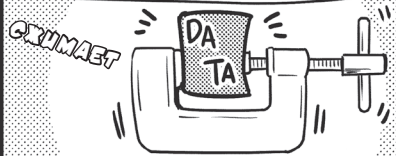
ТАКИМ ОБРАЗОМ, ВХОДНОЙ СЛОЙ ЗАРАНЕЕ РЕГУЛИРУЕТ ВЕС САМОГО БЛИЗКОГО СКРЫТОГО СЛОЯ.



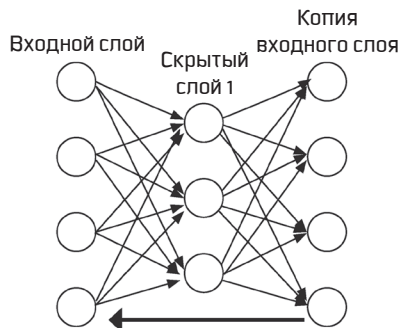
ЗДЕСЬ ВЕСА НЕОБХОДИМЫ ДЛЯ ТОГО, ЧТОБЫ КОНВЕРТИРОВАТЬ ДАННЫЕ ВВОДА В ПРОСТРАНСТВО МЕНЬШЕЙ РАЗМЕРНОСТИ, ЧТОБЫ ИХ МОЖНО БЫЛО КЛАССИФИЦИРОВАТЬ.



ТОГДА ПРЕДСТАВЬ, ЧТО ПРОБЛЕМА В ТОМ, ЧТО "НАДО СЖАТЬ ПРИЗНАКОВОЕ ОПИСАНИЕ ОБЪЕКТА, ЧТОБЫ ОНО УЛОЖИЛОСЬ В НЕБОЛЬШОЕ КОЛИЧЕСТВО ЕДИНИЦ И С МИНИМАЛЬНЫМИ ПОТЕРЯМИ, НАСКОЛЬКО ЭТО ВОЗМОЖНО".



ЭТО КАК?

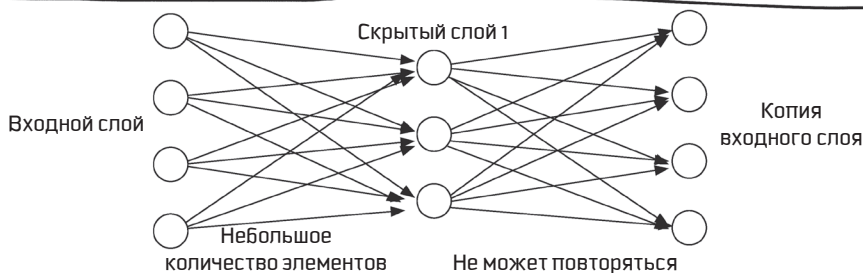


Обучение методом обратного распространения ошибок

ПРЕЖДЕ ВСЕГО СОЗДАЕТСЯ КОПИЯ ЭЛЕМЕНТОВ ВХОДНОГО СЛОЯ НАД СКРЫТЫМ СЛОЕМ, И ОНА СТАНОВИТСЯ ВЫХОДНЫМ СЛОЕМ. ЗАТЕМ ИНФОРМАЦИЯ ИЗ ВХОДНОГО СЛОЯ СНОВА ПОСТУПАЕТ НА ВЫХОДНОЙ И ПРОХОДИТ ОБУЧЕНИЕ. ЭТО ТАК НАЗЫВАЕМЫЙ АВТОКОДИРОВЩИК.

ТО ЕСТЬ КОПИЯ ВХОДНОГО СЛОЯ СТАНОВИТСЯ ВЫХОДНЫМ СЛОЕМ? ОНА ТАКАЯ ЖЕ, КАК И ВХОДНОЙ СЛОЙ?

В ЦЕЛОМ КОЛИЧЕСТВО ЭЛЕМЕНТОВ СКРЫТОГО СЛОЯ ПО СРАВНЕНИЮ С КОЛИЧЕСТВОМ ЭЛЕМЕНТОВ ВХОДНОГО СЛОЯ НЕВЕЛИКО, И ПОЭТОМУ ОНИ КОПИРУЮТ ИНФОРМАЦИЮ ВХОДНОГО СЛОЯ, ЧТОБЫ ВЫХОДНОЙ СЛОЙ НЕ ПОВТОРЯЛСЯ.

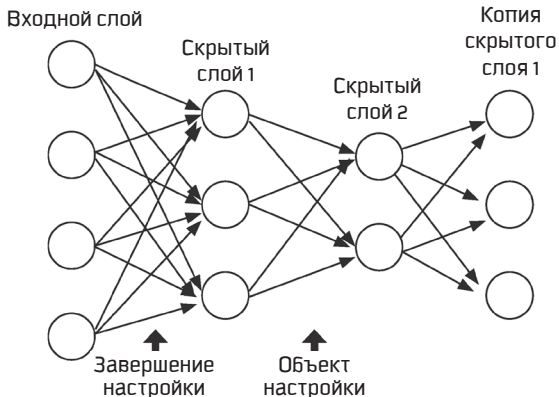


ВОТ КАК.

БЕЗ ЭТОГО НЕЛЬЗЯ ИЗВЛЕЧЬ ПРИЗНАКИ ДАННЫХ?

СКРЫТЫЙ СЛОЙ ПОЛУЧАЕТ ЗАДАЧУ "ИЗВЛЕЧЬ ИНФОРМАЦИЮ ИЗ ПРОСТРАНСТВА ДАННЫХ, СЖАТОГО ДО БОЛЕЕ НИЗКОГО ИЗМЕРЕНИЯ".

ПОСЛЕ ТОГО КАК НАСТРОЕНЫ ВЕСА МЕЖДУ ВХОДНЫМ СЛОЕМ И СКРЫТЫМ СЛОЕМ 1, ОНИ ФИКСИРУЮТСЯ, ЗАТЕМ ОБУЧЕНИЕ ВЕДЕТСЯ МЕЖДУ СКРЫТЫМ СЛОЕМ 1 И СКРЫТЫМ СЛОЕМ 2, И ТАК ПОВТОРЯЕТСЯ ВПЛОТЬ ДО ВЫХОДНОГО СЛОЯ.



Извлечем
важные
признаки!



ЕСЛИ НЕ ИЗВЛЕКАТЬ ПРИЗНАКИ, ТО, ПОСКОЛЬКУ КОЛИЧЕСТВО ЭЛЕМЕНТОВ (УЗЛОВ) УМЕНЬШАЕТСЯ ПО МЕРЕ УДАЛЕНИЯ ОТ ВХОДНОГО СЛОЯ, ИНФОРМАЦИЯ НЕ СОХРАНЯЕТСЯ.

МЕТОД ПРЕДВАРИТЕЛЬНОГО ОБУЧЕНИЯ ПОЗВОЛИЛ РЕШИТЬ ПРОБЛЕМУ ИСЧЕЗАЮЩЕГО ГРАДИЕНТА И ПРОДОЛЖАТЬ ИЗВЛЕКАТЬ АБСТРАКТНЫЕ ИНФОРМАЦИОННЫЕ ВЫРАЖЕНИЯ, СОХРАНЯЯ ИНФОРМАЦИЮ.

**МЕТОД
ПРЕДВАРИТЕЛЬНОГО
ОБУЧЕНИЯ**

ПОСЛЕ РАЗРАБОТКИ ЭТОГО МЕТОДА В 2006 ГОДУ СТАЛА ВОЗМОЖНОЙ РАЗРАБОТКА ГЛУБОКИХ НЕЙРОННЫХ СЕТЕЙ.

ДЕЙСТВИТЕЛЬНО, ПРОРЫВ...

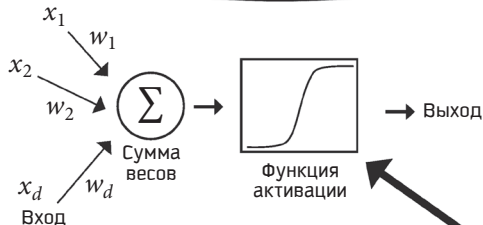


4.3.3. Хитрости многоступенчатого обучения

2. Функция активации

ДРУГОЙ ПОДХОД К ПРОБЛЕМЕ ИСЧЕЗАЮЩЕГО ГРАДИЕНТА - ЭТО МЕТОД НАСТРОЙКИ ФУНКЦИИ АКТИВАЦИИ.

НАДО ОБРАТИТЬ ВНИМАНИЕ, КАК ФУНКЦИЯ АКТИВАЦИИ ОБРАЩАЕТСЯ С СУММОЙ ВЕСОВ.



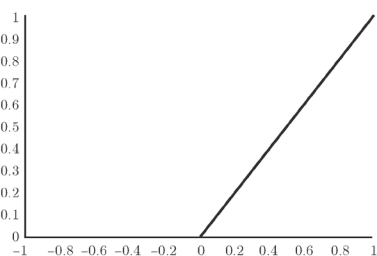
Хочу поправить
функцию активации!



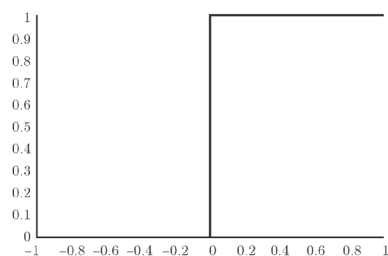
ВМЕСТО СИГМОИДЫ, О КОТОРОЙ МЫ ГОВОРИЛИ, ИСПОЛЬЗУЕТСЯ
УСЕЧЕННАЯ ЛИНЕЙНАЯ ФУНКЦИЯ, ГДЕ $f(x) = \max(0, x)$.
ЭЛЕМЕНТ, КОТОРЫЙ ИСПОЛЬЗУЕТ ЭТУ ФУНКЦИЮ,
НАЗЫВАЕТСЯ **БЛОКОМ ЛИНЕЙНОЙ РЕКТИФИКАЦИИ (RELU)**.



УСЕЧЕННАЯ ЛИНЕЙНАЯ ФУНКЦИЯ - ЭТО ЛИНЕЙНАЯ ФУНКЦИЯ, КОТОРАЯ
ВОЗВРАЩАЕТ ЗНАЧЕНИЕ x , ЕСЛИ x ПОЛОЖИТЕЛЬНО, И 0 В ПРОТИВНОМ СЛУЧАЕ.



Усеченная линейная функция



Производная усеченной линейной функции

В ОТЛИЧИЕ ОТ СИГМОИДА,
ГРАДИЕНТ РАВЕН 1.

ПРИ ИСПОЛЬЗОВАНИИ БЛОКА ЛИНЕЙНОЙ РЕКТИФИКАЦИИ
ГРАДИЕНТ РАВЕН 1, И ОШИБКА НЕ ИСЧЕЗАЕТ.
ОДНАКО ВЫХОД МНОГИХ ЭЛЕМЕНТОВ РАВЕН 0, НЕЙРОННАЯ СЕТЬ
МОЖЕТ С БОЛЬШОЙ СКОРОСТЬЮ СЧИТАТЬ ГРАДИЕНТЫ,
И ОБУЧЕНИЕ МОЖНО ВЕСТИ
БЕЗ ПРЕДВАРИТЕЛЬНОГО ОБУЧЕНИЯ.

КРУТО, ЧТО МОЖНО ОБОЙТИСЬ
БЕЗ ПРЕДВАРИТЕЛЬНОГО ОБУЧЕНИЯ.



4.3.4. Хитрости многоступенчатого обучения

3. КАК ИЗБЕЖАТЬ ПЕРЕОБУЧЕНИЯ

В ГЛУБОКОМ ОБУЧЕНИИ, КРОМЕ ПРОБЛЕМЫ ИСЧЕЗАЮЩЕГО ГРАДИЕНТА, ЕСТЬ ЕЩЕ И ПРОБЛЕМА ПЕРЕОБУЧЕНИЯ.



ЕСЛИ В МОДЕЛИ ЕСТЬ ПАРАМЕТРЫ С БОЛЬШИМИ ВЕСАМИ, ОНИ МОГУТ СЛИШКОМ ХОРОШО ОБЪЯСНЯТЬ ПРИМЕРЫ В ДАННЫХ ДЛЯ ОБУЧЕНИЯ.

ЧТОБЫ ИЗБЕЖАТЬ ПЕРЕОБУЧЕНИЯ, МОЖНО ИСПОЛЬЗОВАТЬ МЕТОД ПРОРЕЖИВАНИЯ (DROPOUT), ПРИ КОТОРОМ ПЕРЕОБУЧЕНИЕ ПРОИСХОДИТ РЕЖЕ, А УНИВЕРСАЛЬНОСТЬ СИСТЕМЫ ПОВЫШАЕТСЯ.

А КАК РАБОТАЕТ ЭТОТ МЕТОД?

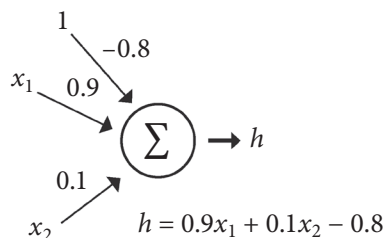


ОБУЧЕНИЕ ВЕДЕТСЯ ПУТЕМ ИСКЛЮЧЕНИЯ СЛУЧАЙНЫХ НЕЙРОНОВ.

ОБУЧЕНИЕ ПУТЕМ ИСКЛЮЧЕНИЯ СЛУЧАЙНЫХ ЭЛЕМЕНТОВ?

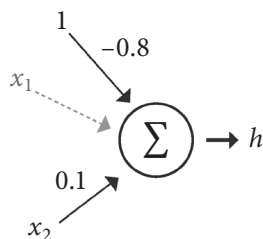
ПЕРВЫМ ДЕЛОМ В КАЖДОМ СЛОЕ СЛУЧАЙНЫМ ОБРАЗОМ УБИРАЮТСЯ ЭЛЕМЕНТЫ В ЗАВИСИМОСТИ ОТ p .

НАПРИМЕР, ПРИ $p = 0,5$ В НЕЙРОННОЙ СЕТИ РАБОТАЕТ ПОЛОВИНА ЭЛЕМЕНТОВ.



Хотя оба элемента – и x_1 , и x_2 – важны, небольшая разница во время обучения приводит к большой разнице в весах в окончательном результате. Если величина x_1 в неизвестных данных будет даже ненамного меньше, это может привести к ошибочной идентификации.

Исключение

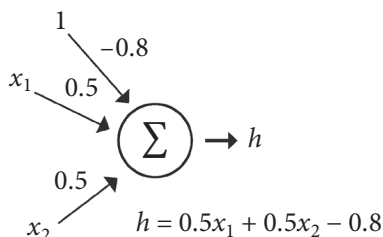


Поскольку любой входной элемент во время обучения может исчезнуть с определенной вероятностью, значения весов выровнены так, что только один из них даст правильный ответ.

Значения на входе в неизвестных данных могут немного меняться, но это не важно.

А ЗАТЕМ ОБУЧЕНИЕ ВЕДЕТСЯ МЕТОДОМ ОБРАТНОГО РАСПРОСТРАНЕНИЯ ОШИБОК НА ОДНОМ "МИНИ-ПАКЕТЕ" ДАННЫХ.

КОГДА ЛОГИСТИЧЕСКАЯ КЛАССИФИКАЦИЯ ВЕДЕТСЯ УЖЕ ОБУЧЕННОЙ НЕЙРОННОЙ СЕТЬЮ, ВЕС УВЕЛИЧИВАЕТСЯ В p РАЗ. ПРОИЗВОДИТСЯ ОБУЧЕНИЕ НЕСКОЛЬКИХ БЛОКОВ НЕЙРОННОЙ СЕТИ, РЕЗУЛЬТАТЫ УСРЕДНЯЮТСЯ.



НО ПОЧЕМУ ПРИ ИСПОЛЬЗОВАНИИ МЕТОДА ПРОРЕЖИВАНИЯ НЕ ПРОИСХОДИТ ПЕРЕОБУЧЕНИЯ?



СУЩЕСТВУЕТ ТЕОРИЯ, ЧТО, ВО-ПЕРВЫХ, БЛАГОДАРЯ СНИЖЕНИЮ СТЕПЕНИ СВОБОДЫ ПРОИСХОДИТ РЕГУЛЯРИЗАЦИЯ, А ВО-ВТОРЫХ, ГРАДИЕНТ ОСТАЕТСЯ ТАКИМ ЖЕ, ДАЖЕ НЕСМОТРЯ НА ТО, ЧТО РАСПРЕДЕЛЕНИЕ ВЗВЕШЕННОЙ СУММЫ ВЕСОВ h НА ВХОДЕ НА УЗЕЛ РАСТЕТ.

КОНЕЧНО, ДО СИХ ПОР СУЩЕСТВУЮТ РАЗНОГЛАСИЯ СРЕДИ УЧЕНЫХ, НО, КАК БЫ ТО НИ БЫЛО, ЕСЛИ ИЗБЕЖАТЬ МНОГОКРАТНОГО ВВОДА ОДНИХ И ТЕХ ЖЕ ДАННЫХ В НЕЙРОННЫХ СЕТЯХ С ТАКОЙ ЖЕ СТРУКТУРОЙ, ТО МОЖНО ИЗБЕЖАТЬ ПЕРЕОБУЧЕНИЯ.

НЕЛЬЗЯ, ЧТОБЫ ГЛУБОКАЯ СЕТЬ ВСЕ ЗАПОМИНАЛА!

4.3.5. Нейронные сети со специализированной структурой

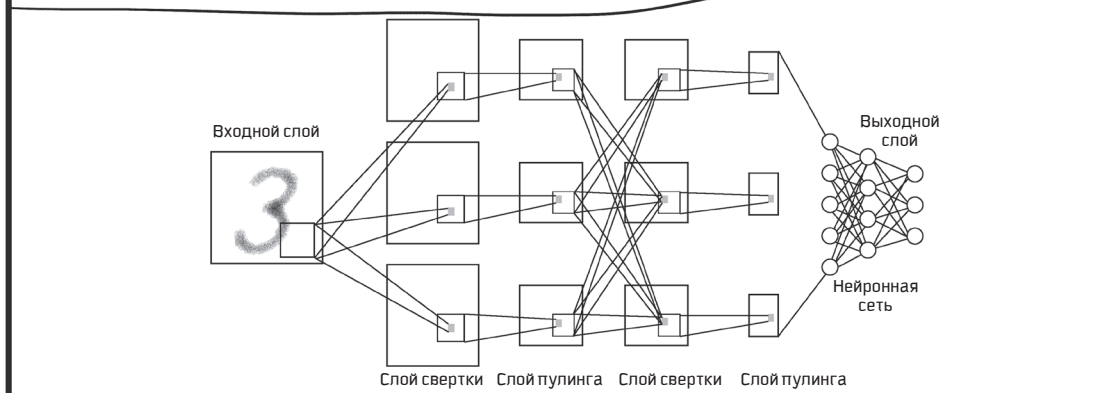
МЫ ГОВОРИЛИ О ХИТРОСТЯХ, КОТОРЫЕ ПОМОГАЮТ РЕШИТЬ ПРОБЛЕМЫ МНОГОСЛОЙНЫХ НЕЙРОННЫХ СЕТЕЙ, А ТЕПЕРЬ ПОГОВОРИМ О ЕЩЕ ОДНОМ МЕТОДЕ ОБУЧЕНИЯ.

И ЧТО ЭТО ЗА МЕТОД?

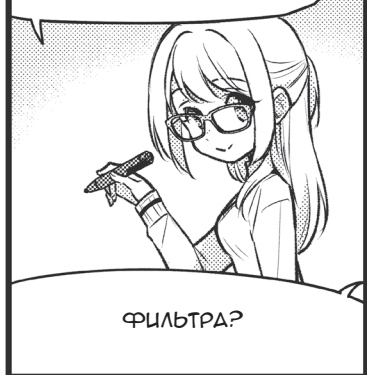
ЭТО МЕТОД СПЕЦИАЛИЗАЦИИ ЗАДАНИЙ В СТРУКТУРЕ НЕЙРОННОЙ СЕТИ. ПРИМЕРОМ СПЕЦИАЛИЗИРОВАННОЙ НЕЙРОННОЙ СЕТИ СЛУЖИТ **СВЕРТОЧНАЯ НЕЙРОННАЯ СЕТЬ**, КОТОРАЯ ИСПОЛЬЗУЕТСЯ ПРИ РАСПОЗНАВАНИИ ИЗОБРАЖЕНИЙ.



НА РИСУНКЕ ПОКАЗАНА НЕЙРОННАЯ СЕТЬ, ГДЕ ЧЕРЕДУЮТСЯ СЛОИ СВЕРТКИ И СЛОИ ПУЛИНГА. ПОСЛЕДНИЙ СЛОЙ ПУЛИНГА ПОЛУЧАЕТ ВЫВОД, И НА НЕМ РАЗМЕЩАЕТСЯ ОБЫЧНАЯ НЕЙРОННАЯ СЕТЬ.



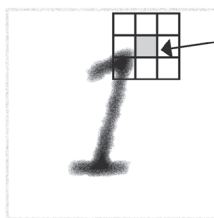
РАССМОТРИМ, КАК РАБОТАЕТ СВЕРТОЧНЫЙ СЛОЙ ПРИ ПОМОЩИ НАЛОЖЕНИЯ ФИЛЬТРА.



ПРИ НАЛОЖЕНИИ ФИЛЬТРА ВЫДЕЛЯЕТСЯ ПАТТЕРН, НАПРИМЕР ФИЛЬТР, КАК НА ИЗОБРАЖЕНИИ, ПРЕДСТАВЛЯЮЩИЙ СОБОЙ МАЛЕНЬКУЮ КАРТИНКУ 3×3, ТАК?

ТАК.

Развертка изображения по пикселям



$$\sum_{p=0}^2 \sum_{q=0}^2 x_{i+p,j+q} h_{pq}$$

-1	0	1
-1	0	1
-1	0	1

Фильтр выделения края h

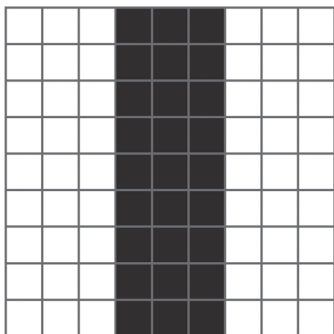
Данные изображения x

ФИЛЬТР ИЗОБРАЖЕНИЯ НАКЛАДЫВАЕТСЯ НА КАЖДЫЙ ПИКСЕЛЬ ЕГО РАЗВЕРТКИ, И ИЗ ИЗОБРАЖЕНИЯ ВЫДЕЛЯЕТСЯ ПАТТЕРН. В ЭТОМ ПРИМЕРЕ ФИЛЬТР ВЫДЕЛЕНИЯ КРАЯ БУДЕТ РЕАГИРОВАТЬ НА ИЗМЕНЕНИЕ ЦВЕТА В ВЕРТИКАЛЬНОМ НАПРАВЛЕНИИ.

В ВЕРТИКАЛЬНОМ НАПРАВЛЕНИИ? НО Я ВИЖУ ТОЛЬКО РЯДЫ ЦИФР: -1, 1 и 0.



НАПРИМЕР, ИСХОДНОЕ ИЗОБРАЖЕНИЕ - МОНОХРОМНОЕ, РАЗМЕРОМ 9x9 И ПРЕДСТАВЛЯЕТ СОБОЙ ВЕРТИКАЛЬНУЮ ПОЛОСУ. ТОГДА МОЖНО ЗАПИСАТЬ ЕГО ЦИФРАМИ ТАК: 0 - БЕЛЫЙ, А 1 - ЧЕРНЫЙ.

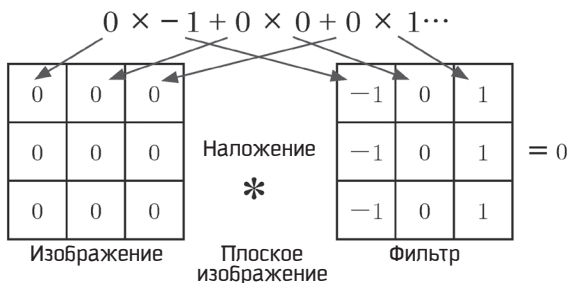


0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0
0	0	0	1	1	1	0	0	0

ТАК?

ДА.

НА ИЗОБРАЖЕНИЕ НАКЛАДЫВАЕТСЯ ФИЛЬТР, И ИЩЕТСЯ ПАТТЕРН. НАПРИМЕР, ЕСЛИ ФИЛЬТР НАКЛАДЫВАЕТСЯ НА КАРТИНКУ СЛЕВА, ТО ЗНАЧЕНИЯ ПЕРЕМНОЖАЮТСЯ И СТАНОВЯТСЯ РАВНЫМИ 0.



ВИЖУ!

А ЧТО, ЕСЛИ ПОМЕНЯТЬ ЗНАЧЕНИЯ
В ЦЕНТРАЛЬНОЙ ВЕРТИКАЛЬНОЙ ПОЛОСЕ
ФИЛЬТРА?

0	1	1
0	1	1
0	1	1

Изображение

Наложение
*

Изображение
с вертикальной
полосой

-1	0	1
-1	0	1
-1	0	1

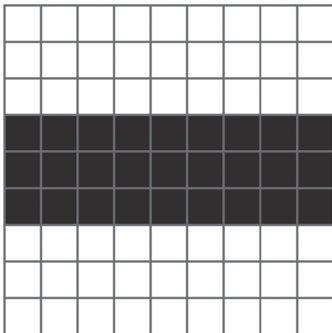
Фильтр

ТОГДА ПРИ УМНОЖЕНИИ
И СЛОЖЕНИИ ПОЛУЧИТСЯ 3.

3 ЖЕ БОЛЬШЕ, ЧЕМ 0?
ПРИ НАЛОЖЕНИИ ФИЛЬТРА ТОЛЬКО ВЕРТИКАЛЬНАЯ
ПОЛОСА БУДЕТ ИМЕТЬ БОЛЬШИЕ ЗНАЧЕНИЯ.
ЗНАЧИТ, МОЖНО ВЫДЕЛИТЬ ПАТТЕРН.

ПОНЯТНО.

А ЧТО ПОЛУЧИТСЯ, ЕСЛИ МЫ ВОЗЬМЕМ ИЗОБРАЖЕНИЕ
С ГОРИЗОНТАЛЬНОЙ ПОЛОСОЙ И НАЛОЖИМ НА НЕГО
ВЕРТИКАЛЬНЫЙ ФИЛЬТР ВЫДЕЛЕНИЯ КРАЯ?



0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

НУ, ЕСЛИ ПЕРЕМНОЖИТЬ,

0	0	0
1	1	1
1	1	1

*

-1	0	1
-1	0	1
-1	0	1

БУДЕТ 0.

ДА!

ДАЖЕ ЕСЛИ У НАС БУДЕТ ГОРИЗОНТАЛЬНАЯ ПОЛОСА, ТО ВСЕ РАВНО ПРИ ПЕРЕМНОЖЕНИИ И СЛОЖЕНИИ УЧАСТКОВ С ОДИНАКОВОЙ МОЩНОСТЬЮ (-1 СПРАВА И -1 СЛЕВА) ПОЛУЧИТСЯ 0. ДРУГИМИ СЛОВАМИ, ФИЛЬТР ВЫДЕЛЕНИЯ КРАЯ НЕ СМОЖЕТ ВЫДЕЛИТЬ ПАТТЕРН.

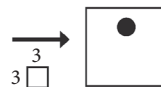
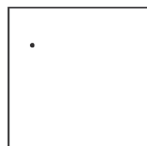
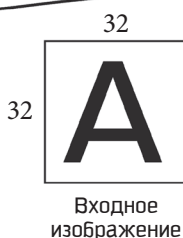


А ЕСЛИ ХОТИМ ВЫДЕЛИТЬ ДРУГОЙ ПРИЗНАК, ТО НАМ ЛУЧШЕ ПОМЕНЯТЬ ФИЛЬТР?

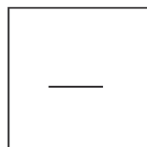


ДА, НА ПЕРВОМ СЛОЕ СВЕРТКИ КАК РАЗ И ОПРЕДЕЛЯЕТСЯ КОЛИЧЕСТВО ФИЛЬТРОВ ТОГО ЖЕ РАЗМЕРА, ЧТО И ИЗОБРАЖЕНИЕ НА ВХОДЕ, НЕОБХОДИМЫХ ДЛЯ ИЗВЛЕЧЕНИЯ ПРИЗНАКОВ.

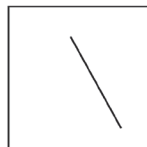
ПОСМОТРИ НА КАРТИНКУ. ЗАДЕСЬ ОПЕРАЦИЯ БУДЕТ ПРОВОДИТЬСЯ В СООТВЕТСТВИИ С ТРЕМЯ ТИПАМИ ФИЛЬТРОВ, КОТОРЫЕ БЫЛИ ОТОБРАНЫ НА ПЕРВОМ СЛОЕ СВЕРТКИ.



Развертка фильтра с самыми большими величинами данной территории



Повторение свертки и пулинга



Информация после понижающей передискретизации

Свертка

Совокупность точек, где был обнаружен паттерн фильтра (после использования функции активации)

Пулинг

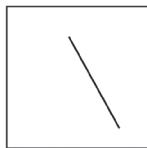
ЧТОБЫ РАСПОЗНАТЬ БУКВУ А, ИСПОЛЗУЮТСЯ ФИЛЬТРЫ УГЛА, ГОРИЗОНТАЛИ И ДИАГОНАЛИ - ВСЕГО ТРИ ФИЛЬТРА.

КАЖДЫЙ ЭЛЕМЕНТ СВЕРТОЧНОГО СЛОЯ СВЯЗАН ТОЛЬКО С ЧАСТЬЮ ВХОДНОГО ИЗОБРАЖЕНИЯ, И ВЕС ЕГО РАСПРЕДЕЛЯЕТСЯ МЕЖДУ ВСЕМИ НЕЙРОНАМИ. ОБЛАСТЬ ЭТОЙ СВЯЗИ СООТВЕТСТВУЕТ РАЗМЕРУ ФИЛЬТРА И НАЗЫВАЕТСЯ РЕЦЕПТИВНЫМ ПОЛЕМ.

ТО ЕСТЬ ФИЛЬТРЫ ТОЧНО
ТАК ЖЕ ВЫДЕЛЯЮТ
ПАТТЕРНЫ.



ДА. НО НА СЛЕДУЮЩЕМ СЛОЕ ПУЛИНГА
НЕОБХОДИМО СЖАТЬ ПОЛУЧЕННЫЕ ДАННЫЕ,
ОТРЕГУЛИРОВАВ ПРИЗНАКИ.



Объединенные
данные
фильтра



Информация после сжатия
(понижающей
передискретизации)

Пулинг

А ПОЧЕМУ?

НАПРИМЕР, ЕСЛИ НЕОБХОДИМО
РАСПОЗНАТЬ РУКОПИСНЫЙ ТЕКСТ,
ВЕДЬ ВСЕ ПИШУТ ПО-РАЗНОМУ, НЕ ТАК ЛИ?

ПОЧЕРК У ВСЕХ ЛЮДЕЙ
РАЗНЫЙ.

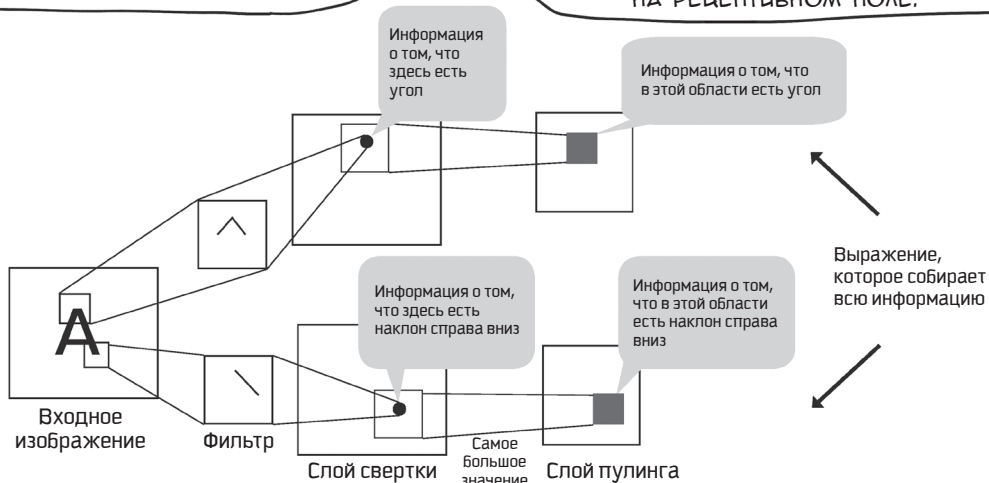


ПОЭТОМУ В СЛОЕ ПУЛИНГА
УБИРАЮТСЯ МЕЛКИЕ
ДЕТАЛИ ИНФОРМАЦИИ,
ЧТОБЫ УМЕНЬШИТЬ
РАЗЛИЧИЯ.

ЭТО КАК?

СЛОЙ ПУЛИНГА ОТЛИЧАЕТСЯ
ОТ СЛОЯ СВЕРТКИ МЕНЬШИМ
КОЛИЧЕСТВОМ ЭЛЕМЕНТОВ,
НО У КАЖДОГО ЭЛЕМЕНТА ТАКИЕ
ЖЕ РЕЦЕПТИВНЫЕ ПОЛЯ.

ГЛЯДЯ НА ЗНАЧЕНИЕ КАЖДОГО
ЭЛЕМЕНТА, МОЖНО ВЫВЕСТИ СРЕДНИЕ
И САМЫЕ БОЛЬШЕ ЗНАЧЕНИЯ.
ТАКИМ ОБРАЗОМ МОЖНО СГЛАДИТЬ
ИЗМЕНЕНИЯ ПОЗИЦИЙ ПАТТЕРНА
НА РЕЦЕПТИВНОМ ПОЛЕ.



ЧТОБЫ ТЫ ПОНЯЛ,
КАК РАБОТАЕТ СВЕРТОЧНАЯ
НЕЙРОННАЯ СЕТЬ,
ПРОДЕМОНСТРИРУЕМ
ВЫСОКОЭФФЕКТИВНОЕ
РАСПОЗНАВАНИЕ КАРТИНОК.

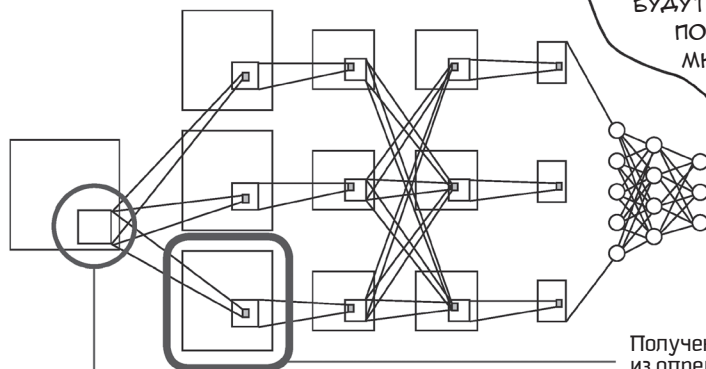


ДЛЯ ЭТОГО НЕОБХОДИМО
БОЛЬШОЕ КОЛИЧЕСТВО ФОТОГРАФИЙ.
И ГДЕ МОЖНО ЕГО НАЙТИ?

ДЕЙСТВИТЕЛЬНО.



В СВЕРТОЧНОЙ НЕЙРОННОЙ СЕТИ, ПО СРАВНЕНИЮ С НЕЙРОННЫМИ СЕТЯМИ,
ГДЕ ВСЕ ЭЛЕМЕНТЫ СОЕДИНЕННЫ, КОЛИЧЕСТВО МЕЖЭЛЕМЕНТНЫХ СВЯЗЕЙ
НЕВЫСОКО, ТАК КАК СУЩЕСТВУЕТ СТРУКТУРНОЕ ОГРАНИЧЕНИЕ -
ОПРЕДЕЛЕННЫЕ ЭЛЕМЕНТЫ ПОЛУЧАЮТ ТОЛЬКО ВЫХОДНЫЕ ДАННЫЕ
С ОПРЕДЕЛЕННОГО ПРЕДЫДУЩЕГО СЛОЯ.



Общий вес всех
рецептивных полей

Получение выходных данных
из определенной области

КРОМЕ ТОГО,
КАЖДАЯ СЛОИ
СВЕРТКИ РАСПРЕДЕЛЯЕТ ВЕСА
ДЛЯ РЕЦЕПТИВНЫХ ПОЛЕЙ, ПОЭТОМУ
ПАРАМЕТРЫ, ПРИ ПОМОЩИ КОТОРЫХ
ДОЛЖНО ВЕСТИСЬ ОБУЧЕНИЕ,
БУДУТ ЗНАЧИТЕЛЬНО УМЕНЬШЕНЫ.
ПОЭТОМУ МОЖНО ПОСТРОИТЬ
МНОГОСЛОЙНУЮ НЕЙРОННУЮ
СЕТЬ, В КОТОРУЮ МОЖНО
ВВОДИТЬ ИЗОБРАЖЕНИЯ
НАПРЯМУЮ, А ЦЕЛЮ
ОБУЧЕНИЯ МОЖЕТ
СТАТЬ ИЗВЛЕЧЕНИЕ
ПРИЗНАКОВ.

А ТЕПЕРЬ ПОПРОБУЕМ
РАСПОЗНАТЬ ИЗОБРАЖЕНИЯ
ИЗ БАЗЫ ДАННЫХ MNIST.

АА!





Для кодирования глубокого обучения в Python лучше всего использовать библиотеки.



Здесь мы будем использовать библиотеку Keras. Keras представляет собой надстройку над фреймворками DeepLearning4j, TensorFlow и Theano. Поскольку она хорошо описана, кодирование для задач глубокого обучения будет несложным.

```
import keras
```



MNIST – база данных рукописных изображений символов. Одно изображение представляет собой квадрат размером 28×28 пикселей со значением насыщенности от 0 до 255.



Мы используем 60 000 изображений для обучения и 10 000 для оценки. MNIST автоматически загрузится в Keras, а для разделения изображения на данные для обучения и для оценки используем метод `numpy array`.

```
from keras.datasets import mnist  
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```



Затем немного обработаем данные. Сначала преобразуем вход в стандартную систему распознавания изображений сверточной нейронной сети.



Поскольку обычно распознаются цветные изображения, в одном изображении будет трехмерный тензор (количество пикселей по вертикали \times количество пикселей по горизонтали \times значение цвета), а в данных для ввода будет четырехмерный тензор, к которому добавится количество изображений.



Так как изображения в градациях серого, которые мы используем сейчас, будут описаны количеством изображений \times количеством пикселей по вертикали \times количеством пикселей по горизонтали, четвертое измерение тензора будет представлено «цветом» со значением 1.



Поскольку на входе нейронной сети будет небольшая область значений от 0 до 1, нет необходимости на первом этапе настраивать вес или разряды коэффициентов для обучения. Поскольку максимальная величина пиксела равна 255, после преобразования целочисленного типа в тип с плавающей запятой выполнится операция деления всех данных на 255.

```
img_rows, img_cols = 28, 28
X_train = X_train.reshape(X_train.shape[0], img_rows, img_cols, 1)
X_test = X_test.reshape(X_test.shape[0], img_rows, img_cols, 1)
input_shape = (img_rows, img_cols, 1)
X_train = X_train.astype('float32') / 255
X_test = X_test.astype('float32') / 255
```



Далее настроим выход. Правильной меткой будет целое число от 0 до 9, которым обозначаются изображения, поэтому нужен десятимерный вектор, называющийся one-hot.



При использовании унитарного кодирования только один выход равен 1, а оставшиеся – 0. На выходном слое нейронной сети количество классов (на этот раз 10) будет преобразовано в формат, который будет легко подаваться в качестве учительского сигнала.


```
from keras.utils import to_categorical
Y_train = to_categorical(y_train)
Y_test = to_categorical(y_test)
```



А теперь определим структуру сверточной нейронной сети. У нас будет по два слоя свертки (фильтр размером 3×3) и пулинга (размер 2×2), выход будет конвертирован в одномерный вектор и перенаправлен на двухслойную нейронную сеть для классификации. В качестве функции активации на выходном слое используется softmax, а на остальных RELU.

```
from keras.models import Sequential
from keras.layers import Conv2D, MaxPooling2D, Flatten, Dense

n_out = len(Y_train[0]) # 10

model = Sequential()
model.add(Conv2D(16, kernel_size=(3, 3),
                 activation='relu',
                 input_shape=input_shape))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Conv2D(32, (3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dense(n_out, activation='softmax'))
model.summary()
```



Применяя метод `summary` Keras, посмотрим структуру полученной нейросети.

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 26, 26, 16)	160
max_pooling2d_1 (MaxPooling2)	(None, 13, 13, 16)	0
conv2d_2 (Conv2D)	(None, 11, 11, 32)	4640
max_pooling2d_2 (MaxPooling2)	(None, 5, 5, 32)	0
fl_atten_1 (Flatten)	(None, 800)	0
dense_1 (Dense)	(None, 128)	102528
dense_2 (Dense)	(None, 10)	1290
Total params: 108,618		
Trainable params: 108,618		
Non-trainable params: 0		



При помощи метода compile найдем categorical cross entropy и rmsprop, а также проведем обучение методом fit.

```

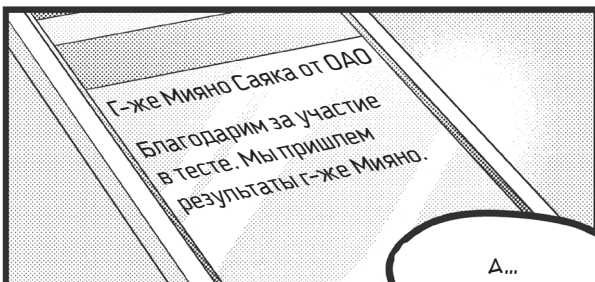
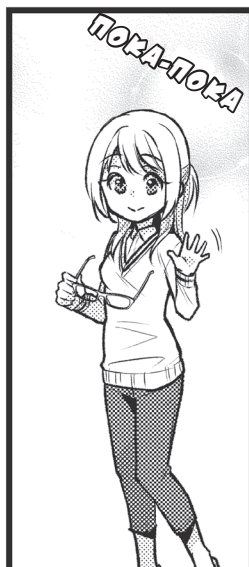
model.compile(loss = 'categorical_crossentropy',
              optimizer = 'rmsprop',
              metrics = ['accuracy'])
model.fit(X_train, Y_train, epochs=5, batch_size=200)
score = model.evaluate(X_test, Y_test, verbose=0)
print('Test loss:', score[0])
print('Test accuracy:', score[1])

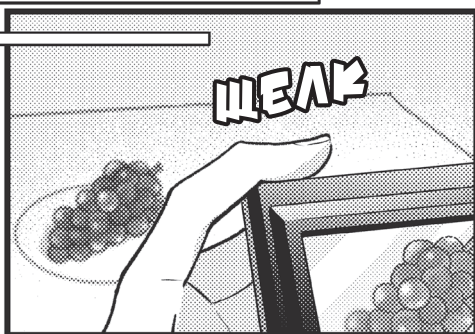
Epoch 1/5
60000/60000 [=====] - 13s 224us/step - loss: 0.2883 -
acc: 0.9130
Epoch 2/5
60000/60000 [=====] - 13s 210us/step - loss: 0.0763 -
acc: 0.9765
Epoch 3/5
60000/60000 [=====] - 14s 239us/step - loss: 0.0516 -
acc: 0.9836
Epoch 4/5
60000/60000 [=====] - 14s 238us/step - loss: 0.0384 -
acc: 0.9874
Epoch 5/5
60000/60000 [=====] - 14s 235us/step - loss: 0.0306 -
acc: 0.9906
Test loss: 0.03475515839108266
Test accuracy: 0.9878

```



Точность распознавания – 98,78 %. Она очень высокая.





СЕЙЧАС СИСТЕМА РАСПОЗНАЕТ
ВИНОГРАД С ТОЧНОСТЬЮ ДО 98 %.

Попросил прислать данные для обучения –
многожество фотографий винограда.

Выделил обучающую, тестовую и контрольную
выборки.

Сделал систему сортировки – сверточную
нейронную сеть на основании данных для
обучения.

Оценил при помощи данных для анализа,
откорректировал фильтры и элементы.

Наконец, протестировал на тестовой выборке.



ТЕПЕРЬ ДВЕ ОШИБКИ ИЗ 100 ФОТОГРАФИЙ,
ЗНАЧИТ, НАДО ПОДНЯТЬ ЭФФЕКТИВНОСТЬ.

УГУ.



Я ЖЕ НЕДАВНО
ТЕБЯ УГОЩАЛ!

КАК-ТО МАЛОВАТО
ЗА МОЮ ПОМОЩЬ...

ЭТО ТЕБЯ.

АА? ЧТО?

Голос
из телефона

УГОВОРИ ИХ
НА 98 %,
ПОМОЩНИК!

ЧТО?

ЭТО ХАСИМОТО!

А-АЛЛО!

ЭТО КИЁХАРА ИЗ МЭРИИ.
ПРОСТИТЕ, ЕСЛИ НЕ ВОВРЕМЯ.

А, КИЁХАРА-САН! ЗАРАВСТВУЙТЕ.
ВЫ О СОРТИРОВКЕ ВИНОГРАДА?
ВСЕ НАСТРОИЛИ?

КУДЗЁ-Ё-Ё-Ё-Ё-Ё-Ё!

АА, НО ТОЧНОСТЬ
СОСТАВЛЯЕТ 98 %.

98 %?!

НУ...
МЫ БЫ ХОТЕЛИ ЕЕ
ЕЩЕ УВЕЛИЧИТЬ,
ЕСЛИ МОЖНО...

ХА-ХА-ХА, КИЁХАРА-САН,
И 98 % ХВАТИТ.

ЧТО?

ЛЮДИ НА 5 % ОШИБАЮТСЯ,
ТАК ЧТО ВСЕ В ПОРЯДКЕ!

СПАСИБО ЗА ВАШУ УСЕРДНУЮ РАБОТУ!
ВЫ НАС ОЧЕНЬ ВЫРУЧИЛИ!

АА НЕТ,
СПАСИБО ВАМ...

МОЛОДЕЦ, КИЁХАРА!
ПОШЛИ УЖИНАТЬ...

НУ ПОШЛИ...

ЖАЛЬ Я НЕ СМОГ
УЛУЧШИТЬ
ПРОИЗВОДИТЕЛЬНОСТЬ
ДО ПРЕДЕЛА...

У МЕНЯ ЕЩЕ МАЛО
ОПЫТА, НО Я УЖЕ
МОГУ ПОМОГАТЬ
ЛЮДЯМ...

БУДУ СТАРАТЬСЯ.

Как важно, что сегодня
я услышал «спасибо»

БЫСТРЕЙ, КИЁХАРА.

ЦАДУ!

Математическое повторение (4)

Киёхара-кун, кажется, исправился. Я слышала, что у него были проблемы с сайтом.



Ну, надо было мне его сначала научить как следует... Но хорошо, что он исправился.

Я только слышала ваш разговор, но не кажется ли тебе, что он уже повзрослел?



Да, когда он учился, положиться на него было нельзя. Но сейчас он, кажется, стал получше...

Я хотела бы с ним встретиться, если можно. Кстати, разве он не симпатичный?



Ну... Да вроде бы и нет, не знаю. Когда можно будет встретиться, я дам знать.

Да, если есть время. (Какая упрямая сестрица. Бедный Киёхара.)



Сегодня мы говорили о глубоком обучении. Но это на практике выглядит страшно.



Если привыкнуть к тензорам, то его код становится простым. Тензор – многомерное пространство на языке программирования.

Я знаю, что в программировании есть двухмерное распределение, но кроме этого ничего не слышала.



Ай-тян, ты писала программу по обработке изображений?

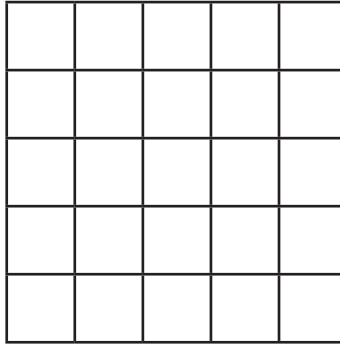
Я этим занималась, но программу я не писала.



Ну, представь себе черно-белое изображение. Пусть черный – 0, а белый – 255, а цвета между ними могут принимать значения от 1 до 254. Чем темнее серый, тем меньше число, а чем светлее, тем больше. И как мы выразим данные, которые расположены по горизонтали и вертикали в виде прямоугольника?

В виде двухмерной фигуры из целых чисел?





Ага. Эта двухмерная фигура, которую можно представить в виде матрицы.



И ее можно записать в фигурных скобках.



А теперь возьмем цветное изображение. Обычно используется формат RGB, где 1 пиксель может иметь значение красного, зеленого или синего цвета. При сложении трех цветов получается еще один цвет, и можно записывать разные выражения.

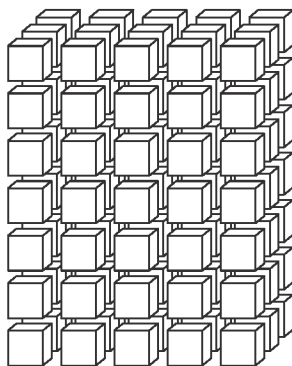


Это интересно.



Как тогда будут представлены данные?

Поскольку красный, зеленый и синий представлены двумерными массивами одного размера, то если их сгруппировать, получится трехмерный массив?



Да, численное выражение структуры данных в таком случае будет выражено 3D-тензором.



А еще есть 4D-тензор, который состоит из нескольких 3D-изображений...

Как нескольких цветных изображений... И это все данные для обучения?



Именно. Какой тензор будет использоваться для видео?

Если у нас видео, то изображения выстроены по оси времени. И чтобы его обработать, 4D-мерный тензор становится... 5D?



Да!

Ого... тензоры все сложнее и сложнее.



Но в машинном обучении используются только они. В глубоком обучении, чтобы упростить последнюю классификацию, преобразуют тензоры, и все становится легко.

Ага... Кстати, тетя спрашивала, нашла ли ты работу?



Да вроде... Но пока еще думаю...

Да... у взрослых много разных дел.

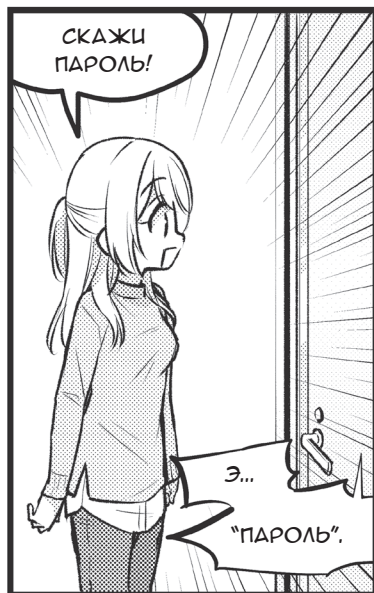
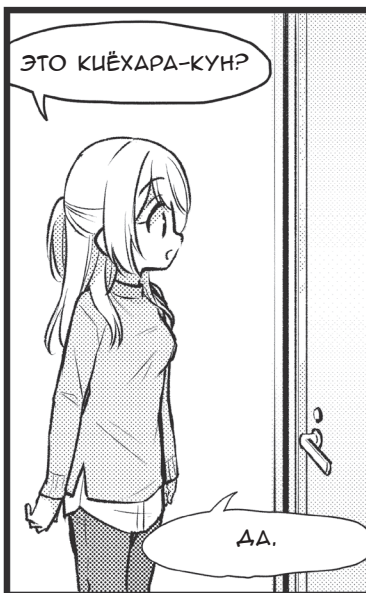
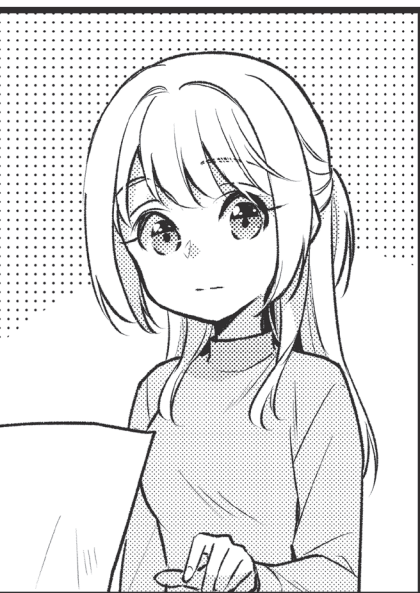
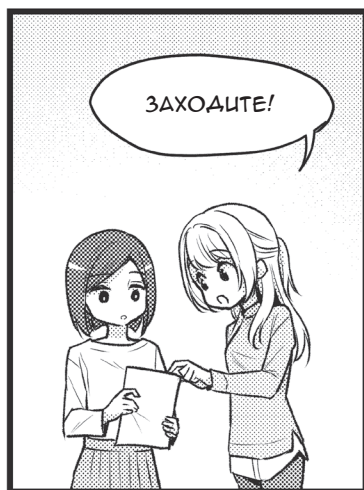
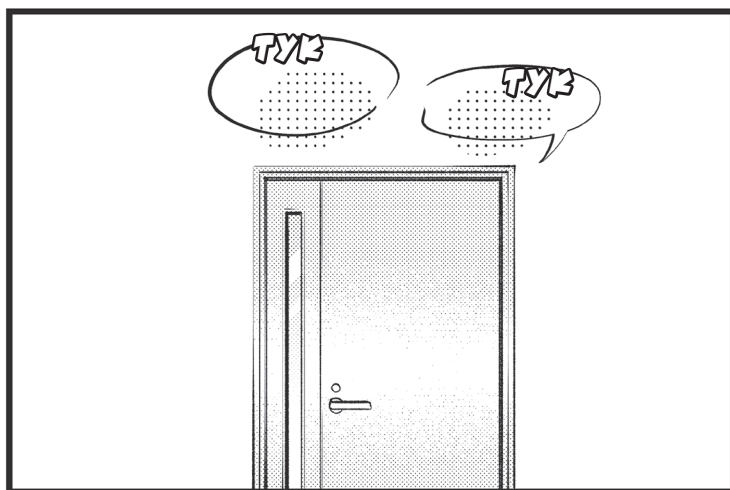


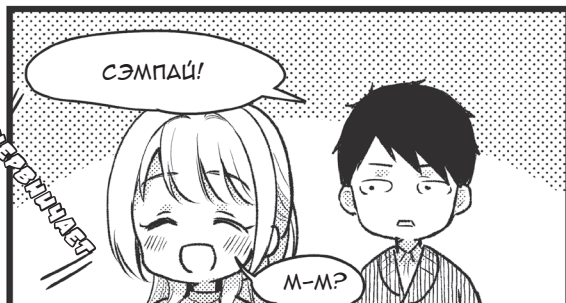
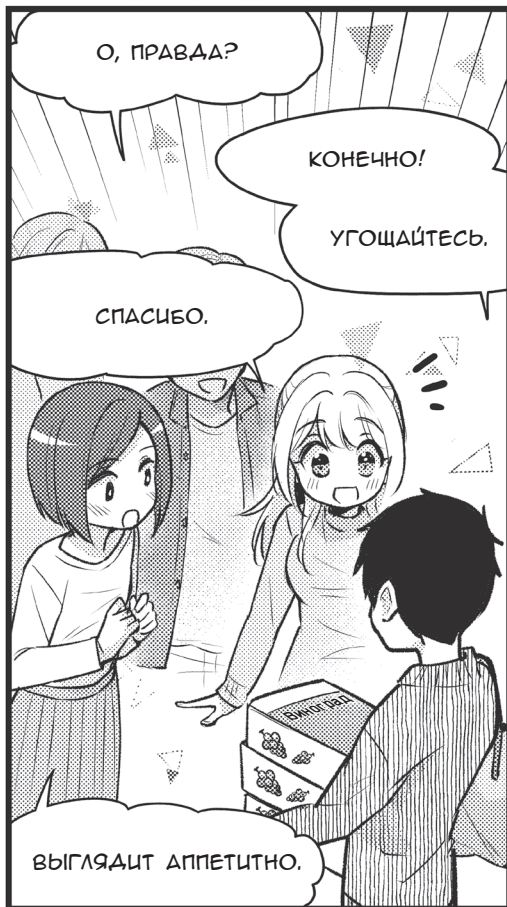
ГЛАВА 5

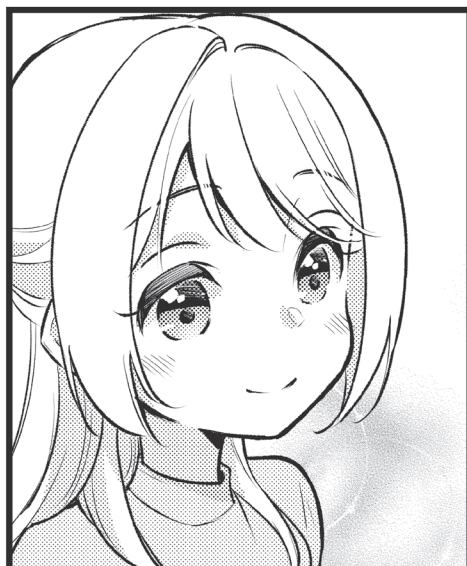
АНСАМБЛЕВЫЕ МЕТОДЫ

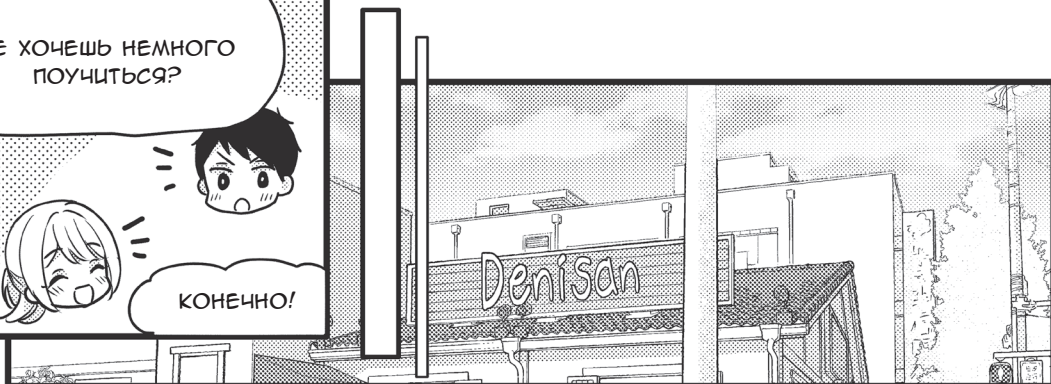
БУДЕМ ИСПОЛЬЗОВАТЬ
НЕСКОЛЬКО
КЛАССИФИКАТОРОВ!

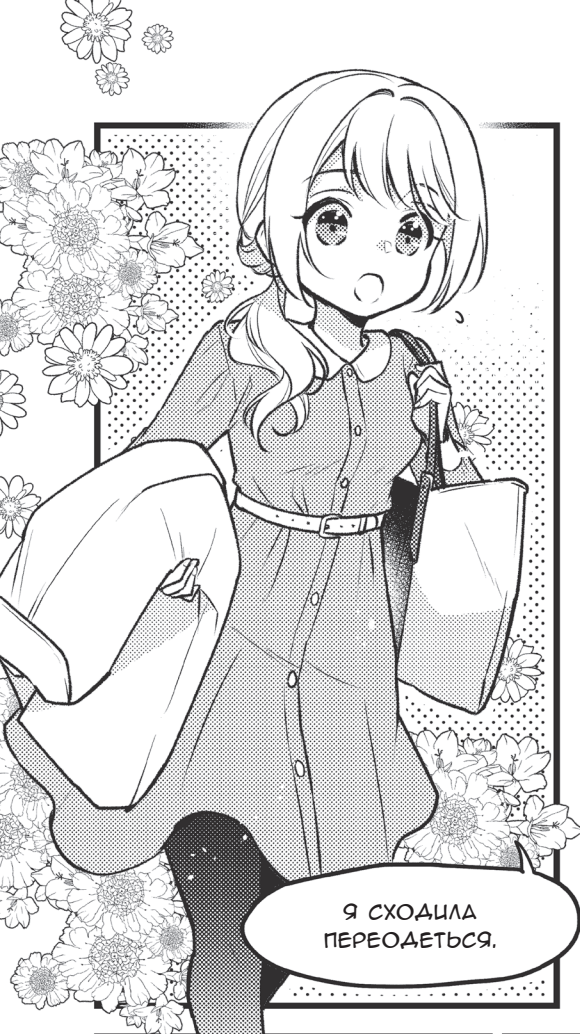












СЕГОДНЯ ПОГОВОРИМ
ОБ АНСАМБЛЕВЫХ
МЕТОДАХ!

АНСАМБЛЕВЫХ?



ДА... ХОТЯ ГЛУБОКОЕ ОБУЧЕНИЕ ВПОЛНЕ
УСПЕШНО СПРАВЛЯЕТСЯ С КЛАССИФИКАЦИЕЙ
ДАННЫХ, СВЯЗАННЫХ СО ЗВУКОМ,
ИЗОБРАЖЕНИЯМИ ИЛИ ЕСТЕСТВЕННОЙ РЕЧЬЮ,
НЕОБХОДИМО СМОТРЕТЬ,
КАК ОНО БУДЕТ РАБОТАТЬ
СО СЛОЖНЫМИ ПРИЗНАКАМИ
В КАЖДОМ ОТДЕЛЬНОМ СЛУЧАЕ.

Данные со множеством признаков

Пол	Возраст	ИМТ	Уровень глюкозы	Дав- ление	Диабет
Ж	65	22	180	135	Нет
М	60	28	200	140	Да
М	75	21	175	120	Нет
Ж	72	25	195	140	Нет

Я съел вкусный пирожок



Отношения между
близкими данными

Данные необязательно должны быть
близкими

КАК ЭТО?

ОДИН ИЗ СПОСОБОВ СПРАВИТЬСЯ
СО СЛОЖНЫМИ ПРИЗНАКАМИ -
ОБУЧЕНИЕ ПРИ ПОМОЩИ
АНСАМБЛЕВЫХ МЕТОДОВ.

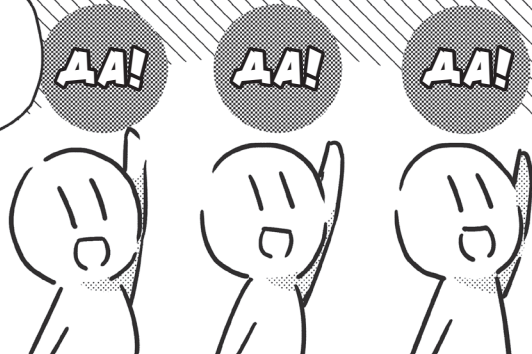
КАКИХ МЕТОДОВ?

ЭТО МЕТОД,
ПРИ КОТОРОМ СОЕДИНЯЮТСЯ
НЕСКОЛЬКО ОБУЧАЮЩИХ АЛГОРИТМОВ,
И ИХ СОЧЕТАНИЕ ОКАЗЫВАЕТСЯ
БОЛЕЕ ЭФФЕКТИВНЫМ.



* Мондзю – бодхисатва, олицетворение высшей мудрости. – Прим. перев.

ОДНАКО ЭТОТ МЕТОД
НАДО ИСПОЛЬЗОВАТЬ С УМОМ -
ЧТО, ЕСЛИ ТРИ ЧЕЛОВЕКА
СКАЖУТ ОДНО И ТО ЖЕ?

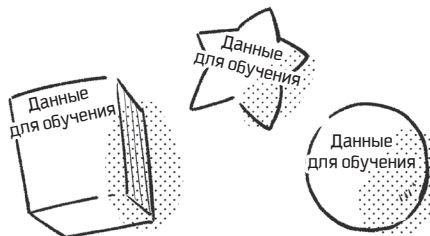




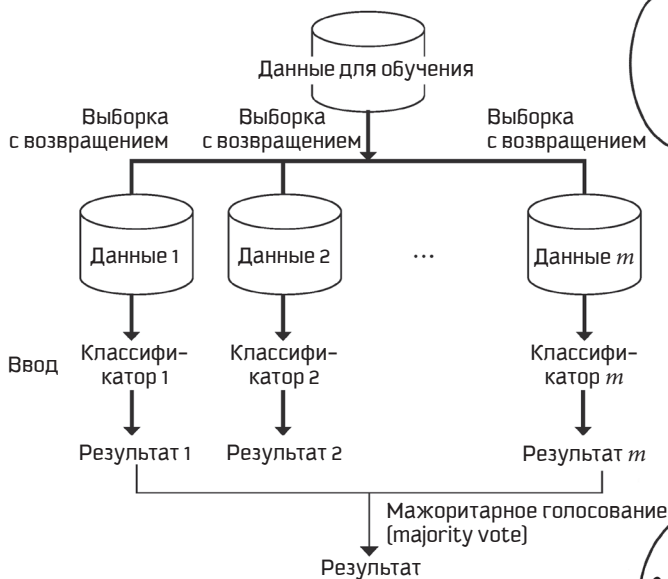
5.1. БЭГГИНГ



ПЕРВАЯ МЫСЛЬ: ЧТОБЫ СДЕЛАТЬ АЛГОРИТМЫ С ОТЛИЧАЮЩИМСЯ ПОВЕДЕНИЕМ, НУЖНО ИСПОЛЬЗОВАТЬ НЕСКОЛЬКО ОТЛИЧАЮЩИХСЯ НАБОРОВ ДАННЫХ ДЛЯ ОБУЧЕНИЯ.



НО, НАВЕРНОЕ, ЭТО ОЧЕНЬ ТРУДНО - ПОДГОТОВИТЬ РАЗНЫЕ ДАННЫЕ ДЛЯ ОБУЧЕНИЯ?



ПРИ БЭГГИНГЕ ИЗ ДАННЫХ ДЛЯ ОБУЧЕНИЯ ДЕЛАЕТСЯ **ВЫБОРКА С ВОЗВРАЩЕНИЕМ**, КОГДА ИЗ ИСХОДНЫХ ДАННЫХ ВЫБИРАЕТСЯ НЕКОТОРОЕ КОЛИЧЕСТВО НАБОРОВ ДАННЫХ ОДИНАКОВОГО РАЗМЕРА. ЗАТЕМ ДЛЯ КАЖДОГО НАБОРА ДАННЫХ СОЗДАЕТСЯ КЛАССИФИКАТОР ПРИ ПОМОЩИ ОДНОГО И ТОГО ЖЕ АЛГОРИТМА.

ВЫБОРКА С ВОЗВРАЩЕНИЕМ?

ПРИ ЭТОМ МЕТОДЕ ДАННЫЕ, ПОПАВШЕ В ВЫБОРКУ, ЗАПИСЫВАЮТСЯ И ВОЗВРАЩАЮТСЯ. КАКИЕ-ТО ДАННЫЕ МОГУТ ПОПАСТЬ В ВЫБОРКУ МНОГО РАЗ, А КАКИЕ-ТО - НИ РАЗУ.

ДАВАЙ ПОСЧИТАЕМ, НАСКОЛЬКО ПРИ ИСПОЛЬЗОВАНИИ ЭТОГО МЕТОДА БУДУТ ОТЛИЧАТЬСЯ ДАННЫЕ В РАЗНЫХ ВЫБОРКАХ!

ДА!

ПЕРЕВОРАЧИВАЕТ СТРАНИЦУ





Допустим, в наборе данных есть N отдельных элементов. Какова вероятность, что один элемент не попадет в выборку?

Если элементов N , то вероятность того, что он попадет в выборку, $-\frac{1}{N}$.
А вероятность, что не попадет, $-(1 - \frac{1}{N})$.



Да. А вероятность того, что данные не попадут в выборку N раз, равна $(1 - \frac{1}{N})^N$. Таким образом, рассчитаем вероятность того, что данные не попадут в выборку...



Если $N = 10$, то она равна 0,349.
Если $N = 100$, то 0,366.
Если $N \rightarrow \infty$, то вероятность равна $1/e = 0,368$.



Число N не слишком влияет на результат.

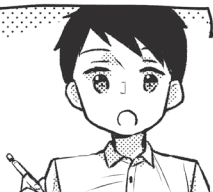


Именно! При таких расчетах ясно, каким бы N ни было, при выборке с возвращением примерно $\frac{1}{3}$ исходных данных не попадут в нее.



Целая треть данных!

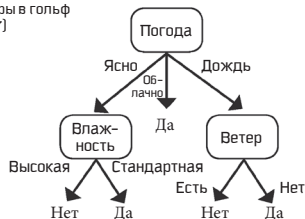
КАКОЙ МЕТОД МАШИННОГО ОБУЧЕНИЯ МЫ ИСПОЛЬЗУЕМ ДЛЯ СОЗДАНИЯ КЛАССИФИКАТОРА?



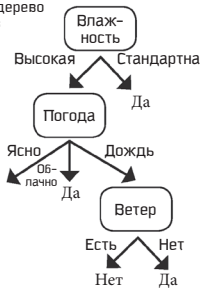
В ПРИНЦИПЕ, ЛЮБОЙ МЕТОД МОЖЕТ ПОДОЙТИ, НО АЛГОРИТМ, КОТОРЫЙ СОЗДАЕТ КЛАССИФИКАТОР, БУДЕТ ЗАВИСЕТЬ ОТ НЕСТАБИЛЬНОСТИ, ИНЫМИ СЛОВАМИ, ОН БУДЕТ ЧУВСТВИТЕЛЕН К РАЗНИЦЕ ДАННЫХ ДЛЯ ОБУЧЕНИЯ.

НАПРИМЕР, ЕСЛИ У РЕШАЮЩЕГО ДЕРЕВА ДАННЫЕ НЕНАМНОГО ОТЛИЧАЮТСЯ, ТО КЛАССИФИКАТОРЫ МОГУТ БЫТЬ РАЗНЫМИ.

Решающее дерево по данным для игры в гольф (стр. 57)



Решающее дерево (6 примеров убрано)



ого!

ПОСКОЛЬКУ КАЖДЫЙ КЛАССИФИКАТОР ОБУЧАЕТСЯ НА ОДИНАКОВОМ КОЛИЧЕСТВЕ ДАННЫХ, ВСЕ КЛАССИФИКАТОРЫ СЧИТАЮТСЯ ОДИНАКОВО НАДЕЖНЫМИ, И В РЕЗУЛЬТАТЕ ОТВЕТ ДАЕТСЯ ПРОСТЫМ БОЛЬШИНСТВОМ ГОЛОСОВ (MAJORITY VOTE).

БОЛЬШИНСТВО ГОЛОСОВ



5.2. СЛУЧАЙНЫЙ ЛЕС

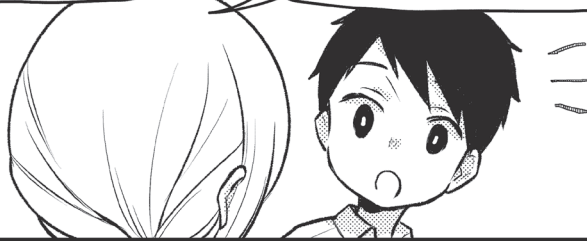
ВТОРОЙ МЕТОД



СУЩЕСТВУЕТ МЕТОД СЛУЧАЙНОГО ЛЕСА, КОТОРЫЙ ОТЛИЧАЕТСЯ ПО КЛАССИФИКАТОРАМ ОТ БЭГГИНГА.

ЧЕМ ОН ОТЛИЧАЕТСЯ ОТ БЭГГИНГА?

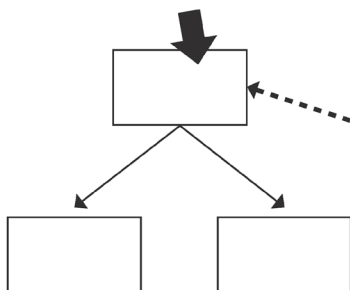
НУ, ТАМ ТОЖЕ ИСПОЛЬЗУЮТСЯ ВЫБОРКИ С ВОЗВРАЩЕНИЕМ НАБОРОВ ДАННЫХ ОДИНАКОВОГО РАЗМЕРА, КОТОРЫЕ ВЫДЕЛЯЮТСЯ ИЗ ДАННЫХ ДЛЯ ОБУЧЕНИЯ.



СТРОИТСЯ ДЕРЕВО РЕШЕНИЙ
ПО КЛАССИФИКАТОРАМ
ДЛЯ КАЖДОГО НАБОРА ДАННЫХ.

ВЫБИРАЕТСЯ ЗАРАНЕЕ ОПРЕДЕЛЕННОЕ КОЛИЧЕСТВО
ПРИЗНАКОВ ИЗ ВСЕХ, А УЖЕ В НИХ ИЩЕТСЯ ПРИЗНАК
С МАКСИМАЛЬНО ЭФФЕКТИВНЫМ РАЗДЕЛЕНИЕМ.

Выбирается условие
для деления данных



Возраст

Давление

ИМТ*

Уровень глюкозы
в крови

↓ Извлекается случайное
количество

Возраст

Давление

ИМТ*

Уровень глюкозы
в крови

Выбираются признаки с высоким
информационным выигрышем

Не используются

* Индекс массы тела.

И КАК ОН ОПРЕДЕЛЯЕТ ЧИСЛО ПРИЗНАКОВ?

КОЛИЧЕСТВО
ПРИЗНАКОВ ДЛЯ
ВЫБОРА ИЗ ОБЩЕГО
ЧИСЛА ПРИЗНАКОВ
(d) – ЭТО ЧАСТО
ЛИБО КВАДРАТНЫЙ
КОРЕНЬ ИЗ d ,
ЛИБО $\log_2 d$.

Общее количество признаков d	$\text{floor}(\sqrt{d})$	$\text{floor} \log_2 d$
5	2	2
10	3	3
50	7	5
100	10	6

$\text{floor}(x)$ – это наибольшее целое число, но не больше x .

ОПЕРАЦИЯ ПРОДОЛЖАЕТСЯ
РЕКУРСИВНО, ПОКА В КАЖДОМ
ЛИСТЕ НЕ ОСТАНУТСЯ
ПРЕДСТАВИТЕЛИ ОДНОГО КЛАССА.

ИТАК, КИЁХАРА-КУН,
ВОПРОС!

НА ЧТО НУЖНО ОБРАЩАТЬ
ВНИМАНИЕ, КОГДА МЫ
СТРОИМ РЕШАЮЩЕЕ
ДЕРЕВО?

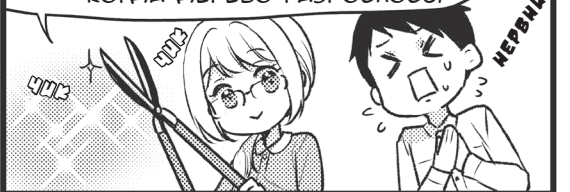


ПЕ-ПЕРЕОБУЧЕНИЕ!



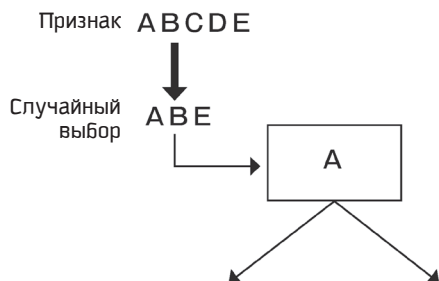
ИМЕННО!

РАНЬШЕ МЫ ГОВОРИЛИ, ЧТО В КАЧЕСТВЕ
МЕРЫ ПРОТИВ ПЕРЕОБУЧЕНИЯ МОЖНО
ЛИБО ОСТАНОВИТЬ РОСТ, КОГДА
КОЛИЧЕСТВО ДАННЫХ В ЛИСТЬЯХ УПАДЕТ
НИЖЕ ОПРЕДЕЛЕННОГО УРОВНЯ,
ЛИБО ЖЕ ОБРЕЗАТЬ ВЕТВИ,
КОГДА ДЕРЕВО РАЗРОСЛОСЬ.

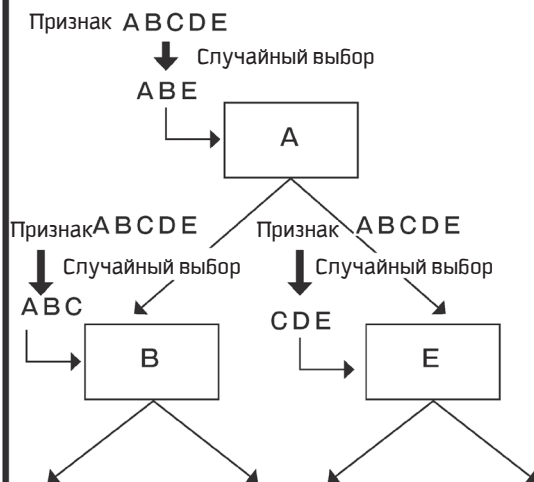




СНАЧАЛА ПОДСЧИТЫВАЕТСЯ
КЛАССИФИЦИРУЮЩАЯ
СПОСОБНОСТЬ КАЖДОГО
ПРИЗНАКА, И ЗАТЕМ ВЫБИРАЕТСЯ
ПРИЗНАК С САМОЙ ВЫСОКОЙ.



ПОТОМ ДЛЯ ОТДЕЛЬНЫХ НАБОРОВ
ДАННЫХ СЛУЧАЙНЫМ ОБРАЗОМ
ОПРЕДЕЛЯЕТСЯ НАБОР ПРИЗНАКОВ,
ИЗ НИХ ВЫБИРАЕТСЯ ПРИЗНАК
С САМОЙ ВЫСОКОЙ СПОСОБНОСТЬЮ
КЛАССИФИКАЦИИ, И ДЕРЕВО РАСТЕТ.



ТАК ИЗ ПОХОЖИХ ДАННЫХ
МОЖНО ПОСТРОИТЬ ОТЛИЧАЮЩИЕСЯ
РЕШАЮЩИЕ ДЕРЕВЬЯ.

5.3. БУСТИНГ

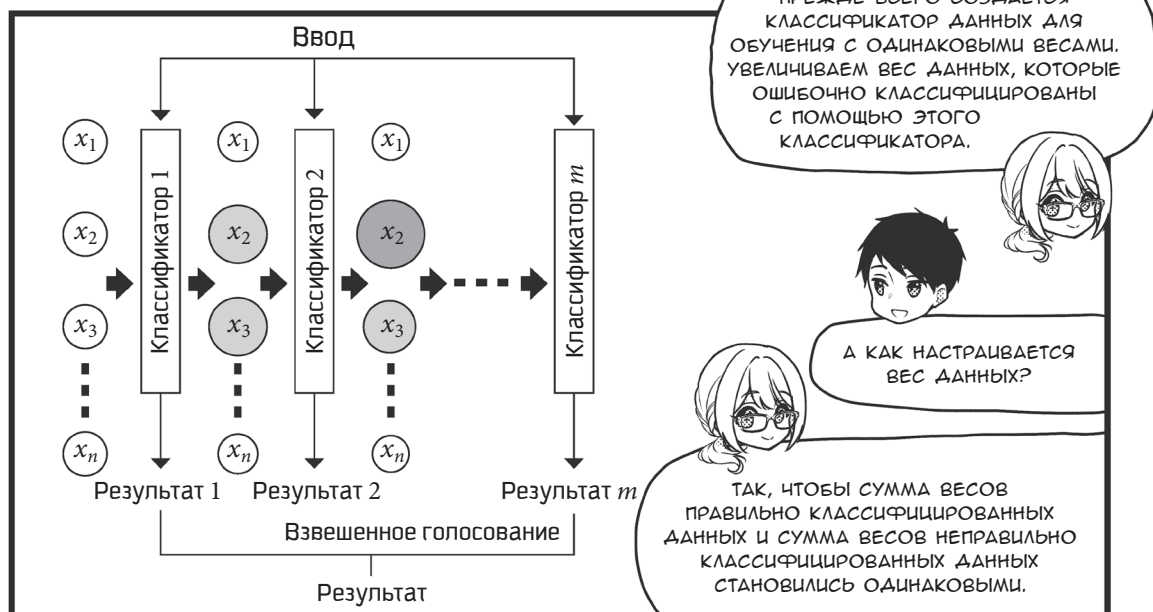
В БЭГГИНГЕ И МЕТОДЕ СЛУЧАЙНОГО
ЛЕСА ИЗМЕНЯЮТСЯ ИСПОЛЬЗУЕМЫЕ
НАБОРЫ ДАННЫХ, ИЗМЕНЯЮТСЯ УСЛОВИЯ
ВЫБОРА КЛАССИФИКАТОРОВ, И ПОЭТОМУ
МОЖНО СТРОИТЬ ОТЛИЧАЮЩИЕСЯ
КЛАССИФИКАТОРЫ.
КРОМЕ ТОГО, В МЕТОДЕ БУСТИНГА
НАБОР КЛАССИФИКАТОРОВ,
КОТОРЫЕ ВЕДУТ СЕБЯ
ПО-РАЗНОМУ, СОЗДАЕТСЯ ПУТЕМ
ПОСЛЕДОВАТЕЛЬНОГО
ДОБАВЛЕНИЯ КЛАССИФИКАТОРОВ,
КОТОРЫЕ СПЕЦИАЛИЗИРУЮТСЯ
НА УМЕНЬШЕНИИ ОШИБОК.

**ТРЕТИЙ
МЕТОД**

КЛАССИФИКАТОР, КОТОРЫЙ
СПЕЦИАЛИЗИРУЕТСЯ
НА УМЕНЬШЕНИИ
КОЛИЧЕСТВА ОШИБОК?

ДА.

ДЛЯ ЭТОГО НУЖНО
ОПРЕДЕЛИТЬ ВЕС
КАЖДОГО ЭЛЕМЕНТА
ДАННЫХ.



ЗАТЕМ ДЛЯ НАБОРА ДАННЫХ, В КОТОРОМ БЫЛИ ИЗМЕНЕНЫ ВЕСА, ПРОВОДИТСЯ ОБУЧЕНИЕ ПРИ ПОМОЩИ СЛЕДУЮЩЕГО КЛАССИФИКАТОРА.

и ТАК ПОЭТАПНО СОЗДАЮТСЯ НОВЫЕ КЛАССИФИКАТОРЫ. НОВЫЙ КЛАССИФИКАТОР, В ОТЛИЧИЕ ОТ РАНЕЕ СОЗДАННОГО, КЛАССИФИЦИРУЮЩЕГО ОШИБОЧНЫЕ ДАННЫЕ, ВОСПОЛНЯЕТ ЕГО СЛАБЫЕ СТОРОНЫ.

ЭТО МЕТОД ADABOOST.

Остались ошибки...

Оставь их мне!

УЧЕБНЫЙ АЛГОРИТМ КЛАССИФИКАТОРА, КОТОРЫЙ ИСПОЛЬЗУЕТСЯ В БУСТИНГЕ, ДОЛЖЕН В ОСНОВНОМ ПРИМЕНЯТЬ ВЕСА ДАННЫХ В КАЧЕСТВЕ КРИТЕРИЯ ДЛЯ СОЗДАНИЯ НОВОГО КЛАССИФИКАТОРА.

ТО ЕСТЬ ЭТОТ МЕТОД НЕЛЬЗЯ ИСПОЛЬЗОВАТЬ, ЕСЛИ НЕ ДУМАТЬ О ВЕСАХ С САМОГО НАЧАЛА.

НЕТ, МОЖНО СПРАВИТЬСЯ, СОЗДАВАЯ НАБОРЫ ДАННЫХ ПУТЕМ ВЫБОРКИ С ВОЗВРАЩЕНИЕМ, ЗАНОВО РАСПРЕДЕЛЯЯ ВЕСА.

УДОБНО! А РЕШЕНИЕ ПРИНИМАЕТСЯ БОЛЬШИНСТВОМ ГОЛОСОВ, КАК В БЭГГИНГЕ?

В СЛУЧАЕ МЕТОДА ADABOOST КЛАССИФИКАТОР СОЗДАЕТСЯ НА ОСНОВАНИИ ОШИБОК НА ПРЕДЫДУЩЕМ ЭТАПЕ, ТАК?

ЦИНЫМИ СЛОВАМИ, ЕСЛИ В ИСХОДНЫХ ДАННЫХ ДЛЯ ОБУЧЕНИЯ МНОГО ОШИБОК, ТО НАДЕЖНОСТЬ КЛАССИФИКАТОРА, СОЗДАННОГО НА ОСНОВЕ ЭТИХ ДАННЫХ, ДЛЯ НЕИЗВЕСТНЫХ ДАННЫХ НА ВХОДЕ БУДЕТ УВЕЛИЧИВАТЬСЯ.

НАДЕЖНОСТЬ

КАК ОПРЕДЕЛЯЕТСЯ РЕЗУЛЬТАТ КЛАССИФИКАЦИИ?

ПОСКОЛЬКУ В ADABOOST К КЛАССИФИКАТОРАМ НА ПРЕДЫДУЩИХ ЭТАПАХ ОДИН ЗА ДРУГИМ ДОБАВЛЯЮТСЯ КЛАССИФИКАТОРЫ, КОТОРЫЕ ПРАВИЛЬНО КЛАССИФИЦИРОВАЛИ ОШИБОЧНЫЕ ДАННЫЕ, ЕГО ЭФФЕКТИВНОСТЬ ПОВЫШАЕТСЯ.

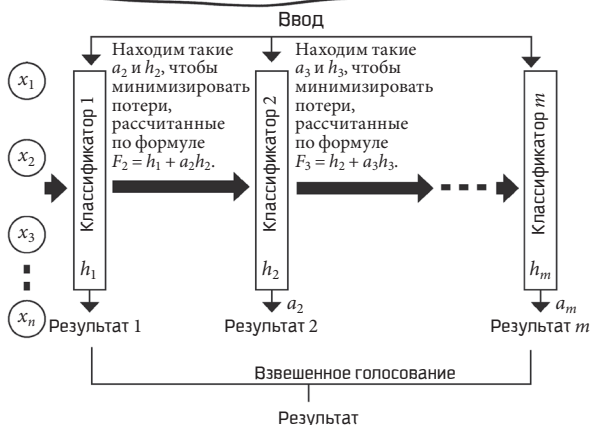
ВЫЧИСЛЯЕТСЯ ВЕС НА ОСНОВАНИИ ВЕЛИЧИНЫ ФУНКЦИИ ОШИБОК ДЛЯ КАЖДОГО КЛАССИФИКАТОРА, ПРОВОДИТСЯ ГОЛОСОВАНИЕ ВЕСОВ, И ПОЛУЧАЕТСЯ РЕЗУЛЬТАТ.

ЕСТЬ СПОСОБ, ПРИ КОТОРОМ В КАЧЕСТВЕ ЕЩЕ ОДНОГО УРАВНЕНИЯ МОЖНО ИСПОЛЬЗОВАТЬ ФУНКЦИЮ ПОТЕРИ. СОСТАВНОЙ КЛАССИФИКАТОР, КОТОРЫЙ ЯВЛЯЕТСЯ РЕЗУЛЬТАТОМ БУСТИНГА, МОЖЕТ ОПРЕДЕЛИТЬ ФУНКЦИЮ ПОТЕРИ.

МОЖНО ПРЕДСТАВИТЬ УРАВНЕНИЕ, ПРИ КОТОРОМ ПОСЛЕДУЮЩИЕ КЛАССИФИКАТОРЫ БУДУТ ВЫБИРАТЬСЯ ТАК, ЧТОБЫ МАКСИМАЛЬНО УМЕНЬШИТЬ ФУНКЦИЮ ПОТЕРИ.

БУСТИНГ, ОСНОВАННЫЙ НА ЭТОЙ ИДЕЕ, НАЗЫВАЕТСЯ ГРАДИЕНТНЫМ БУСТИНГОМ.

А ТЕПЕРЬ НЕМНОГО ПОКОАИРУЕМ!





Возьмем данные diabetes.arff из инструментов для машинного обучения Weka. Они близки к данным, необходимым для сайта с тестом определения вероятности диабета.



Хотя в scikit-learn есть такой же набор данных, но он нужен для задач регрессии, и объяснить результат будет трудно, поэтому для задач классификации мы используем набор diabetes.arff.



diabetes.arff содержит в себе результаты осмотров и диагнозы. Признаки – возраст, давление, ИМТ и т. п.



Для начала надо скачать diabetes.arff. Можно скачать несколько наборов данных с <https://www.cs.waikato.ac.nz/ml/weka/datasets.html> и найти там нужный файл.

```
import numpy as np
from scipy.io import arff
from sklearn.ensemble import BaggingClassifier, RandomForestClassifier,
AdaBoostClassifier, GradientBoostingClassifier
from sklearn.model_selection import cross_val_score
```



Файл формата arff будет прочитан как модуль arff scipy. Поскольку признаковые описания и метки правильных ответов находятся в одной строке, каждый можно сохранить в отдельный массив numpy.

```
data, meta = arff.loadarff('diabetes.arff')
X = np.empty((0,8), np.float)
y = np.empty((0,1), np.str)
```

```
for e in data:
    e2 = list(e)
    X = np.append(X, [e2[0:8]], axis=0)
    y = np.append(y, e2[8:9])
```



В scikit-learn ансамблевые методы используются так же, как и обучение и оценка данных с классификатором. Используя перекрестную проверку по 10 % данных, определим среднюю точность и дисперсию.

```
clf1 = BaggingClassifier()
scores = cross_val_score(clf1, X, y, cv=10)
print("{0:4.2f} +/- {1:4.2f} %".format(scores.mean() * 100, scores.std() * 100))
73.69 +/- 5.11 %

clf2 = RandomForestClassifier()
scores = cross_val_score(clf2, X, y, cv=10)
print("{0:4.2f} +/- {1:4.2f} %".format(scores.mean() * 100, scores.std() * 100))
74.72 +/- 5.72 %

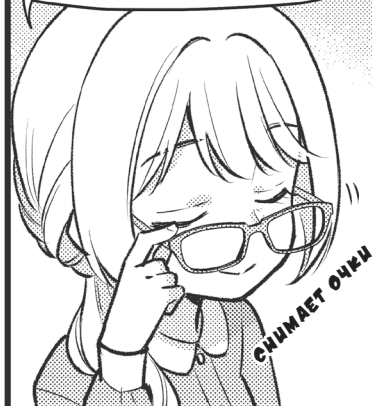
clf3 = AdaBoostClassifier()
scores = cross_val_score(clf3, X, y, cv=10)
print("{0:4.2f} +/- {1:4.2f} %".format(scores.mean() * 100, scores.std() * 100))
75.52 +/- 5.71 %

clf4 = GradientBoostingClassifier()
scores = cross_val_score(clf4, X, y, cv=10)
print("{0:4.2f} +/- {1:4.2f} %".format(scores.mean() * 100, scores.std() * 100))
76.30 +/- 5.11 %
```



Используя параметры по умолчанию, можно получить хорошие результаты при помощи градиентного бустинга.

НУ, НА СЕГОДНЯ ВСЕ
ПРО АНСАМБЛЕВЫЕ МЕТОДЫ.



УДАЧИ
С ОТКРЫТИЕМ САЙТА!



ЛУЧШЕ БЫ Я В ТОТ РАЗ
ВСЕ ДОСЛУШАЛ...

НАВЕРНОЕ, НАДО
БЕЖАТЬ, ПОТОМУ ЧТО
БОЮСЬ СНОВА
УСЛЫШАТЬ, ЧТО Я ЕЁ
КАК "БРАТИК".



НЕТ-НЕТ,
ЭТО Я САМА
ВЫЗВАЛАСЬ ТЕБЯ
УЧИТЬ!

КСТАТИ, МНЕ
НЕДАВНО ПРИШЛО
ПИСЬМО...

ПИСЬМО?..

МЕНЯ ВЗЯЛИ НА РАБОТУ!



НУ ЭТО ЖЕ ПРЕКРАСНО!
ПОЗДРАВЛЯЮ!

УСТАЛО

ДА...
НО Я ДУМАЛА, ЧТО НЕ ВОЗЬМУТ.

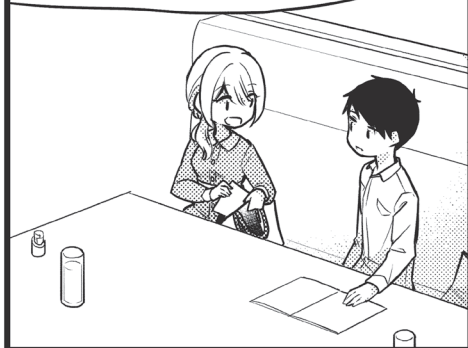


И ПОСЛЕ ВЫПУСКА
Я ПЕРЕЕДУ В ТОКИО.



КОНЕЧНО, БУДЕТ ТРУДНО
ВСТРЕЧАТЬСЯ, И Я НЕ СМОГУ
ТЕБЯ БОЛЬШЕ УЧИТЬ.

Я ЖАЛЕЮ ОБ ЭТОМ,
ПОТОМУ ЧТО ТЫ ХОРОШИЙ
УЧЕНИК.



А ТЕПЕРЬ Я ЗАПЛАЧУ.



ЭТО Я ДОЛЖЕН,
ПЛАТА ЗА ОБУЧЕНИЕ!

НО ТЫ ЖЕ УЖЕ ЗАПЛАТИЛ
ДО ЭТОГО, НЕ ХОЧУ,
ЧТОБЫ ТЫ ПОСТОЯННО
МЕНЯ УГОЩАЛ!



НЕТ!



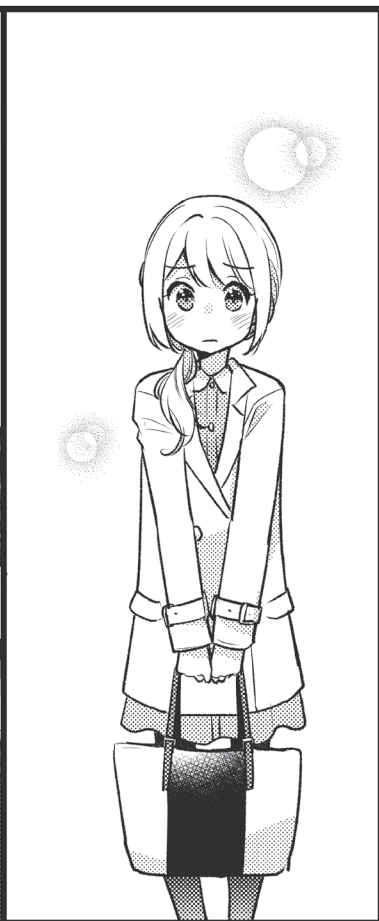
ТО МЛАДШИЙ
ТОВАРИЩ,
ТО БРАТИК,

ТО УЧЕНИК.



ПОЗВОЛЬ МНЕ ЗАПЛАТИТЬ,
КАК МУЖЧИНЕ!





Математическое повторение (5)

Сегодня вы много говорили о классификаторах и их эффективности, но мне не совсем понятно, что это такое.



Что было непонятно?

Ну, я поняла, что даже если много людей дают одинаковый ответ, они не умнее одного, который ответит так же. Но если мы соберем много людей, умных и не очень, и они будут давать разные ответы, они не поссорятся?



Конечно, могут поссориться. Но чтобы они не поссорились, мы используем голосование по методу большинства, и оно может численно подтвердить наилучший ответ.



Прежде всего пусть для одинаковых данных для обучения мы будем использовать L разных классификаторов.



То есть у нас есть L умных людей.



Ну, можно взять и неумных людей. Хоть это немногим лучше, чем случайные ответы, но математическая формулировка не изменится.

Правда?



Предположим, что процент ошибок классификатора e один и тот же, и ошибки независимы.



Под независимостью ошибок имеется в виду то, что вероятность ошибки классификатора для каждого элемента данных независима. Иными словами, можно предположить, что нет данных, для которых вместе ошибется множество классификаторов.



Основываясь на этой гипотезе, рассмотрим вероятность того, что из L классификаторов m классификаторов ошиблись. Пусть $m = 1$.

То есть если вероятность того, что ошибся один классификатор, равна e , то вероятность того, что остальные $(L - 1)$ классификаторы не ошиблись, равна $(1 - e)^{(L-1)}$.



Так как ошибиться может любой из L классификаторов, то при перемножении получим $Le(1 - e)^{(L-1)}$.



А если ошиблись два?

Тогда число способов выбрать 2 из L будет LC_2 , и, значит, у нас будет $LC_2 e_2 (1 - e)^{(L-2)}$.

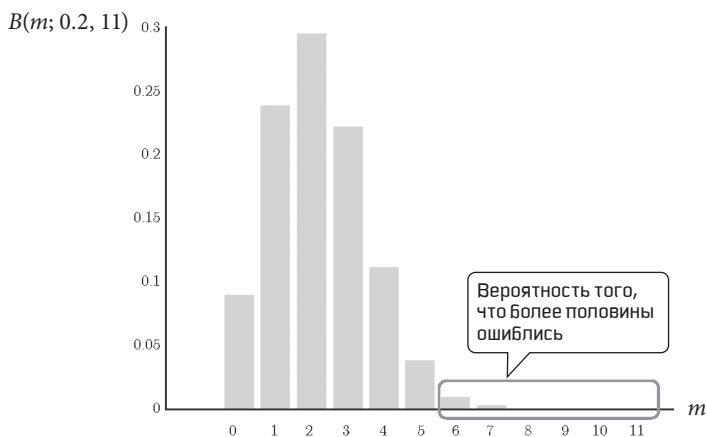


Так. Вероятность, что m из L классификаторов, вероятность ошибки каждого из которых e , находится при помощи биномиального распределения $B(m; e, L)$.

$$B(m; e, L) = {}_L C_m e^m (1 - e)^{L-m}.$$



Предположим, что количество классификаторов $L = 11$, коэффициент ошибки каждого $e = 0.2$, тогда биномиальное распределение $B(m; 0.2, 11)$ будет представлено на графике ниже.



Если результат классификации определяется большинством голосов, то вероятность полностью неправильного решения будет равняться сумме вероятностей ошибок шести классификаторов и больше. При подсчете получится 1,2 %.



Если вероятность ошибки каждого классификатора по отдельности 20 %, но мы используем все 11 классификаторов, то вероятность ошибки составит всего 1,2 %.

Ого! Ансамблевые методы – это здорово!



Однако и у этого метода есть недостатки. Данные, которые трудно классифицировать даже людям, могут вызвать ошибки у многих классификаторов, и поэтому на самом деле эффективность не так высока.



Но чтобы по возможности разрешить эту проблему, можно делать классификаторы с разным поведением, для чего и нужны ансамблевые методы. Понятно в целом?

В целом да. Жаль, что ты, Сая, уезжаешь в Токио. Там будет тебе одиноко...



Конечно, уезжать из дома всегда немного грустно...

А что ты сказала Киёхара-куну?





Ну... ну что он мне как братик, но это было еще тогда...



И почему ты это вспомнила?

ГЛАВА 6

ОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ

КЛАСТЕРИЗАЦИЯ?
РАЗЛОЖЕНИЕ МАТРИЦЫ?



Отдел здравоохранения
и благосостояния

ТАК...

КИЁХАРА-КУН, МОЖНО
НА МИНУТОЧКУ?

ДА, ДИРЕКТОР.

ЧТО СЛУЧИЛОСЬ?

ВСТАЕТ

ИЗВИНИТЕ, ЧТО ОТВЛЕКАЮ,
МОЖНО С ВАМИ НЕМНОГО
ПОВОРОЧИТЬ?

ХОЧУ ВАС КОЕ О ЧЕМ ПОПРОСИТЬ,
КИЁХАРА-КУН. ВЫ ЖЕ ХОРОШО
РАЗБИРАЕТЕСЬ В МАШИННОМ
ОБУЧЕНИИ?

НУ-У-У-У-У-У-У...

У САЙТА, КОТОРЫЙ
ОПРЕДЕЛЯЕТ ВЕРОЯТНОСТЬ
ЗАБОЛЕВАНИЯ ДИАБЕТОМ,
ХОРОШИЕ ОТЗЫВЫ.

ТАК СКАЗАТЬ

С-СПАСИБО!

НЕ СТОИТ. В ЭТОМ ГОДУ,
КАК ВАМ ИЗВЕСТНО,
Я ВЫХОЖУ НА ПЕНСИЮ.

И ЕСТЬ ОДНО ДЕЛО,
СЛИШКОМ СЛОЖНОЕ
ДЛЯ МОЕГО ПРЕЕМНИКА.

ЧТО ЗА ДЕЛО?



ТАК...



ЧТО-ТО ТЫ НЕДАВНО СТАЛ
УЧИТЬСЯ ВО ВРЕМЯ ОБЕДА...

УГУ.



ДИРЕКТОР СКАЗАЛ, ЧТО ОТЗЫВЫ НА САЙТ
С ПРОГНОЗАМИ ДИАБЕТА ХОРОШИЕ, И ПОПРОСИЛ
СДЕЛАТЬ КОЕ-ЧТО, ДЛЯ ЧЕГО НАДО ОБУЧЕНИЕ
БЕЗ УЧИТЕЛЯ...



НУ... МОЖЕТ, ТЫ К СЭМПАЮ
СВОЕМУ ОБРАТИШЬСЯ?



РАССКАЖИ МНЕ, МОЛ,
КАК НАДО СДЕЛАТЬ ТО И ЭТО.



КСТАТИ, ЧТО-ТО ПИСЬМА ОТ НЕЕ
КАКИЕ-ТО ФОРМАЛЬНЫЕ!



МОЖЕТ, Я ЧТО-ТО НЕ ТО СКАЗАЛ?

НЕРВНИЧАЕТ



НУ, Я МОГУ СПРОСИТЬ
У ТОВАРИЩА ПО ИГРЕ, КОТОРЫЙ
ПОМОГ МНЕ С ЖУРНАЛОМ
ПРО МАШИННОЕ ОБУЧЕНИЕ.



АА? СПАСИБО!

Неделю спустя



ТАК...

ЧТОБЫ ПРОАНАЛИЗИРОВАТЬ МОДЕЛИ
ПОВЕДЕНИЯ ПОЖИЛЫХ, НУЖНА
КЛАСТЕРИЗАЦИЯ, А ЧТОБЫ ВЫБРАТЬ
НУЖНУЮ ИНФОРМАЦИЮ -
РАЗЛОЖЕНИЕ МАТРИЦЫ.



НО Я СЛИШКОМ
МНОГОГО НЕ ЗНАЮ...



АЗЫНЬ

М-М-М?

СООБЩЕНИЕ.



ТЫК
ТЫК

О, ЭТО ОТ
САЯКА-СЭМПАЙ!

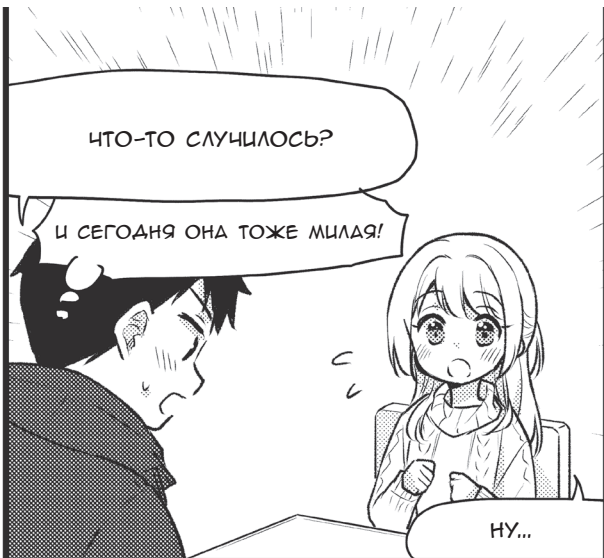


Извини, что отвлекаю
Не хочешь встретиться
в кафе после работы?

АЗЫНЬ!

ДОБРО
ПОЖАЛОВАТЬ!





Я ПОПЫТАЛСЯ САМ НАЙТИ
ИНФОРМАЦИЮ ПРО КЛАСТЕРИЗАЦИЮ
И РАЗЛОЖЕНИЕ МАТРИЦЫ.

НАПРАВЛЕНИЕ ВЕРНОЕ!

ХОРОШО!

СЭМПАЙ,
ВЫ МНЕ РАССКАЖЕТЕ
ПРО ОБУЧЕНИЕ БЕЗ
УЧИТЕЛЯ, ПРАВДА?

КОНЕЧНО!

УЛЫБАЕТСЯ

Я ЖЕ РАДА ВИДЕТЬ,
ЧТО ТЫ ЧЕМУ-ТО УЧИШЬСЯ!

АГА!

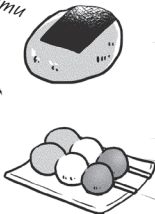
РАДОСТНО

6.1. КЛАСТЕРИЗАЦИЯ

ОДИН ИЗ МЕТОДОВ ОБУЧЕНИЯ БЕЗ УЧИТЕЛЯ – **КЛАСТЕРНЫЙ АНАЛИЗ** – ИСПОЛЬЗУЕТСЯ ДЛЯ РАЗДЕЛЕНИЯ ДАННЫХ ПО ГРУППАМ.

Например,

сладости



ЭТО ОЧЕНЬ ШИРОКО. ПОПРОБУЙ ОПРЕДЕЛИТЬ КЛАСТЕРИЗАЦИЮ, ИСПОЛЬЗУЯ СЛОВО "РАССТОЯНИЕ".

НУ...

ДАННЫЕ, КОТОРЫЕ ОБЪЕДИНЯЮТСЯ В ОДНУ ГРУППУ, ДОЛЖНЫ НАХОДИТЬСЯ НА БЛИЗКОМ РАССТОЯНИИ ДРУГ ОТ ДРУГА.

И МОЖНО СКАЗАТЬ, ЧТО ИХ ЭТО ОБЪЕДИНЯЕТ.



рядом



Расстояния между элементами велики, нельзя объединить

Расстояния малы, можно объединить



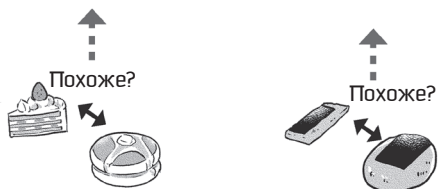
НЕОБХОДИМО ЕЩЕ, ЧТОБЫ ДАННЫЕ, КОТОРЫЕ НАХОДЯТСЯ В РАЗНЫХ ГРУППАХ, БЫЛИ НА ДАЛЕКОМ РАССТОЯНИИ ДРУГ ОТ ДРУГА.

ИМЕННО!

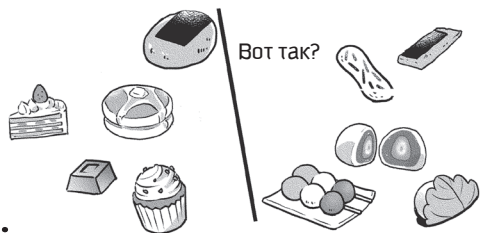
КЛАСТЕРИЗАЦИЯ ДЕЛИТСЯ НА **ИЕРАРХИЧЕСКУЮ**, КОГДА НОВЫЕ КЛАСТЕРЫ СОЗДАЮТСЯ ПУТЕМ ОБЪЕДИНЕНИЯ БОЛЕЕ МЕЛКИХ КЛАСТЕРОВ И ДЕРЕВО СОЗДАЕТСЯ ОТ ЛИСТЬЕВ К СТВОЛУ; И НА **РАЗДЕЛЯЮЩУЮ (ДИВИЗИВНУЮ)**...

...КОГДА НОВЫЕ КЛАСТЕРЫ СОЗДАЮТСЯ ПУТЕМ ДЕЛЕНИЯ БОЛЕЕ КРУПНЫХ КЛАСТЕРОВ НА БОЛЕЕ МЕЛКИЕ И ДЕРЕВО СОЗДАЕТСЯ ОТ СТВОЛА К ЛИСТЬЯМ.

Иерархическая кластеризация



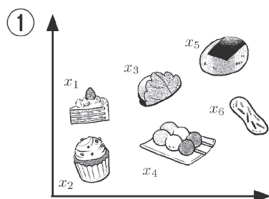
Разделяющая кластеризация



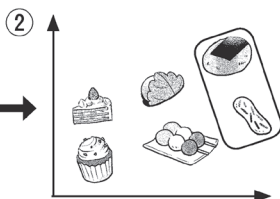
6.1.1. Иерархическая кластеризация

ИЕРАРХИЧЕСКАЯ КЛАСТЕРИЗАЦИЯ – ЭТО ПРОЦЕСС, ПРИ КОТОРОМ СОЗДАЮТСЯ НОВЫЕ КЛАСТЕРЫ ПУТЕМ ОБЪЕДИНЕНИЯ БОЛЕЕ МЕЛКИХ КЛАСТЕРОВ. КЛАСТЕРЫ ПОСТЕПЕННО РАСТУТ.

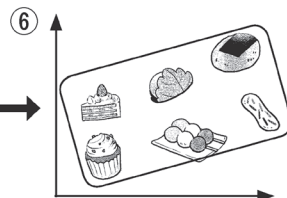
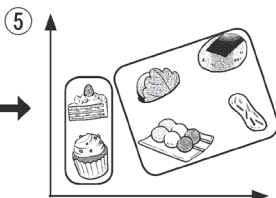
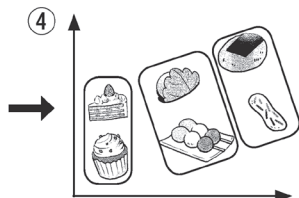
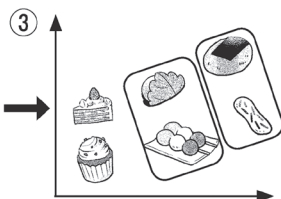
АЛГОРИТМ ВЫГЛЯДИТ ТАК:



① Данные объединяются в кластеры



②–⑤ Кластеры, которые находятся рядом, объединяются – так получают новые кластеры



⑥ В конце концов данные объединяются в один кластер

ПОНЯТНО, КОГДА ДАННЫЕ НАХОДЯТСЯ РЯДОМ, НО КАК ВЫЧИСЛИТЬ РАССТОЯНИЕ МЕЖДУ ДАННЫМИ И КЛАСТЕРОМ, А ТАКЖЕ МЕЖДУ ДВУМЯ КЛАСТЕРАМИ?



ЕСЛИ ПРЕДСТАВИТЬ, ЧТО В КЛАСТЕР ВХОДИТ ТОЛЬКО ОДИН ЭЛЕМЕНТ ДАННЫХ, ТО РАССТОЯНИЕ МЕЖДУ КЛАСТЕРАМИ МОЖНО ОПРЕДЕЛИТЬ ПО СТЕПЕНИ ИХ СХОДСТВА.



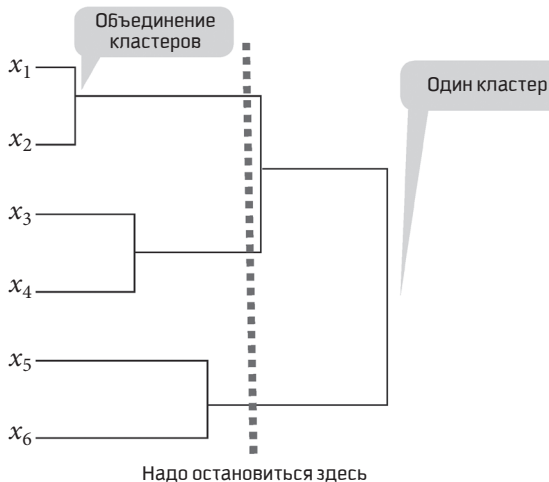
А ЧТОБЫ ОПРЕДЕЛИТЬ СТЕПЕНЬ СХОДСТВА, СУЩЕСТВУЮТ СЛЕДУЮЩИЕ МЕТОДЫ:

Метод одиночной связи	Метод полной связи	Центроидный метод	Метод Уорда
Сходство определяется величиной расстояния между самыми близкими элементами	Сходство определяется величиной расстояния между самыми дальними элементами	Сходство определяется расстоянием между центроидами кластеров	После объединения кластеров вычисляется квадрат среднего расстояния между данными и центром кластера, из него вычитается эта величина до объединения
Есть тенденция к удлинению кластеров в одном направлении	Есть тенденция к избеганию удлинения кластеров в одном направлении	Тенденция к удлинению кластеров находится между одиночной и полной связью	Часто получаются сравнительно хорошие кластеры

КАКИМ МЕТОДОМ ЛУЧШЕ ВСЕГО РАЗДЕЛИТЬ ДАННЫЕ НА ТРИ КЛАСТЕРА?



ЕСЛИ ЗАПИСАТЬ ОПЕРАЦИЮ ОБЪЕДИНЕНИЯ КЛАСТЕРОВ В ВИДЕ ДЕРЕВА, КАК ПОКАЗАНО НА РИСУНКЕ, ТО СНАЧАЛА КОЛИЧЕСТВО КЛАСТЕРОВ БУДЕТ РАВНЯТЬСЯ КОЛИЧЕСТВУ ДАННЫХ N , А ДАЛЕЕ УМЕНЬШАТЬСЯ ПО ОДНОМУ ЗА КАЖДУЮ ОПЕРАЦИЮ. КЛАСТЕРЫ ПОСТЕПЕННО ОБЪЕДИНЯЮТСЯ В ОДИН. В ПРИНЦИПЕ, МОЖНО ПОЛУЧИТЬ ЛЮБОЕ ИХ КОЛИЧЕСТВО. ЕСЛИ НАДО ПОЛУЧИТЬ ТРИ КЛАСТЕРА, ТО ЛУЧШЕ ОСТАНОВИТЬСЯ ЗДЕСЬ.



6.1.2. Разделяющая кластеризация



А что за метод – разделяющая кластеризация?

Разделяющая кластеризация – метод, который определяет функцию оценки качества разделения данных и выдает результат лучшего значения этой функции.



А чем он отличается от иерархической кластеризации?

В иерархической кластеризации деление идет снизу вверх, поэтому с точки зрения целого могут создаваться искаженные кластеры. С этой точки зрения кластеры, которые создаются при разделяющей кластеризации, более оптимальны.



Так. Получается, что разделяющая кластеризация всегда дает лучшие результаты, чем иерархическая?..

Ну...



Например, какие вычисления надо провести, если надо перебрать и найти лучший результат функции оценки при разделении N элементов данных на два кластера?

Один элемент данных может быть разделен по кластерам одним из двух способов. Два – из четырех. А три, получается, из восьми.



Получается 2 в степени N . Если N – большое число, то на практике невозможно оценить все значения функции.

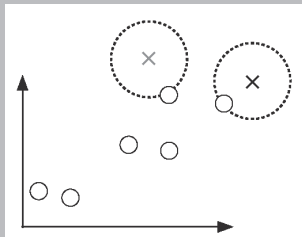


В таком случае в качестве обычного метода используют поиск оптимального разделения. Основной метод разделяющей кластеризации – метод k -средних, можешь рассказать о нем?

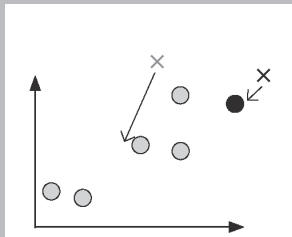
Да. Выглядит он так. У нас изначально должно быть задано количество кластеров k .



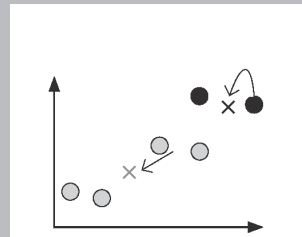
① Определяем средний вектор для каждого из k кластеров, сформированных из выбранных случайным образом элементов данных, определяем центр кластера.



② Элемент данных становится частью того кластера, центр которого ближе всего.
③ Средний вектор каждого кластера перевычисляется.



④ Повторяем пункты 2 и 3, пока средние векторы всех кластеров не перестанут меняться.



Повторяем пункты ② и ③



А какая у этого метода функция оценки?

Ну, сумма расстояний между каждым элементом данных и средним вектором кластера, к которому он принадлежит.



А почему при этом величина функции оценки уменьшается?

На втором этапе изменение кластера означает, что найден средний вектор, который находится еще ближе, поэтому величина функции оценки уменьшится.

Далее, на третьем этапе, снова рассчитывается средний вектор, и так повторяется, пока он не перестанет меняться; это будет означать, что найдена наименьшая сумма расстояний данных в кластере, поэтому величина функции оценки становится меньше.



Да. Однако при использовании этого метода мы находим локальное оптимальное решение. Я специально говорю «локальное», потому что этот метод не позволяет сказать, является ли оно общим для целого.

Значит, можно найти ответ и получше, но метод останавливается на полпути.



Да. Однако если провести несколько итераций с разными входными данными, то кластеризация методом k -средних даст хороший результат.

КСТАТИ, СУЩЕСТВУЕТ ТАК НАЗЫВАЕМЫЙ **ЕМ-АЛГОРИТМ**, ОСНОВАННЫЙ НА ФУНКЦИИ РАСПРЕДЕЛЕНИЯ ВЕРОЯТНОСТИ, КОТОРЫЙ ПОМОГАЕТ НЕ ТОЛЬКО РАЗДЕЛИТЬ ДАННЫЕ ПО КЛАСТЕРАМ, НО И ГЕНЕРИРОВАТЬ НОВЫЕ ДАННЫЕ В КАЖДОМ КЛАССЕ.



ФУНКЦИЯ РАСПРЕДЕЛЕНИЯ ВЕРОЯТНОСТИ? СЛОЖНОВАТО.

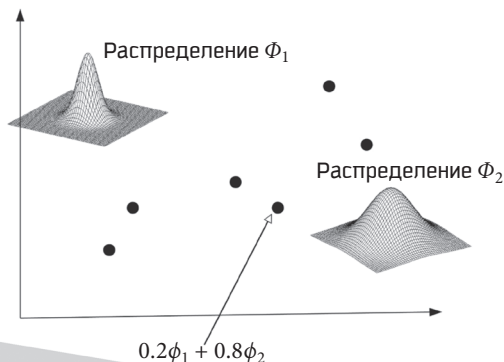


В ОЩЕМ, ЕЕ МОЖЕТ ЛЕГКО ОБЪЯСНИТЬ ГИПОТЕЗА РАСПРЕДЕЛЕНИЯ ПРАВИЛЬНЫХ ОТВЕТОВ.

ПОРЯДОК ТАКОЙ ЖЕ, КАК И В МЕТОДЕ *k*-СРЕДНИХ: СНАЧАЛА МЫ СТРОИМ СРЕДНИЕ ВЕКТОРЫ И НАХОДИМ МАТРИЦУ РАСПРЕДЕЛЕНИЯ. ЭТО ЭКВИВАЛЕНТНО ПОМЕЩЕНИЮ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ В СООТВЕТСТВУЮЩЕМ МЕСТЕ В ПРОСТРАНСТВЕ ПРИЗНАКОВ.

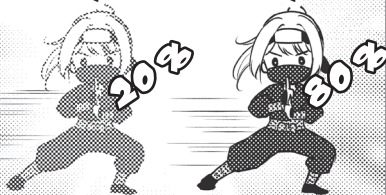


При перевычислении распределения добавляется вес 0,2



ЗАТЕМ С ПОМОЩЬЮ МЕТОДА *k*-СРЕДНИХ МЫ ОПРЕДЕЛЯЕМ, К КАКОМУ КЛАСТЕРУ ПРИНАДЛЕЖИТ КАЖДАЯ ЭЛЕМЕНТ ДАННЫХ. С ПОМОЩЬЮ ЕМ-АЛГОРИТМА МОЖНО НАЙТИ РАСПРЕДЕЛЕНИЕ: НАПРИМЕР, В ОДНОМ КЛАСТЕРЕ 20 %, А В ДРУГОМ 80 %.

БАМ! БАМ!



И КАК УДАЛОСЬ НАЙТИ ЭТО РАСПРЕДЕЛЕНИЕ?


ЕГО МОЖНО РАССЧИТАТЬ НА ОСНОВАНИИ ПОДХОДЯЩЕГО РАСПРЕДЕЛЕНИЯ ВСЕХ ДАННЫХ.

ЭТО КАК-ТО СЛУЧАЙНО!

НО СНАЧАЛА, ИСПОЛЬЗУЯ МЕТОД *k*-СРЕДНИХ, МЫ СЛУЧАЙНЫМ ОБРАЗОМ ОПРЕДЕЛЯЛИ СРЕДНИЙ ВЕКТОР И НА ОСНОВАНИИ ЕГО СТРОИЛИ КЛАСТЕРЫ, ПРАВАА?

НУ... АА.







ЗАТЕМ, КАК И В МЕТОДЕ k -СРЕДНИХ, ВЫЧИСЛЯЮТСЯ ПАРАМЕТРЫ КЛАСТЕРОВ, ТО ЕСТЬ СРЕДНИЙ ВЕКТОР И КОВАРИАЦИОННАЯ МАТРИЦА.

ТОГДА КОЛИЧЕСТВО ЭЛЕМЕНТОВ ДАННЫХ В КЛАСТЕРЕ БУДЕТ ВЫЧИСЛЯТЬСЯ НА ОСНОВАНИИ ТОГО, ЧТО БЫЛО НА ПРЕДЫДУЩЕМ ЭТАПЕ?


ВСЕ ТАК.



РАЗНИЦА ОЦЕНКИ ПЛОТНОСТИ ВЕРОЯТНОСТИ В МЕТОДЕ k -СРЕДНИХ И ЕМ-АЛГОРИТМЕ ВЫГЛЯДИТ ВОТ ТАК:



Метод k -средних	Случайным образом определяем k средних векторов	Формируем кластеры в зависимости от расстояния объекта данных до среднего вектора	Перевычисляем средние векторы каждого нового кластера
ЕМ-алгоритм	Случайным образом определяем k нормальных распределений	Рассчитываем вероятность, с которой каждый объект принадлежит к каждому кластеру, а затем определяем, где она выше	Пересчитываем параметры каждого распределения, рассматривая вероятность принадлежности каждого элемента (к кластеру) как вес



ТАКИМ ОБРАЗОМ МОЖНО СОЗДАТЬ КЛАСТЕР, КУДА ВХОДЯТ ЛЮДИ С ПОХОЖИМИ МОДЕЛЯМИ ПОВЕДЕНИЯ, И СОВЕТОВАТЬ ИМ ПОДХОДЯЩИЕ ДЛЯ НИХ СОБЫТИЯ.

6.2. РАЗЛОЖЕНИЕ МАТРИЦЫ

ПОГОВОРИМ О РАЗЛОЖЕНИИ МАТРИЦЫ.

ПОЧЕМУ ТЫ ХОТЕЛ ИСПОЛЬЗОВАТЬ ЭТОТ МЕТОД, ДЛЯ РЕКОМЕНДАЦИИ СОБЫТИЙ ПОЖИЛЫМ ЛЮДЯМ?



Я ОБНАРУЖИЛ, ЧТО ЭТО ЧАСТО ИСПОЛЬЗУЕТСЯ В РЕКОМЕНДАТЕЛЬНЫХ СИСТЕМАХ.

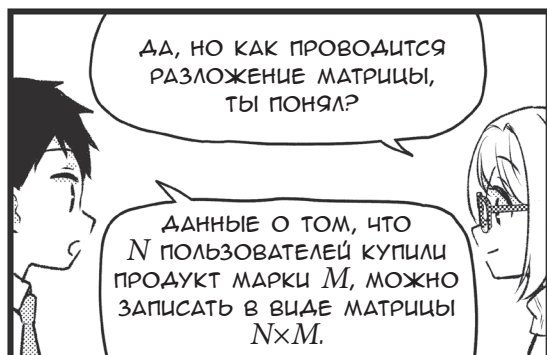
000 Куплю!

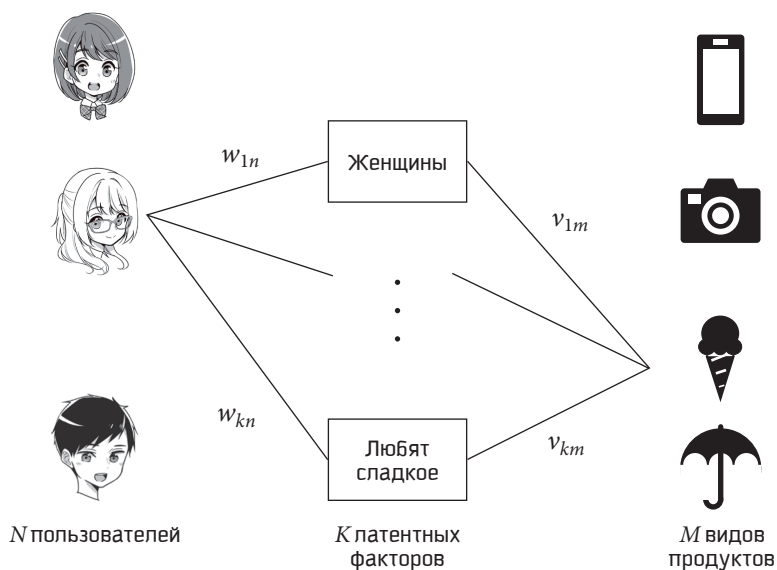


ЧТО ИЗ ЭТОГО?



НА ОСНОВАНИИ ИСТОРИИ ПОКУПОК ПОЛЬЗОВАТЕЛЯМ РЕКОМЕНДУЮТ ТОВАРЫ, КОТОРЫМИ ОНИ МОГУТ ЗАИНТЕРЕСОВАТЬСЯ.





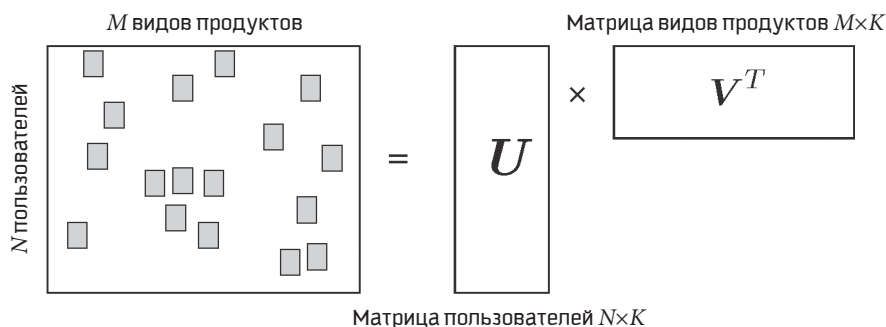
У нас есть K латентных факторов, и это число K меньше как количества пользователей N , так и товаров M .

Если мы будем следовать данной гипотезе, то величина элементов изначальной матрицы будет определяться по формуле ниже:

$$x_{nm} = w_1 v_{1m} + w_2 v_{2m} + \dots + w_K v_{Km}.$$

А как с ее помощью разложить матрицу?

Разложение матрицы при помощи этой формулы – на следующей странице.



Если матрица $N \times M$ большая, то информация о пользователях собирается в матрицу U , в которой N столбцов и K строк, а информация о товарах собирается в матрицу V , в которой M столбцов и K строк.



А зачем в целом U и V нужны?

Ну, если изначальная матрица – X , то можно сравнить X с UV^T и, обозначив ошибку за E , можно решить задачу оптимизации $E = X - UV^T$. Учитывая совместимость с методом градиента, который используется в качестве процедуры оптимизации, мы рассмотрим квадрат ошибки.



Для вычисления ошибки матрицы мы вычисляем квадратный корень из суммы квадратов каждого элемента матрицы. С использованием нормы Фробениуса оптимальный вариант выглядит так:



$$\min_{U,V} \frac{1}{2} \|E\|_{Fro}^2 = \min_{U,V} \frac{1}{2} \|X - UV^T\|_{Fro}^2$$

Это можно решить с использованием сингулярного разложения, которое изучается в линейной алгебре, но для этого надо заменить X на 0, если значение X не указано. Несмотря на то что никакой информации поначалу нет, далее она появится, и она будет отличаться от изначально введенных данных.



И что мы будем делать?

Так. Поскольку величина X существует, надо максимально уменьшить квадрат ошибки. Этим задача близка к регрессии. И, как и в задаче регрессии, тут появляются те же проблемы.



Переобучение. Тогда используем регуляризацию.

Да. Мы проводим оптимизацию с помощью алгоритма Alternating Least Squares.

$$\min_{U, V} \sum_{(i, j) \in \Omega} (x_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda_1 \|\mathbf{U}\|_{Fro}^2 + \lambda_2 \|\mathbf{V}\|_{Fro}^2$$

Здесь Ω – это индекс элемента матрицы X , \mathbf{u}_i – k -мерный вектор в i -й строке матрицы U , то же самое верно для \mathbf{v}_j – k -мерного вектора в j -й строке матрицы V .

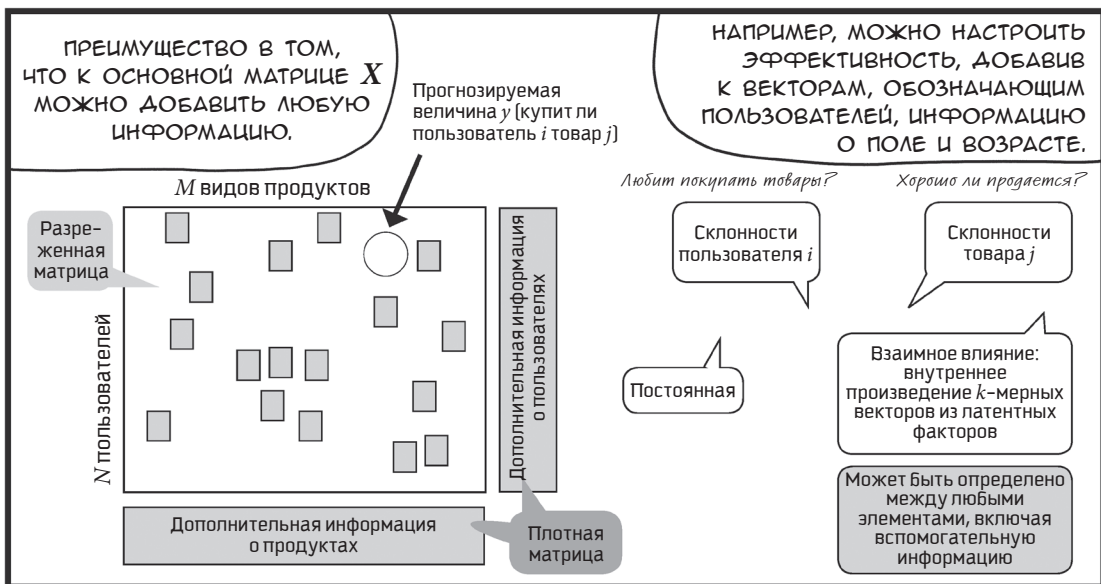


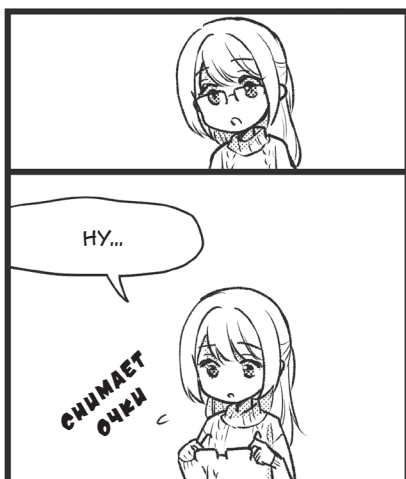
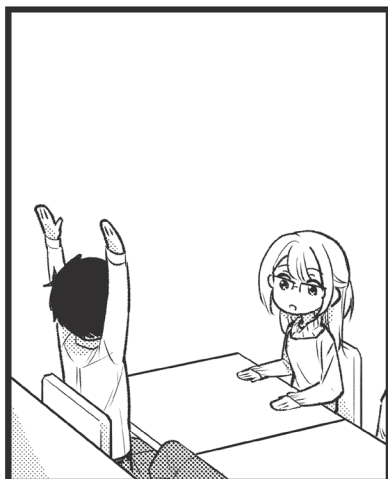
Ты говорил, что будешь использовать метод градиентов, но можно ли его использовать, если у нас два объекта оптимизации, U и V ?

U и V могут взаимно оптимизироваться, по всей видимости.











СЕГОДНЯ ПЛАТЫ ЗА ОБУЧЕНИЕ
НЕ БУДЕТ, МЫ РАЗДЕЛИМ СЧЕТ
ПОПОЛАМ!

ДА!

Спасибо

Хотя вы же меня многому научили сегодня

ДЭЫНЫ!

КСТАТИ, СЭМПАЙ,
У ВАС ЖЕ БЫЛИ ДРУГИЕ ОЧКИ.



ЭТИ МНЕ
ПОДАРИЛИ ДРУЗЬЯ
В ЧЕСТЬ НОВОЙ РАБОТЫ.

ЕЙ ОЧЕНЬ ЦАЕТ

В-А-АМ

ЕЙ ОЧЕНЬ ЦАЕТ

О-О-Ч-

ЕЙ ОЧЕНЬ ЦАЕТ

ОЧЕ-ЕНЬ

НА САМОМ ДЕЛЕ Я НЕ НАСТОЛЬКО
БЛИЗОРУКАЯ.... КОГДА Я НАДЕВАЮ ОЧКИ,
ТО БУДО ТО БЫ ПЕРЕКЛЮЧАЮСЬ НА РАБОТУ.

ХОРОШО, КОГДА
ЕСТЬ ТАКОЙ
ПЕРЕКЛЮЧАТЕЛЬ.

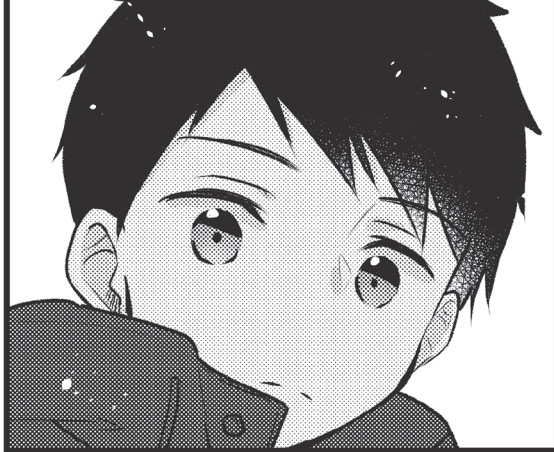
ЕЙ ОЧЕНЬ ЦАЕТ!

КИЁХАРА-КУН,
У ТЕБЯ ТОЖЕ ЕСТЬ ТАКОЙ
ПЕРЕКЛЮЧАТЕЛЬ?

ЧТО?

КОГДА ТЫ ХОЧЕШЬ СДЕЛАТЬ
ЧТО-ТО РАДИ КОГО-ТО ЕЩЕ!





ПРАВДА,
ЧТО ЛИ?

ДА.

ТЫ ВЕДЬ СТАРАЛСЯ НА
УРОКАХ ПО МАШИННОМУ
ОБУЧЕНИЮ, ПОТОМУ ЧТО
ЭТИ ПРОЕКТЫ БЫЛИ
НУЖНЫ КОМУ-ТО ЕЩЕ.

Возможен ли
у вас диалог?

ДА
ИЛИ
НЕТ
Выберите
лучший
ответ

НАЧАТЬ



ПОМНИШЬ, КАК
НА ФЕСТИВАЛЕ КУЛЬТУРЫ
ТЫ С ДРУЗЬЯМИ ЖАРИЛ
ЯКИСОБУ? ВЫ ОШИБЛИСЬ
В ПРОПОРЦИЯХ, И ЕЕ БЫЛО
НЕМНОГО БОЛЬШЕ,
ЧЕМ НУЖНО.

КОНЕЧНО!

ТЫ ОБЫЧНО НЕ ШЕВЕЛИШЬСЯ, НО ТОГДА ТЫ НАЧАЛ
ДЕЛАТЬ ОБЪЯВЛЕНИЯ: "ЛАПША ЗА ПОЛЦЕНЫ!",
"ПРИБЫЛЬ БУДЕТ ПОЖЕРТВОВАНА!"
ОБЪЯВЛЕНИЯ БЫЛИ И В СОЦСЕТЯХ,
ПОЭТОМУ ВЫ СМОГЛИ ВСЕ ПРОДАТЬ.



ДА... ПОМНЮ...

Я ИНОГДА ВОЛНУЮСЬ, МОЖЕТ
БЫТЬ, ТЫ НЕ ХОЧЕШЬ ЧТО-ТО
ДЕЛАТЬ ДЛЯ СЕБЯ?



НУ...

И ЭТО, МОЖЕТ БЫТЬ,
НЕ СОВСЕМ ТО, ЧТО НАДО,
НО МНЕ КАЖЕТСЯ, В ЭТОМ
И ЕСТЬ ТВОЯ Сильная
СТОРОНА.

Я ТАК ДУМАЮ
ИНОГДА.

СПАСИБО.



Ого...
драгоценные слова

И это награда за то,
что я старался весь год

НУ, ХВАТИТ, ПОЖАЛУЙ...

Я, НАВЕРНОЕ, НЕ СМОГУ
ПРОВОДИТЬ ВАС, СЭМПАЙ.

КАК РАЗ ЖЕ КОНЕЦ ГОДА...
СМОТРИ, НЕ ПЕРЕРЕБОТАЙ!

АГА!

НО ЕСЛИ ЧТО-ТО СЛУЧИТСЯ,
ПИШИ. ПУСТЬ ТЫ МНЕ БОЛЬШЕ
И НЕ УЧЕНИК, И НИКАКОЙ ПЛАТЫ
ЗА ОБУЧЕНИЕ ТОЖЕ НЕ НАДО!

КОНЕЧНО.



НУ ЧТО, ПОКА!

СЭМПАЙ!

ПРИЯТНОЙ ПОЕЗДКИ!

СПАСИБО!

ДО СВИДАНИЯ!

Математическое повторение (6)

В сегодняшней беседе меня заботит норма Фробениуса в матрице. Норма – это же величина вектора?



При правильном определении для d -мерного вектора x по отношению к p , если $1 \leq p < \infty$,

$$\sqrt[p]{|x_1|^p + \dots + |x_d|^p},$$

то это называется L_p -нормой x . Если $p = 2$, то норма L_2 обычно имеет смысл величины вектора.



Трудно. А зачем тогда нужно p ?



Представь, что $p = 1$. Разве ты не видела нигде L_1 -норму?



Если $p = 1$, то L_1 -норма – это сумма абсолютных значений всех элементов... Ага, если мы заменим x весом w , то получим дополнительный член лассо-регрессии!



Да. Изменение стандарта размера также меняет эффект, который влияет на размер. Ридж-регрессия с нормой L_2 в качестве члена регуляризации и лассо-регрессия с нормой L_1 в качестве члена регуляризации будут по-разному влиять на коэффициенты.





Если помнить это, то можно связать многие вещи.



Таким образом, L_2 -норма вектора – это квадратный корень из суммы квадратов всех его элементов. Точно так же норма Фробениуса матрицы – это квадратный корень из суммы квадратов всех элементов матрицы.

Я могу еще представить величину вектора, а величину матрицы – нет.



Нет нужды ее представлять. Надо просто понимать, что это все нужно только для того, чтобы минимизировать ошибку E до нулевой матрицы.

Я знала, что можно минимизировать ошибку способом, похожим на тот, что описан в главе 1, но не думала, что найду такой ответ.



Кстати, допустим, на сайте есть 100 пользователей и 50 товаров, и матрица будет размером 100×50 , то есть в ней будет 5000 элементов. А если мы добавим 10 латентных факторов, то количество элементов будет:

$$100 \times 10 + 50 \times 10 = 1500.$$

Но 5000 единиц информации не могут быть представлены 1500 единицами!





Ты говоришь правильно, но идея разложения матрицы не в том, чтобы можно было восстановить предыдущую матрицу, а сделать приближение более низкого ранга.



В случае если значение каждого элемента исходной матрицы полностью независимо от других элементов, этот метод не дает результата разложения, близкого к исходной матрице.



Но в данном случае предполагается, что пользователи с похожими характеристиками демонстрируют схожее покупательское поведение, или же пользователи, приобретающие продукты с аналогичными характеристиками, имеют сходные тенденции покупок. Говоря техническими терминами, предполагается, что в данных, у которых много измерений, есть структуры с более низкими измерениями.

Я, наверное, лучше промолчу...



Начнем с разложения по собственным числам. Для матрицы $d \times d$, т. е. для квадратной матрицы d -порядка M , рассмотрим пару из действительного числа λ и d -мерного вектора x , удовлетворяющих следующим условиям:

$$Ax = \lambda x, \quad x \neq 0.$$

Тогда формула преобразуется в $(A - \lambda I)x = 0$.



I – единичная матрица, если есть обратная матрица $A - \lambda I$, то $x = 0$, что противоречит условиям. Значит, обратной матрицы $A - \lambda I$ нет. В этом случае определитель равен 0, или же $\det(A - \lambda I)x = 0$.

Если у нас квадратная матрица второго порядка $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$, то ее определитель равен $ad - bc$. В первой главе мы узнали, как можно найти обратную матрицу $-1/(ad - bc)$, но если знаменатель равен 0, то найти обратную матрицу не получится!



Да. В случае квадратной матрицы d -порядка $\det(A - \lambda I)x = 0$ представляет собой многочлен и будет иметь d решений. λ – собственное число, а соответствующий вектор x называется собственным вектором.

Угу.



Тогда, используя собственные числа и собственные векторы, матрица M будет записана как произведение матриц по формуле ниже:

$$M = U \text{diag}(\lambda_1, \dots, \lambda_d) U^{-1},$$

где U – вектор-строка с d собственными векторами, которые находятся рядом, а diag – диагональная матрица, в которой упорядоченные числа расположены в диагональных элементах.



Мне пока понятно. Если в M d строк и d столбцов, а в U и U^{-1} – тоже d строк и d столбцов, то если выбрать правильные значения, можно провести преобразование.



А разве то, что M должна быть квадратной матрицей, – это не необычное условие? Разве, когда на сайте количество пользователей и товаров одинаковое, это не странно?



Конечно. Поэтому это разложение по собственным числам превращается в сингулярное разложение.

$$M = U\Sigma V^T.$$

Здесь M – матрица $n \times m$, U – матрица $n \times n$, V – матрица $m \times m$, Σ – матрица $n \times m$. Таким образом, при перемножении $U\Sigma V^T$ получится матрица $n \times m$.



Σ – не квадратная матрица, и диагональной матрицы из нее не получится.



Σ – матрица с r сингулярными числами, лежащими на главной диагонали (r меньше m и n), а оставшиеся элементы дополняются нулями до n строк и m столбцов.



Сингулярные числа рассчитываются по собственному значению, и здесь мы можем думать о них как о числах, полученных в результате разложения.

$$\Sigma = \begin{array}{|c|c|} \hline \begin{array}{c} \sigma_1 \\ \\ \cdot \\ \cdot \\ \cdot \\ \\ \sigma_r \end{array} & \begin{array}{c} \\ \\ \\ 0 \end{array} \\ \hline \begin{array}{c} 0 \end{array} & \begin{array}{c} 0 \end{array} \\ \hline \end{array}$$



Здесь важно знать, что сингулярные значения от σ_1 до σ_r расположены по возрастанию.



Каждое сингулярное значение умножается на элементы U и V , в результате чего получается элемент M , и, конечно, большая величина сингулярного значения оказывает большое влияние на определение значения M . Это сразу станет ясно, если представить, что первое сингулярное число настолько велико, что будто возвышается над остальными.

А если сингулярные значения со второго и ниже маленькие, то M будет меняться незначительно. Тогда они будут играть какую-то роль?



Да. Можно выбрать несколько больших сингулярных значений и затем смотреть, как меняется значение M .



Если сумма нескольких выбранных сингулярных значений составляет большую часть от суммы всех сингулярных значений, то получится матрица, которая не сильно отличается от исходной M .

Ага. Если мы выберем k элементов из больших сингулярных чисел, чтобы сделать Σ , то U будет матрицей $n \times k$, $V - t \times k$, а $\Sigma - k \times k$. Если k меньше, чем n или t , то произведением небольших матриц мы сможем получить большую.



Да. Как я и говорила, в данных, у которых много измерений, есть структуры с более низкими измерениями.

Теперь я поняла. Даже когда ты будешь в Токио, Саяка, то учи меня математике! (Хотя не представляю, как можно спокойно работать в Токио...)



ЭПИЛОГ

ЗНАНИЯ!
ПРАКТИКА!
ПРИМЕНЕНИЕ



НАЧАЛЬНИК!

ЧЕГО?

КУДА ЭТО СТАВИТЬ?

СЮДА ВНИЗ,
ПОЖАЛУЙСТА.

МОЖЕТ, ХВАТИТ МЕНЯ ТАК НАЗЫВАТЬ?

СТЕСНЯЕШЬСЯ?
ПРИБЫКАЙ, ПРИБЫКАЙ...

НУ, ТЫ ЖЕ СМОГ СОЗДАТЬ
КОМПАНИЮ, КИЁХАРА!

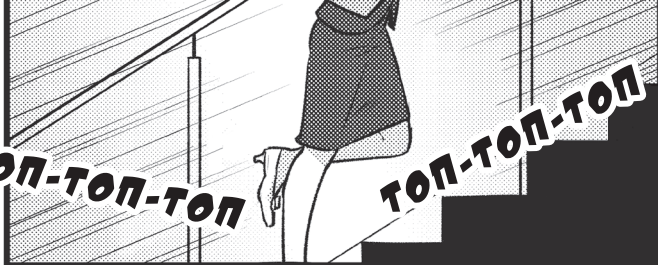
САМ УДИВЛЕН!

Количество проектов, связанных с машинным обучением, увеличилось, поэтому Киёхара больше не мог концентрироваться на работе в администрации и создал свою компанию.

НО Я НЕ ДУМАЛ, КУДЗЁ-САН,
ЧТО ВЫ ПОЙДАЕТЕ СО МНОЙ.

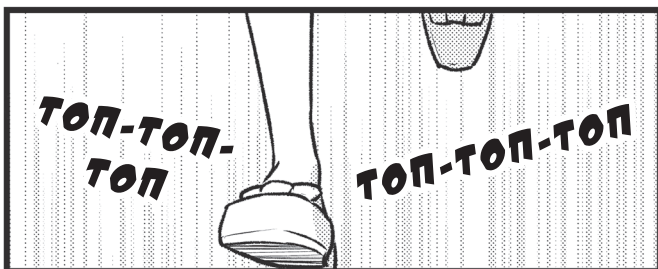
Я ПОДУМАЛ,
С КИЁХАРА СКУЧНО
НЕ БУДЕТ.

СПАСИБО БОЛЬШОЕ ТЕБЕ!



НАДО НЕ ТОЛЬКО МЕНЯ
БЛАГОДАРИТЬ, НО И ВСЕХ В МЭРЦИ,
ЧТО НАШЛИ ОФИС ПОДЕШЕВЛЕ
И ДАЛИ СТАРУЮ МЕБЕЛЬ.

СОГЛАСЕН.



ПОЙДУ-КА Я ПОРАБОТАЮ.

КОНЕЧНО!

ТОП-
ТОП-
ТОП

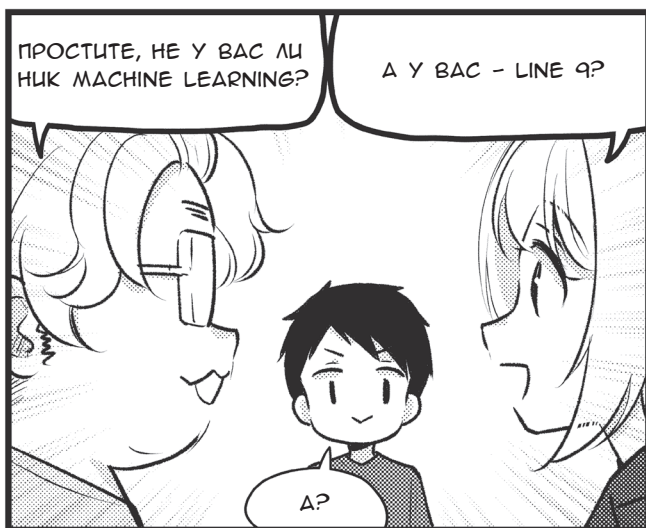
КИЁХАРА-КУН?

БУХ!

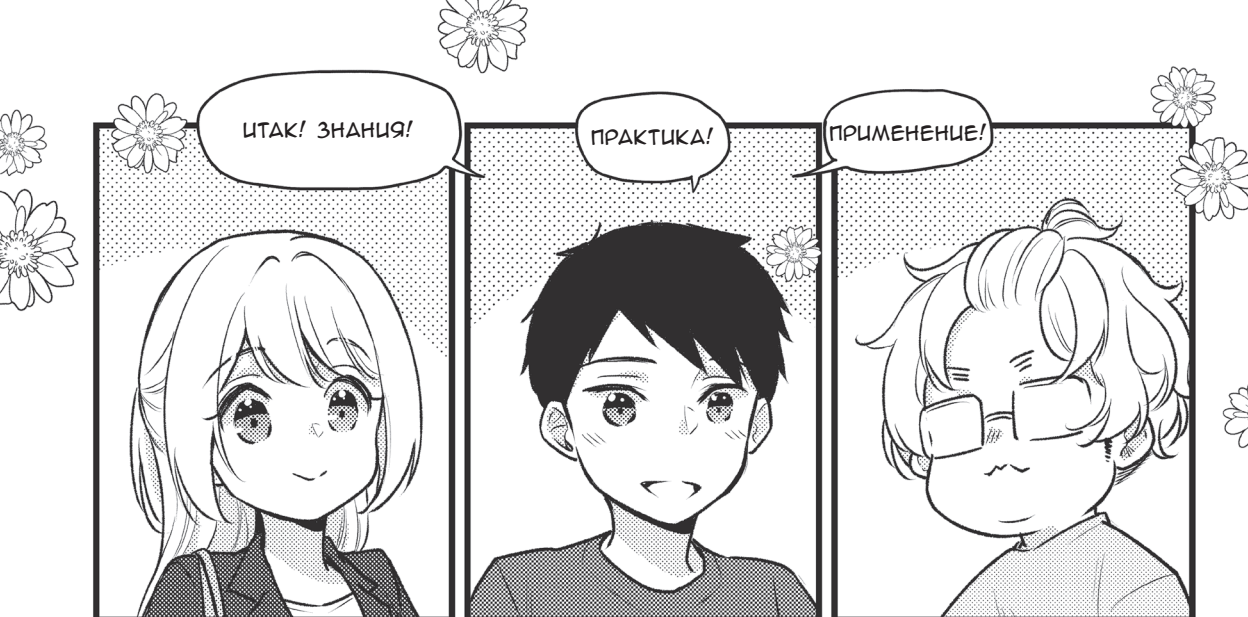












ИТАК! ЗНАНИЯ!

ПРАКТИКА!

ПРИМЕНЕНИЕ!

ПОСТАРАЕМСЯ
ИМ СЛЕДОВАТЬ!

АА!

АГА!

ДЛЯ НАЧАЛА
СДЕЛАЕМ ПРИЛИЧНЫЙ ОФИС
ИЗ ЭТОЙ КОМНАТЫ!

АА! Я ПОЙДУ ПЕРЕОДЕНУСЬ
В СПОРТИВНУЮ ОДЕЖДУ.

ВПЕРЕД!

ОТДОХНУ-КА Я!

КУАЗЁ-САН!

ТАК, КУЁХАРА, КОТОРЫЙ
НАШЕЛ РАБОТУ В КОМПАНИИ
СВОИХ ХОРОШИХ ДРУЗЕЙ,
СТАЛ ЗЛИТЬСЯ...

НЕ НАДО ТУТ
ЭТИХ РАЗГОВОРОВ!



ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ

Ada Boost.....	153	Метод кросс-валидации (перекрестная проверка).....	85
Alternating Least Squares, метод.....	183	Метод наименьших квадратов.....	20
data mining.....	10	Метод обратного распространения ошибки.....	107
dropout, прореживание.....	117	Метод предварительного обучения.....	113
EM-алгоритм.....	178	Метод проверки на зарезервированных данных.....	83
Factorization Machine.....	185	Независимая переменная.....	17
F-мера.....	92	Нейронная сеть.....	103
ID3-алгоритм.....	56	Нейронная сеть с прямым распространением.....	105
mini-batch метод.....	54	Неопределенность.....	59
Автокодировщик.....	11	Норма Фробениуса.....	182
Ансамблевые методы.....	145	Обучение без учителя.....	10
Бинарная классификация.....	48	Обучение с подкреплением.....	12
Блок линейной ректификации.....	116	Обучение с учителем.....	10
Бритва Оккама.....	65	Пакетный метод.....	54
Бустинг.....	152	Переобучение.....	65
Бэггинг.....	146	Полнота.....	91
Вес.....	18	Пороговый логический элемент.....	104
Входной слой.....	105	Правдоподобие.....	52
Выборка с возвращением.....	147	Проблема исчезновения градиента.....	112
Выходной слой.....	105	Промежуточные способы обучения.....	10
Глубокое обучение.....	11, 111	Разделяющая кластеризация.....	173
Градиентный бустинг.....	154	Разделяющая поверхность.....	50
Решающее дерево.....	55	Разложение матрицы.....	179
Дискретизация.....	66	Регрессия.....	16
Зависимая переменная.....	17	Регуляризация.....	23
Иерархическая кластеризация.....	173	Ридж-регрессия.....	23
Информационная энтропия.....	59	Сверточная нейронная сеть.....	119
Информационный выигрыш.....	64	Сигмоидная функция.....	51
Истинный пример.....	87	Скрытый слой.....	105
Квадрат ошибки.....	20	Слой пулинга.....	119
Классификация.....	47	Слой свертки.....	119
Кластеринг.....	172	Случайный лес.....	149
Количество правильно предсказанных ответов.....	88	Стохастический градиентный спуск.....	54
Контроль по отдельным объектам.....	86	Тестовый блок.....	85
Лассо-регрессия.....	23	Точность.....	90
Латентный фактор.....	180	Узел.....	55
Линейная регрессия.....	18	Усеченная линейная функция.....	116
Лист.....	55	Учительский сигнал.....	107
Логистическая классификация.....	49	Функция Softmax.....	106
Ложный пример.....	87	Функция активации.....	104
Матрица неточностей.....	88		
Машинное обучение.....	9		
Метод k -средних.....	176		

Книги издательства «ДМК ПРЕСС» можно купить оптом и в розницу в книготорговой компании «Галактика» (представляет интересы издательств «ДМК ПРЕСС», «СОЛОН ПРЕСС», «КТК Галактика»).

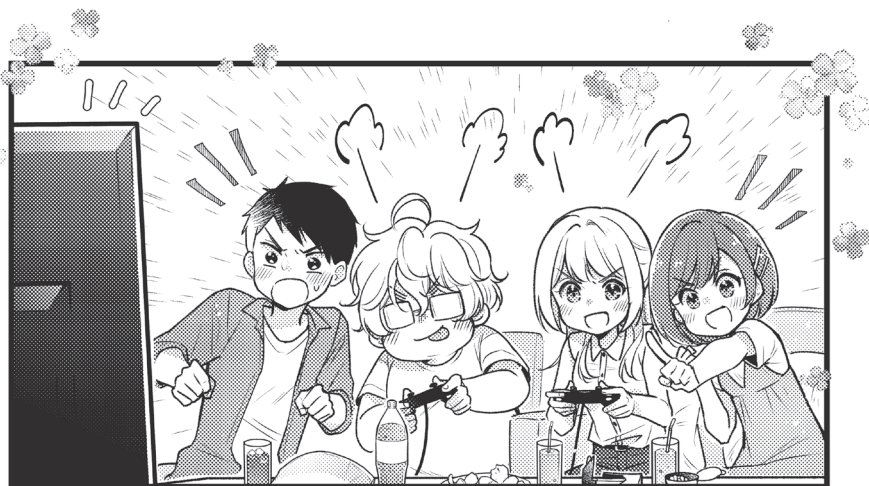
Адрес: г. Москва, пр. Андропова, 38;

тел.: **(499) 782-38-89**, электронная почта: **books@aliens-kniga.ru**.

При оформлении заказа следует указать адрес (полностью), по которому должны быть высланы книги; фамилию, имя и отчество получателя.

Желательно также указать свой телефон и электронный адрес.

Эти книги вы можете заказать и в интернет-магазине: **www.a-planet.ru**.



Араки Масахиро (автор), Ватари Макана (художник)

Занимательное программирование

Машинное обучение

Манга

Главный редактор *Д. А. Мовчан*

dmkpress@gmail.com

Научный редактор *М. Е. Петровичева*

Переводчик *А. С. Слащева*

Корректор *Г. И. Синяева*

Верстальщик *А. А. Чаннова*

Формат 70×100 1/16.

Гарнитура Anime Ace. Печать офсетная.

Усл. п. л. 17,39. Тираж 500 экз.

Веб-сайт издательства www.dmkpress.com