

Д. П. Ветров, Д. А. Кропотов

Байесовские методы машинного обучения
Учебное пособие

Пособие создано при поддержке программы «Формирование системы инновационного образования в МГУ им. М.В. Ломоносова»

Москва, 2007г

Цели курса

- Ознакомление с классическими методами обработки данных, особенностями их применения на практике и их недостатками
- Представление современных проблем теории машинного обучения
- Введение в байесовские методы машинного обучения
- Изложение последних достижений в области практического использования байесовских методов
- Напоминание основных результатов из смежных дисциплин (теория кодирования, анализ, матричные вычисления, статистика, линейная алгебра, теория вероятностей, случайные процессы)

Структура курса

- 1 семестр, 12 лекций, 24 аудиторных часа + 12 часов для самостоятельной работы
- В каждой лекции секция ликбеза, содержащая краткое напоминание полезных фактов из смежных областей математики
- В конце курса экзамен. Три вопроса в билете, один из секции ликбеза + задача
- Каждая лекция сопровождается показом презентации
- Методические материалы (включая презентации), а также большая часть рекомендуемой литературы доступна на сайте <http://mmphome.lgb.ru>
- Лекторы: Дмитрий Ветров (VetrovD@yandex.ru) и Дмитрий Кропотов (DKropotov@yandex.ru)

Оглавление

1	Различные задачи машинного обучения	3
1.1	Некоторые задачи машинного обучения	4
1.1.1	Задача классификации	4
1.1.2	Задача восстановления регрессии	5
1.1.3	Задача кластеризации (обучения без учителя)	6
1.1.4	Задача идентификации	6
1.1.5	Задача прогнозирования	7
1.1.6	Задача извлечения знаний	8
1.2	Основные проблемы машинного обучения	9
1.2.1	Малый объем обучающей выборки	9
1.2.2	Некорректность входных данных	10
1.2.3	Переобучение	10
1.3	Ликбез: Основные понятия мат. статистики	12
2	Вероятностная постановка задачи распознавания образов	14
2.1	Ликбез: Нормальное распределение	15
2.2	Статистическая постановка задачи машинного обучения	16
2.2.1	Вероятностное описание	16
2.2.2	Байесовский классификатор	17
2.3	Методы восстановления плотностей	18
2.3.1	Общие замечания	18
2.3.2	Парzenовские окна	19
2.3.3	Методы ближайшего соседа	19
2.4	EM-алгоритм	21
2.4.1	Параметрическое восстановление плотностей	21
2.4.2	Задача разделения смеси распределений	22
2.4.3	Разделение гауссовской смеси	23
3	Обобщенные линейные модели	25
3.1	Ликбез: Псевдообращение матриц и нормальное псевдорешение	26
3.2	Линейная регрессия	27
3.2.1	Классическая линейная регрессия	27
3.2.2	Метод наименьших квадратов	28
3.2.3	Вероятностная постановка задачи	29
3.3	Применение регрессионных методов для задачи классификации	30
3.3.1	Логистическая регрессия	30
3.3.2	Метод IRLS	31

4	Метод опорных векторов и безпризнаковое распознавание образов	34
4.1	Ликбез: Условная оптимизация	35
4.2	Метод опорных векторов для задачи классификации	37
4.2.1	Метод потенциальных функций	37
4.2.2	Случай линейно разделимых данных	38
4.2.3	Случай линейно неразделимых данных	43
4.2.4	Ядровой переход	44
4.2.5	Заключительные замечания	46
4.3	Метод опорных векторов для задачи регрессии	48
4.4	Безпризнаковое распознавание образов	50
4.4.1	Основная методика безпризнакового распознавания образов	50
4.4.2	Построение функции, задающей скалярное произведение	51
5	Задачи выбора модели	55
5.1	Ликбез: Оптимальное кодирование	56
5.2	Постановка задачи выбора модели	56
5.2.1	Общий характер проблемы выбора модели	56
5.2.2	Примеры задач выбора модели	57
5.3	Общие методы выбора модели	59
5.3.1	Кросс-валидация	59
5.3.2	Теория Вапника-Червоненкиса	61
5.3.3	Принцип минимальной длины описания	62
5.3.4	Информационные критерии	63
6	Байесовский подход к теории вероятностей. Примеры байесовских рассуждений	65
6.1	Ликбез: Формула Байеса	66
6.1.1	Sum- и Product- rule	66
6.1.2	Формула Байеса	67
6.2	Два подхода к теории вероятностей	67
6.2.1	Частотный подход	67
6.2.2	Байесовский подход	68
6.3	Байесовские рассуждения	69
6.3.1	Связь между байесовским подходом и булевой логикой	69
6.3.2	Пример вероятностных рассуждений	70
7	Решение задачи выбора модели по Байесу. Обоснованность модели	73
7.1	Ликбез: Бритва Оккама и Ad Hoc гипотезы	74
7.2	Полный байесовский вывод	74
7.2.1	Пример использования априорных знаний	74
7.2.2	Сопряженные распределения	75
7.2.3	Иерархическая схема Байеса	77
7.3	Принцип наибольшей обоснованности	77
7.3.1	Обоснованность модели	77
7.3.2	Примеры использования	79
8	Метод релевантных векторов	82
8.1	Ликбез: Матричные тождества обращения	83
8.2	Метод релевантных векторов для задачи регрессии	84
8.3	Метод релевантных векторов для задачи классификации	88

9	Недиагональная регуляризация обобщенных линейных моделей	94
9.1	Ликбез: Неотрицательно определенные матрицы и Лапласовское распределение	95
9.2	Метод релевантных собственных векторов	96
9.2.1	RVM и его ограничения	96
9.2.2	Регуляризация степеней свободы	97
9.2.3	Оптимизация обоснованности для различных семейств априорных распределений	98
10	Общее решение для недиагональной регуляризации	102
10.1	Ликбез: Дифференцирование по вектору и по матрице	103
10.2	Общее решение для недиагональной регуляризации	104
10.2.1	Получение выражения для обоснованности с произвольной матрицей регуляризации	104
10.2.2	Получение оптимальной матрицы регуляризации в явном виде	105
11	Методы оценки обоснованности	109
11.1	Ликбез: Дивергенция Кульбака-Лейблера и Гамма-распределение	110
11.2	Вариационный метод	111
11.2.1	Идея метода	111
11.2.2	Вариационная линейная регрессия	113
11.3	Методы Монте-Карло	115
11.3.1	Простейшие методы	115
11.3.2	Схема Метрополиса-Гиббса	116
11.3.3	Гибридный метод Монте-Карло	117
12	Графические модели. Гауссовские процессы в машинном обучении	119
12.1	Ликбез: Случайные процессы и условная независимость	120
12.1.1	Случайные процессы	120
12.1.2	Условная независимость	121
12.2	Графические модели	121
12.2.1	Ориентированные графы	121
12.2.2	Три элементарных графа	124
12.2.3	Неориентированные графы	125
12.3	Гауссовские процессы в машинном обучении	126
12.3.1	Гауссовские процессы в задачах регрессии	126
12.3.2	Гауссовские процессы в задачах классификации	128
12.3.3	Подбор ковариационной функции	129

Глава 1

Различные задачи машинного обучения

В главе рассматриваются различные постановки задачи машинного обучения и вводятся основные обозначения, используемые в последующих главах. Также приведены основные проблемы, возникающие при обучении ЭВМ.

1.1 Некоторые задачи машинного обучения

Концепция машинного обучения

- Решение задач путем обработки прошлого опыта (case-based reasoning)
- Альтернатива построению математических моделей (model-based reasoning)
- Основное требование – наличие обучающей информации
- Как правило в качестве таковой выступает выборка **прецедентов** – ситуационных примеров из прошлого с известным исходом
- Требуется построить алгоритм, который позволял бы обобщить опыт прошлых наблюдений/ситуаций для обработки новых, не встречавшихся ранее случаев, исход которых неизвестен.

1.1.1 Задача классификации

Классификация

- Исторически возникла из задачи машинного зрения, поэтому часто употребляемый синоним – распознавание образов
- В классической задаче классификации обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- В качестве исхода объекта \mathbf{x} фигурирует переменная t , принимающая конечное число значений, обычно из множества $\mathcal{T} = \{1, \dots, l\}$
- Требуется построить алгоритм (классификатор), который по вектору признаков \mathbf{x} вернул бы метку класса \hat{t} или вектор оценок принадлежности (апостериорных вероятностей) к каждому из классов $\{p(s|\mathbf{x})\}_{s=1}^l$ (см. рис. 1.1)

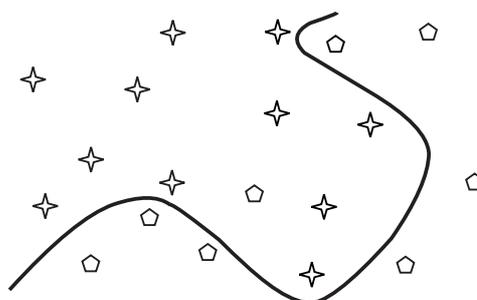


Рис. 1.1. Пример двухклассовой задачи классификации. Звездочками обозначены объекты из одного класса, пятиугольниками – объекты из другого класса. Черная линия соответствует разделяющей поверхности, обеспечивающей качественную классификацию новых объектов

Примеры задач классификации

- Медицинская диагностика: по набору медицинских характеристик требуется поставить диагноз
- Геологоразведка: по данным зондирования почв определить наличие полезных ископаемых

- Оптическое распознавание текстов: по отсканированному изображению текста определить цепочку символов, его формирующих
- Кредитный скоринг: по анкете заемщика принять решение о выдаче/отказе кредита
- Синтез химических соединений: по параметрам химических элементов спрогнозировать свойства получаемого соединения

1.1.2 Задача восстановления регрессии

Регрессия

- Исторически возникла при исследовании влияния одной группы непрерывных случайных величин на другую группу непрерывных случайных величин
- В классической задаче восстановления регрессии обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- В качестве исхода объекта \mathbf{x} фигурирует непрерывная вещественнозначная переменная t
- Требуется постросить алгоритм (регрессор), который по вектору признаков \mathbf{x} вернул бы точечную оценку значения регрессии \hat{t} , доверительный интервал (t_-, t_+) или апостериорное распределение на множестве значений регрессионной переменной $p(t|\mathbf{x})$



Рис. 1.2. Пример задачи восстановления регрессии. Звездочками обозначены прецеденты, черная линия показывает пример восстанавливаемой функции регрессии

Примеры задач восстановления регрессии

- Оценка стоимости недвижимости: по характеристике района, экологической обстановке, транспортной связности оценить стоимость жилья
- Прогноз свойств соединений: по параметрам химических элементов спрогнозировать температуру плавления, электропроводность, теплоемкость получаемого соединения
- Медицина: по постоперационным показателям оценить время заживления органа
- Кредитный скоринг: по анкете заемщика оценить величину кредитного лимита
- Инженерное дело: по техническим характеристикам автомобиля и режиму езды спрогнозировать расход топлива

1.1.3 Задача кластеризации (обучения без учителя)

Кластеризация

- Исторически возникла из задачи группировки схожих объектов в единую структуру (кластер) с последующим выявлением общих черт
- В классической задаче кластеризации обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- Требуется построить алгоритм (кластеризатор), который разбил бы выборку на непересекающиеся группы (кластеры) $X = \bigcup_{j=1}^k C_k$, $C_j \subset \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, $C_i \cap C_j = \emptyset$
- В каждый класс должны попасть объекты в некотором смысле похожие друг на друга (см. рис. 1.3)

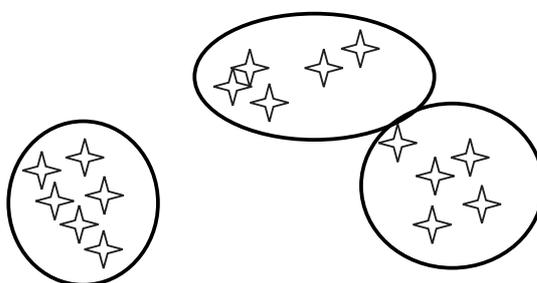


Рис. 1.3. Пример задачи кластеризации. Звездочками обозначены прецеденты. Группы объектов, обведенные кружками, образуют отдельные кластеры

Примеры задач кластерного анализа

- Экономическая география: по физико-географическим и экономическим показателям разбить страны мира на группы схожих по экономическому положению государств
- Финансовая сфера: по сводкам банковских операций выявить группы «подозрительных», нетипичных банков, сгруппировать остальные по степени близости проводимой стратегии
- Маркетинг: по результатам маркетинговых исследований среди множества потребителей выделить характерные группы по степени интереса к продвигаемому продукту
- Социология: по результатам социологических опросов выявить группы общественных проблем, вызывающих схожую реакцию у общества, а также характерные фокус-группы населения

1.1.4 Задача идентификации

Идентификация

- Исторически возникла из классификации, необходимости отделить объекты, обладающие определенным свойством, от «всего остального»
- В классической задаче идентификации обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$, обладающих некоторым свойством $\chi_A(\mathbf{x}) = 1$
- Особенностью задачи является то, что все объекты принадлежат одному классу, причем не существует возможности сделать репрезентативную выборку из класса «все остальное»

- Требуется построить алгоритм (идентификатор), который по вектору признаков \mathbf{x} определил бы наличие свойства A у объекта \mathbf{x} , либо вернул оценку степени его выраженности $p(\chi_A(\mathbf{x}) = 1|\mathbf{x})$ (см. рис. 1.4)

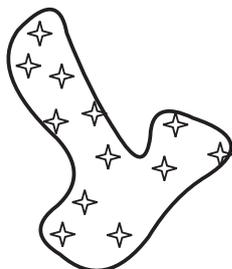


Рис. 1.4. Пример задачи идентификации. Объекты, обладающие определенным свойством, отделены от всех остальных объектов

Примеры задач идентификации

- Медицинская диагностика: по набору медицинских характеристик требуется установить наличие/отсутствие конкретного заболевания
- Системы безопасности: по камерам наблюдения в подъезде идентифицировать жильца дома
- Банковское дело: определить подлинность подписи на чеке
- Обработка изображений: выделить участки с изображениями лиц на фотографии
- Искусствоведение: по характеристикам произведения (картины, музыки, текста) определить, является ли его автором тот или иной автор

1.1.5 Задача прогнозирования

Прогнозирование

- Исторически возникла при исследовании временных рядов и попытке предсказания их значений через какой-то промежуток времени
- В классической задаче прогнозирования обучающая выборка представляет собой набор измерений $X = \{\mathbf{x}[i]\}_{i=1}^n$, представляющих собой вектор вещественнозначных величин $\mathbf{x}[i] = (x_1[i], \dots, x_d[i])$, сделанных в определенные моменты времени
- Требуется построить алгоритм (предиктор), который вернул бы точечную оценку $\{\hat{\mathbf{x}}[i]\}_{i=n+1}^{n+q}$, доверительный интервал $\{(\mathbf{x}_-[i], \mathbf{x}_+[i])\}_{i=n+1}^{n+q}$ или апостериорное распределение $p(\mathbf{x}[n+1], \dots, \mathbf{x}[n+q]|\mathbf{x}[1], \dots, \mathbf{x}[n])$ прогноза на заданную глубину q (см. рис. 1.5)
- В отличие от задачи восстановления регрессии, здесь осуществляется прогноз **по времени**, а не по признакам

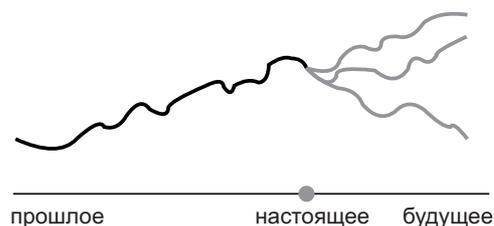


Рис. 1.5. Пример задачи прогнозирования. Черная кривая представляет собой предысторию (известное поведение характеристики). Серые кривые показывают различные варианты прогнозирования поведения характеристики в будущем

Примеры задач прогнозирования

- Биржевое дело: прогнозирование биржевых индексов и котировок
- Системы управления: прогноз показателей работы реактора по данным телеметрии
- Экономика: прогноз цен на недвижимость
- Демография: прогноз изменения численности различных социальных групп в конкретном ареале
- Гидрометеорология: прогноз геомагнитной активности

1.1.6 Задача извлечения знаний

Извлечение знаний

- Исторически возникла при исследовании взаимозависимостей между косвенными показателями одного и того же явления
- В классической задаче извлечения знаний обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- Требуется постросить алгоритм, генерирующий набор объективных закономерностей между признаками, имеющих место в генеральной совокупности (см. рис. 1.6)
- Закономерности обычно имеют форму предикатов «ЕСЛИ ... ТО ...» и могут выражаться как в цифровых терминах $((0.45 \leq x_4 \leq 32.1) \& (-6.98 \leq x_7 \leq -6.59) \Rightarrow (3.21 \leq x_2 \leq 3.345))$, так и в текстовых («ЕСЛИ Давление – низкое И (Реакция – слабая ИЛИ Реакция – отсутствует) ТО Пульс – нитевидный»)

Примеры задач извлечения знаний

- Медицина: поиск взаимосвязей (синдромов) между различными показателями при фиксированной болезни
- Социология: определение факторов, влияющих на победу на выборах
- Генная инженерия: выявление связанных участков генома
- Научные исследования: получение новых знаний об исследуемом процессе
- Биржевое дело: определение закономерностей между различными биржевыми показателями

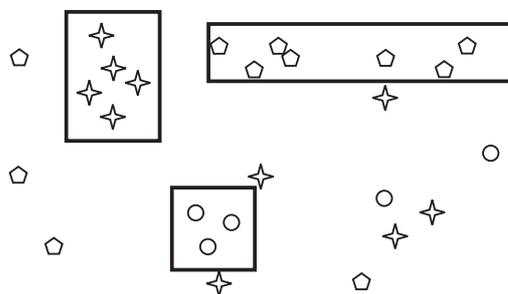


Рис. 1.6. Пример задачи извлечения знаний. В выборке представлены объекты из трех классов. Закономерности представляют собой области признакового пространства, в которых концентрация объектов из одного класса существенно превалирует над концентрациями объектов из других классов

1.2 Основные проблемы машинного обучения

1.2.1 Малый объем обучающей выборки

Объем выборки I

- Основным объектом работы любого метода машинного обучения служит обучающая выборка
- Большой объем выборки позволяет
 - Получить более надежные результаты
 - Использовать более сложные модели алгоритмов
 - Оценить точность обучения
 - **НО:** Время обучения быстро растет
- При малых выборках
 - Можно использовать **только** простые модели алгоритмов
 - Скорость обучения максимальна – можно использовать методы, требующие много времени на обучение
 - Высока вероятность переобучения при ошибке в выборе модели

Объем выборки II

- Одна и та же выборка может являться большой для простых моделей алгоритмов и малой для сложных моделей.
- Для методов с т.н. бесконечной емкостью Вапника-Червоненкиса **любая** выборка является малой.
- С ростом числа признаков увеличивается количество объектов, необходимое для корректного анализа данных
- Часто рассматривается т.н. эффективная размерность выборки $\frac{n}{d}$
- При объемах данных порядка десятков и сотен тысяч встает проблема уменьшения выборки с сохранением ее репрезентативности (active learning)

1.2.2 Некорректность входных данных

Неполнота признакового описания

- Отдельные признаки могут отсутствовать у некоторых объектов. Это может быть связано с отсутствием данных об измерении данного признака для данного объекта, а может быть связано с принципиальным отсутствием данного свойства у данного объекта
- Такое часто встречается в медицинских и химических данных
- Необходимы специальные процедуры, позволяющие корректно обрабатывать пропуски в данных
- Одним из возможных способов такой обработки является замена пропусков на среднее по выборке значение данного признака
- По возможности, пропуски следует игнорировать и исключать из рассмотрения при анализе соответствующего объекта

Противоречивость данных

- Объекты с одним и тем же признаковым описанием могут иметь разные исходы (принадлежать к разным классам, иметь отличные значения регрессионной переменной и т.п.)
- Многие методы машинного обучения не могут работать с такими наборами данных
- Необходимо заранее исключать или корректировать противоречащие объекты
- Использование вероятностных методов обучения позволяет корректно обрабатывать противоречивые данные
- При таком подходе предполагается, что исход t для каждого признакового описания \mathbf{x} есть случайная величина, имеющая некоторое условное распределение $p(t|\mathbf{x})$

Разнородность признаков

- Хотя формально предполагается, что признаки являются вещественнозначными, они могут быть дискретными и номинальными
- Номинальные признаки отличаются особенностями метрики между значениями
- Стандартная практика состоит в замене номинальных признаков на набор бинарных переменных по числу значений номинального признака
- Текстовые признаки, признаки-изображения, даты и пр. необходимо заменить на соответствующие номинальные либо числовые значения

1.2.3 Переобучение

Идея машинного обучения

- Задача машинного обучения заключается в восстановлении зависимостей по конечным выборкам данных (прецедентов)
- Пусть $(X, \mathbf{t}) = (\mathbf{x}_i, t_i)_{i=1}^n$ – обучающая выборка, где $\mathbf{x}_i \in \mathbb{R}^d$ – признаковое описание объекта, а $t \in \mathcal{T}$ – значение скрытой компоненты (классовая принадлежность, значение прогноза, номер кластера и т.д.)
- При статистическом подходе к решению задачи МО предполагается, что **обучающая выборка является выборкой из некоторой генеральной совокупности** с плотностью $p(\mathbf{x}, t)$
- Требуется восстановить $p(t|\mathbf{x})$, т.е. знание о скрытой компоненте объекта по измеренным признакам

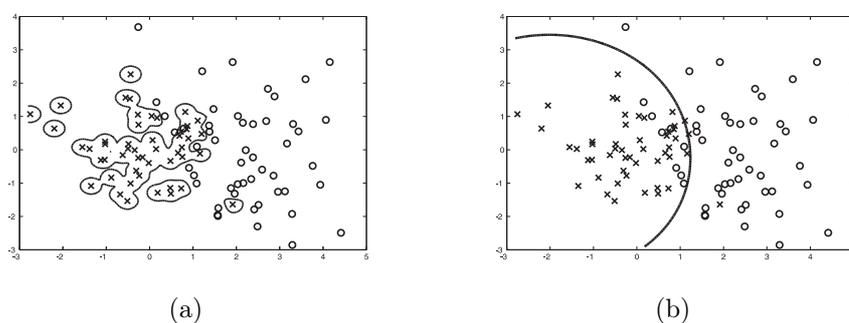


Рис. 1.7. Пример двухклассовой задачи классификации. На рисунке (a) представлено решающее правило, которое способно объяснить только объекты обучающей выборки. Решающее правило на рисунке (b) улавливает общую тенденцию в данных и обладает более высокой обобщающей способностью, чем решающее правило (a)

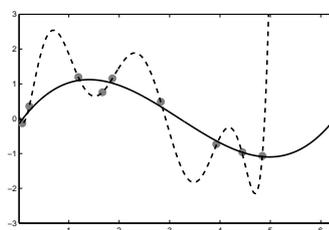


Рис. 1.8. Пример задачи восстановления регрессии. Пунктирная функция регрессии в точности предсказывает объекты обучающей выборки, однако обладает слабой экстраполирующей способностью. Функция регрессии, представленная черной линией, не так точно объясняет объекты обучения, однако хорошо улавливает общую тенденцию в данных

Проблема переобучения

Прямая минимизация невязки на обучающей выборке ведет к получению решающих правил, способных объяснить все что угодно и найти закономерности даже там, где их нет (см. рис. 1.7 и 1.8).

Способы оценки и увеличения обобщающей способности

- На сегодняшний день единственным универсальным способом оценивания обобщающей способности является кросс-валидация
- Все попытки предложить что-нибудь отличное от метода проб и ошибок пока **не привели к общепризнанному решению**. Наиболее известны из них следующие:
 - Структурная минимизация риска (В. Вапник, А. Червоненкис, 1974)
 - Минимизация длины описания (Дж. Риссанен, 1978)
 - Информационные критерии Акаике и Байеса-Шварца (Акаике, 1974, Шварц, 1978)
 - Максимизация обоснованности (МакКай, 1992)
- Последний принцип позволяет надеяться на конструктивное решение задачи выбора модели

Примеры задач выбора модели

- Определение числа кластеров в данных

- Выбор коэффициента регуляризации в задаче машинного обучения (например, коэффициента затухания весов (weight decay) в нейронных сетях)
- Установка степени полинома при интерполяции сплайнами
- Выбор наилучшей ядерной функции в методе опорных векторов (SVM)
- Определение количества ветвей в решающем дереве
- и многое другое...

1.3 Ликбез: Основные понятия мат. статистики

Краткое напоминание основных вероятностных понятий

- $X : \Omega \rightarrow \mathbb{R}$ – случайная величина
- Вероятность попадания величины в интервал (a, b) равна

$$P(a \leq X \leq b) = \int_a^b p(x) dx,$$

где $p(x)$ – плотность распределения X ,

$$p(x) \geq 0, \quad \int_{-\infty}^{\infty} p(x) dx = 1$$

- Если поведение случайной величины определяется некоторым параметром, возникают условные плотности $p(x|\theta)$. Если рассматривать условную плотность как функцию от параметра

$$f(\theta) = p(x|\theta),$$

то принято говорить о т.н. функции правдоподобия

Основная задача мат. статистики

- Распределение случайной величины X известно с точностью до параметра θ
- Имеется выборка значений величины X , $\mathbf{x} = (x_1, \dots, x_n)$
- Требуется оценить значение θ
- Метод максимального правдоподобия

$$\hat{\theta}_{ML} = \arg \max f(\theta) = \arg \max p(\mathbf{x}|\theta) = \arg \max \prod_{i=1}^n p(x_i|\theta)$$

- Можно показать, что ММП является асимптотически оптимальным при $n \rightarrow \infty$
- Увы, мир несовершенен. Величина n конечна и обычно не слишком велика
- Необходима регуляризация метода

Пример некорректного использования метода максимального правдоподобия

- $X \sim w_1 \mathcal{N}(x|\mu_1, \sigma_1^2) + \dots + w_m \mathcal{N}(x|\mu_m, \sigma_m^2)$
- Необходимо определить $\theta = (m, \mu_1, \sigma_1^2, \dots, \mu_m, \sigma_m^2, w_1, \dots, w_m)$
- Применяем ММП

$$p(\mathbf{x}|\theta) = \prod_{i=1}^n p(x_i|\theta) = \prod_{i=1}^n \sum_{j=1}^m \frac{w_j}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{\|x_i - \mu_j\|^2}{2\sigma_j^2}\right) \rightarrow \max_{\theta}$$

- Решение

$$\hat{m}_{ML} = n, \quad \hat{\mu}_{j,ML} = x_j, \quad \hat{\sigma}_{j,ML}^2 = 0, \quad \hat{w}_{ML,j} = \frac{1}{n}$$

Выводы

- Не все параметры можно настраивать в ходе обучения
- Существуют специальные параметры (будем называть их структурными), которые должны быть зафиксированы до начала обучения
- *!! В данном случае величина m (количество компонент смеси) является структурным параметром!!*
- Основной проблемой машинного обучения является проблема выбора структурных параметров, позволяющих избегать переобучения

Глава 2

Вероятностная постановка задачи распознавания образов

В главе вводится статистическая постановка задачи машинного обучения. Доказывается оптимальность байесовского классификатора. Подробно рассматриваются вопросы восстановления плотности распределения по выборке данных. Описываются несколько простых подходов к непараметрическому оцениванию плотности. В конце главы приведена общая схема EM-алгоритма, позволяющего восстанавливать смеси распределений с помощью введения скрытой (латентной) переменной.

2.1 Ликбез: Нормальное распределение

Нормальное распределение

- Нормальное распределение играет важнейшую роль в математической статистике

$$X \sim \mathcal{N}(x|\mu, \sigma^2) \Leftrightarrow p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\mu = \mathbb{E}X, \quad \sigma^2 = \mathbb{D}X \triangleq \mathbb{E}(X - \mathbb{E}X)^2$$

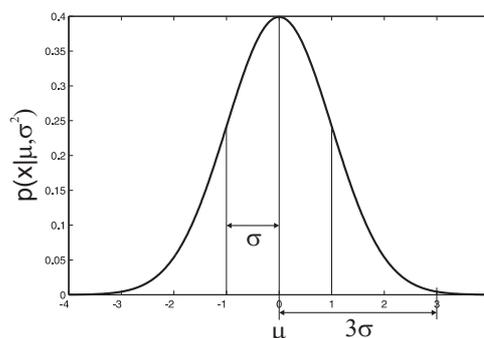


Рис. 2.1. Функция плотности нормального распределения

- Из центральной предельной теоремы следует, что сумма независимых случайных величин с ограниченной дисперсией стремится к нормальному распределению
- На практике, многие случайные величины можно считать приближенно нормальными

Многомерное нормальное распределение

- Многомерное нормальное распределение имеет вид

$$X \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) \Leftrightarrow p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi}^n \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

где $\boldsymbol{\mu} = \mathbb{E}X$, $\Sigma = \mathbb{E}(X - \boldsymbol{\mu})(X - \boldsymbol{\mu})^T$ — вектор математических ожиданий каждой из n компонент и матрица ковариаций соответственно

- Матрица ковариаций показывает, насколько сильно связаны (коррелируют) компоненты многомерного нормального распределения

$$\Sigma_{ij} = \mathbb{E}(X_i - \mu_i)(X_j - \mu_j) = \text{Cov}(X_i, X_j)$$

- Если мы поделим ковариацию на корень из произведений дисперсий, то получим коэффициент корреляции

$$\rho(X_i, X_j) \triangleq \frac{\text{Cov}(X_i, X_j)}{\sqrt{\mathbb{D}X_i \mathbb{D}X_j}} \in [-1, 1]$$

Особенности нормального распределения

- Нормальное распределение **полностью задается** первыми двумя моментами (мат. ожидание и матрица ковариаций/дисперсия)
- Матрица ковариаций неотрицательно определена, причем на диагоналях стоят дисперсии соответствующих компонент
- Нормальное распределение имеет очень легкие хвосты: большие отклонения от мат. ожидания практически невозможны. Это обстоятельство нужно учитывать при приближении произвольных случайных величин нормальными

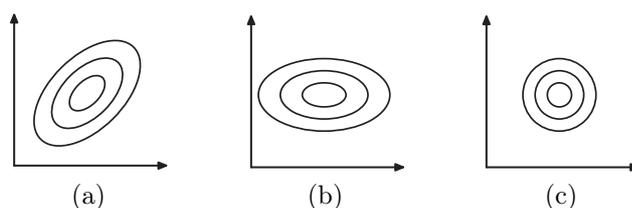


Рис. 2.2. Линии уровня для двухмерного нормального распределения. На рисунке (a) показано нормальное распределение с произвольной матрицей ковариации, случай (b) соответствует диагональной матрице ковариации, а случай (c) соответствует единичной матрице ковариации, умноженной на некоторую константу

2.2 Статистическая постановка задачи машинного обучения

2.2.1 Вероятностное описание

Основные обозначения

- В дальнейшем будут рассматриваться преимущественно задачи классификации и восстановления регрессии
- В этих задачах обучающая выборка представляет собой набор отдельных объектов $X = \{\mathbf{x}_i\}_{i=1}^n$, характеризующихся вектором вещественнозначных признаков $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- Каждый объект также обладает скрытой переменной $t \in \mathcal{T}$
- Предполагается, что существует зависимость между признаками объекта и значением скрытой переменной
- Для объектов обучающей выборки значение скрытой переменной известно $\mathbf{t} = \{t_i\}_{i=1}^n$

Статистическая постановка задачи

- Каждый объект описывается парой (\mathbf{x}, t)
- При статистической (вероятностной) постановке задачи машинного обучения предполагается, что **обучающая выборка является набором независимых, одинаково распределенных случайных величин, взятых из некоторой генеральной совокупности**
- В этом случае уместно говорить о плотности распределения объектов $p(\mathbf{x}, t)$ и использовать вероятностные термины (математическое ожидание, дисперсия, правдоподобие) для описания и решения задачи
- Заметим, что это не единственная возможная постановка задачи машинного обучения

Качество обучения

- Качество обучения определяется точностью прогноза на генеральной совокупности
- Пусть $S(t, \hat{t})$ – функция потерь, определяющая штраф за прогноз \hat{t} при истинном значении скрытой переменной t
- Разумно ожидать, что минимум этой функции достигается при $\hat{t} = t$
- Примерами могут служить $S_r(t, \hat{t}) = (t - \hat{t})^2$ для задачи восстановления регрессии и $S_c(t, \hat{t}) = I\{\hat{t} \neq t\}$ для задачи классификации

Абсолютный критерий качества

- Если бы функция $p(\mathbf{x}, t)$ была известна, задачи машинного обучения не существовало
- В самом деле абсолютным критерием качества обучения является мат. ожидание функции потерь, взятое по генеральной совокупности

$$\mathbb{E}S(t, \hat{t}) = \int S(t, \hat{t}(\mathbf{x}))p(\mathbf{x}, t)d\mathbf{x}dt \rightarrow \min,$$

где $\hat{t}(\mathbf{x})$ – решающее правило, возвращающее величину прогноза для вектора признаков \mathbf{x}

- Вместо методов машинного обучения сейчас бы активно развивались методы оптимизации и взятия интегралов от функции потерь
- Так как распределение объектов генеральной совокупности неизвестно, то абсолютный критерий качества обучения не может быть подсчитан

2.2.2 Байесовский классификатор

Идеальный классификатор

- Итак, одна из основных задач теории машинного обучения — это разработка способов косвенного оценивания качества решающего правила и выработка новых критериев для оптимизации в ходе обучения
- Рассмотрим задачу классификации с функцией потерь вида $S_c(t, \hat{t}) = I\{\hat{t} \neq t\}$ и гипотетический классификатор $t_B(\mathbf{x}) = \arg \max_{t \in \mathcal{T}} p(\mathbf{x}, t)$
- Справедлива следующая цепочка неравенств

$$\begin{aligned} \mathbb{E}S(t, \hat{t}) &= \int \int S(t, \hat{t}(\mathbf{x}))p(\mathbf{x}, t)d\mathbf{x}dt = \\ &= \sum_{s=1}^l \int S(s, \hat{t}(\mathbf{x}))p(\mathbf{x}, s)d\mathbf{x} = 1 - \int p(\mathbf{x}, \hat{t}(\mathbf{x}))d\mathbf{x} \geq \\ &\geq 1 - \int \max_t p(\mathbf{x}, t)d\mathbf{x} = 1 - \int p(\mathbf{x}, t_B(\mathbf{x}))d\mathbf{x} = \mathbb{E}S(t, t_B) \end{aligned}$$

Особенности байесовского классификатора

- Таким образом, знание распределения объектов генеральной совокупности приводит к получению оптимального классификатора **в явной форме**
- Такой оптимальный классификатор называется байесовским классификатором
- Если бы удалось с высокой точностью оценить значение плотности генеральной совокупности в каждой точке пространства, задачу классификации можно было бы считать решенной
- На этом основан один из существующих подходов к машинному обучению

2.3 Методы восстановления плотностей

2.3.1 Общие замечания

Идея методов восстановления плотностей

- Восстановление плотностей является примером задачи **непараметрической** статистики
- Основной идеей методов восстановления плотностей является учет количества объектов обучающей выборки, попавших в некоторую окрестность рассматриваемой точки
- Чем больше объектов находится в окрестности рассматриваемой точки, тем выше значение плотности в этой точке
- Окрестность можно зафиксировать либо по ширине, либо по числу объектов, в нее попавших. Возникают два подхода к задаче восстановления плотности

Гистограммы

- Простейшим подходом к восстановлению плотностей является построение гистограмм
- Разбиваем область значений переменной на k областей (обычно прямоугольных) $\Delta_1 \dots \Delta_k$ и считаем сколько точек n_i попало в каждую область Δ_i (см. рис. 2.3). Оценка плотности в этом случае

$$\hat{p}(\mathbf{x}) = \frac{n_i}{n|\Delta_i|} I\{\mathbf{x} \in \Delta_i\}$$

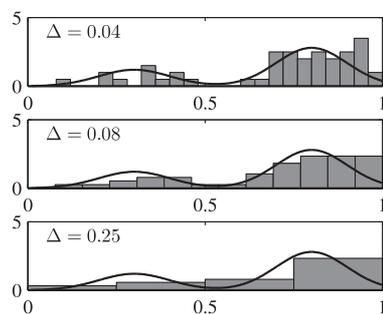


Рис. 2.3. Пример восстановления одномерной плотности с помощью гистограммы для различных значений ширины областей Δ

Недостатки гистограммного подхода

- Получающаяся оценка плотности не является непрерывной функцией
- Оценка зависит от выбора узлов гистограммы (центров областей Δ_i)
- Оценка сильно зависит от выбора ширин областей Δ_i
- **!!Ширина области является структурным параметром!!**
- С ростом размерности пространства d вычислительная сложность возрастает как k^d

2.3.2 Парзеновские окна

Парзеновское оценивание

- Рассмотрим достаточно малую область пространства \mathcal{D} , содержащую рассматриваемую точку \mathbf{x} . Вероятность попадания в эту область равна $P = \int_{\mathcal{D}} p(\mathbf{x}) d\mathbf{x}$. Из n точек обучающей выборки, в эту область попадет k точек, приблизительно равное $k \approx nP$
- Считая, что область \mathcal{D} достаточно мала, и плотность в ней постоянна, получаем $P \approx p(\mathbf{x})V$, где V – объем области \mathcal{D} , т.е.

$$\hat{p}(\mathbf{x}) = \frac{k}{nV}$$

- Отметим две существенные детали
 - Область \mathcal{D} должна быть достаточно широка, чтобы можно было использовать формулу для приближенного оценивания k , т.е. в область \mathcal{D} должно попадать достаточно много точек
 - Область \mathcal{D} должна быть достаточно узка, чтобы значение плотности в ней можно было приблизить константой

Парзеновское окно

- Рассмотрим функцию $k(\mathbf{u}) = k(-\mathbf{u}) \geq 0$, такую что $\int_{-\infty}^{\infty} k(\mathbf{u}) d\mathbf{u} = 1$. Такая функция называется **парзеновским окном**
- Тогда плотность может быть приближена как $\hat{p}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x} - \mathbf{x}_i)$
- В самом деле, легко показать, что $\hat{p}(\mathbf{x}) \geq 0$ и

$$\int \hat{p}(\mathbf{x}) d\mathbf{x} = \frac{1}{n} \sum_{i=1}^n \int k(\mathbf{x} - \mathbf{x}_i) d\mathbf{x} = 1$$

- Обычно в качестве парзеновских окон берут колоколообразные функции с максимумом в нуле, например, $k(\mathbf{u}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\mathbf{u}^T \mathbf{u}}{2}\right)$
- Парзеновские окна зависят от выбранного масштаба, т.е. характерной ширины окрестности: если $k(\mathbf{u})$ окно, то $\frac{1}{h^d} k\left(\frac{\mathbf{u}}{h}\right)$ – тоже окно

2.3.3 Методы ближайшего соседа

Недостатки парзеновского оценивания

- *!! Ширина парзеновского окна h является структурным параметром!!*
- Очевидно, что при фиксированной ширине в некоторые участки пространства будет попадать избыточное количество объектов, и там ширину можно было бы уменьшить и получить более точную оценку плотности
- В то же время, найдутся участки, в которых число объектов, попавших в окно, будет слишком мало, и оценка плотности будет неустойчивой
- Вывод: необходимо сделать ширину окна переменной, потребовав, чтобы в него попадало определенное количество объектов
- Методы этой группы называются методами ближайших соседей

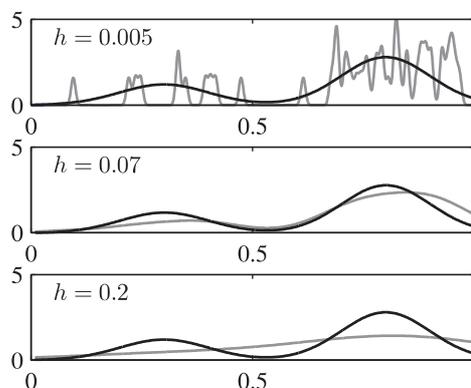


Рис. 2.4. Пример парзеновского оценивания

Идея метода

- Как было получено ранее, значение плотности может быть приближено величиной

$$\hat{p}(\mathbf{x}) = \frac{k}{nV},$$

где V – объем области \mathcal{D} , содержащей точку \mathbf{x}

- Потребуем, чтобы \mathcal{D} являлась сферой минимального радиуса с центром в точке \mathbf{x} , содержащей ровно k точек обучающей выборки
- Можно показать, что для сходимости $\hat{p}(\mathbf{x})$ к $p(\mathbf{x})$ при $n \rightarrow \infty$ необходимыми условиями являются требования

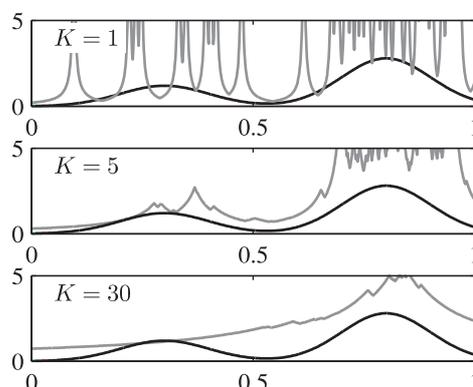
✓ Упр.

$$\lim_{n \rightarrow \infty} k = \infty \quad \lim_{n \rightarrow \infty} \frac{k}{n} = 0$$

- В случае нарушения первого требования оценка плотности будет почти всюду ноль, а при нарушении второго – почти всюду бесконечность, т.к. $V \rightarrow 0$ при неограниченном возрастании n

Особенности использования методов ближайших соседей

- *!! Количество ближайших соседей k является структурным параметром!!*
- В большинстве приложений выбор $k \sim \sqrt{n}$ является адекватным
- Серьезным недостатком метода является тот факт, что получающаяся оценка $\hat{p}(\mathbf{x})$ не является плотностью, т.к. интеграл от нее, вообще говоря, расходится
- Так в одномерном случае функция $\hat{p}(x)$ имеет хвост порядка $\frac{1}{x}$
- С ростом размерности пространства и объема выборки этот недостаток становится менее критическим

Рис. 2.5. Пример восстановления плотности с помощью метода ближайшего соседа с различным значением K

2.4 EM-алгоритм

2.4.1 Параметрическое восстановление плотностей

Параметрический подход

- При параметрическом восстановлении плотностей предполагается, что искомая плотность нам известна с точностью до параметра $p(\mathbf{x}|\theta)$
- Оценка плотности происходит путем оптимизации по θ некоторого функционала
- Обычно в качестве последнего используется правдоподобие или его логарифм

$$p(X|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta) \rightarrow \max_{\theta}$$

- В курсе математической статистики доказывается, что оценки максимального правдоподобия являются состоятельными, асимптотически несмещенными и эффективными (т.е. имеют наименьшую возможную дисперсию)

Пример использования

- Пусть имеется выборка из нормального распределения $\mathcal{N}(x|\mu, \sigma^2)$ с неизвестными мат. ожиданием и дисперсией
- Выписываем логарифм функции правдоподобия

$$L(X|\mu, \sigma) = - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} - n \log \sigma - \frac{n}{2} \log(2\pi) \rightarrow \max_{\mu, \sigma}$$

$$\frac{\partial L}{\partial \mu} = - \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2} = 0 \Rightarrow \mu_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial L}{\partial \sigma} = \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} - \frac{n}{\sigma} = 0 \Rightarrow \sigma_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

2.4.2 Задача разделения смеси распределений

Более сложная ситуация

- Часто возникают задачи, в которых генеральная совокупность является смесью нескольких элементарных распределений
- Дискретная смесь: $X \sim \sum_{j=1}^l w_j p(\mathbf{x}|\theta_j)$, $\sum_{j=1}^l w_j = 1$, $w_j \geq 0$
- Непрерывная смесь: $X \sim \int w(\xi) p(\mathbf{x}|\theta(\xi)) d\xi$, $\int w(\xi) d\xi = 1$, $w(\xi) \geq 0$
- Даже в случае, когда коэффициенты смеси известны, оптимизация правдоподобия крайне затруднительна
- Существует специальный алгоритм, позволяющий итеративно проводить оценивание коэффициентов смеси и параметров каждого из распределений

Особенности функционала качества

- Основная трудность при применении стандартного метода максимального правдоподобия заключается в необходимости оптимизации по θ величины

$$L(X|\theta) = \sum_{i=1}^n \log \left(\sum_{j=1}^l w_j p(\mathbf{x}_i|\theta_j) \right)$$

- Этот функционал имеет вид «логарифм суммы» и крайне сложен для оптимизации
- Необходимо перейти к функционалу вида «сумма логарифмов», который легко оптимизируется в явном виде

Суть EM-алгоритма

- Если бы мы знали, какой объект из какой компоненты смеси взят, т.е. функцию $j(i)$ (будем считать, что соответствующая информация содержится в ненаблюдаемой переменной z), то правдоподобие выглядело бы куда проще

$$L(X, Z|\theta) = \sum_{i=1}^n \log w_{j(i)} p(\mathbf{x}_i|\theta_{j(i)})$$

- Идея EM-алгоритма заключается в введении **латентных**, т.е. ненаблюдаемых переменных Z , позволяющих упростить вычисление правдоподобия
- Оптимизация проводится итерационно методом покоординатного спуска: на каждой итерации последовательно уточняются возможные значения Z (E-шаг), а потом пересчитываются значения θ (M-шаг)

Схема EM-алгоритма

- На входе: выборка X , зависящая от набора параметров θ, \mathbf{w}
- Инициализируем θ, \mathbf{w} некоторыми начальными приближениями
- E-шаг: Оцениваем распределение скрытой компоненты

$$p(Z|X, \theta, \mathbf{w}) = \frac{p(X, Z|\theta, \mathbf{w})}{\sum_Z p(X, Z|\theta, \mathbf{w})}$$

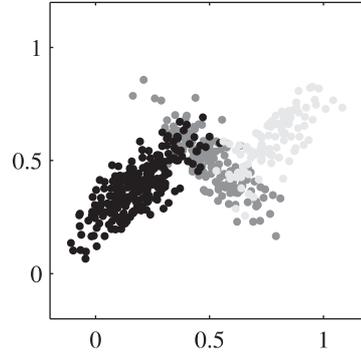


Рис. 2.6. Выборка из смеси трех нормальных распределений

- М-шаг: Оптимизируем

$$\mathbb{E}_Z \log p(X, Z | \boldsymbol{\theta}, \boldsymbol{w}) = \sum_Z p(Z | X, \boldsymbol{\theta}, \boldsymbol{w}) \log p(X, Z | \boldsymbol{\theta}, \boldsymbol{w}) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{w}}$$

Если бы мы **точно знали значение** $Z = Z_0$, то вместо мат. ожидания по всевозможным (с учетом наблюдаемых данных) Z , мы бы оптимизировали $\log p(X, Z_0 | \boldsymbol{\theta}, \boldsymbol{w})$

- Переход к Е-шагу, пока процесс не сойдется

2.4.3 Разделение гауссовской смеси

Смесь гауссовских распределений

- Имеется выборка $X \sim \sum_{j=1}^l w_j \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_j, \Sigma_j) \in \mathbb{R}^d$ (см. рис. 2.6)
- Требуется восстановить плотность генеральной совокупности

EM-алгоритм

- Выбираем начальное приближение $\boldsymbol{\mu}_j, w_j, \Sigma_j$
- Е-шаг: Вычисляем распределение скрытых переменных $z_i \in \{0, 1\}$, $\sum_j z_{ij} = 1$, которые определяют, к какой компоненте смеси принадлежит объект \boldsymbol{x}_i

$$\gamma(z_{ij}) = \frac{w_j \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_j, \Sigma_j)}{\sum_{k=1}^l w_k \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_k, \Sigma_k)}$$

- М-шаг: С учетом новых вероятностей на z_i , пересчитываем параметры смеси

$$\boldsymbol{\mu}_j^{new} = \frac{1}{N_j} \sum_{i=1}^n \gamma(z_{ij}) \boldsymbol{x}_i, \quad w_j^{new} = \frac{N_j}{n}, \quad N_j = \sum_{i=1}^n \gamma(z_{ij})$$

$$\Sigma_j^{new} = \frac{1}{N_j} \sum_{i=1}^n \gamma(z_{ij}) (\boldsymbol{x}_i - \boldsymbol{\mu}_j^{new})^T (\boldsymbol{x}_i - \boldsymbol{\mu}_j^{new})$$

- Переход к Е-шагу, пока не будет достигнута сходимость

Недостатки EM-алгоритма

- В зависимости от выбора начального приближения может сходиться к разным точкам
- EM-алгоритм находит локальный экстремум, в котором значение правдоподобия может оказаться намного ниже, чем в глобальном максимуме
- EM-алгоритм **не позволяет определить количество компонентов смеси l**
- *!! Величина l является структурным параметром!!*

Глава 3

Обобщенные линейные модели

Данная глава посвящена описанию простейших методов решения задачи восстановления регрессии и классификации. Изложение ведется в контексте т.н. обобщенных линейных моделей. Дается подробное описание метода построения линейной регрессии. Особое внимание уделено вопросам регуляризации задачи. Приводится описание метода наименьших квадратов с итеративным перевзвешиванием, используемого для обучения логистической регрессии

3.1 Ликбез: Псевдообращение матриц и нормальное псевдорешение

Псевдообращение матриц

- Предположим, нам необходимо решить СЛАУ вида $Ax = b$
- Если бы матрица A была квадратной и невырожденной (число уравнений равно числу неизвестных и все уравнения линейно независимы), то решение задавалось бы формулой $x = A^{-1}b$
- Предположим, что число уравнений больше числа неизвестных, т.е. матрица A прямоугольная. Домножим обе части уравнения на A^T слева

$$A^T Ax = A^T b$$

- В левой части теперь квадратная матрица и ее можно перенести в правую часть

$$x = (A^T A)^{-1} A^T b$$

- Операция $(A^T A)^{-1} A^T$ называется псевдообращением матрицы A , а x – псевдорешением

Нормальное псевдорешение

- Если матрица $A^T A$ вырождена, псевдорешений бесконечно много, причем найти их на компьютере нетривиально
- Для решения этой проблемы используется ридж-регуляризация матрицы $A^T A$

$$A^T A + \lambda I,$$

где I – единичная матрица, а λ – коэффициент регуляризации. Такая матрица невырождена для любых $\lambda > 0$

- Величина

$$x = (A^T A + \lambda I)^{-1} A^T b$$

называется нормальным псевдорешением. Оно всегда единственно и при небольших положительных λ определяет псевдорешение с наименьшей нормой

Графическая иллюстрация

- Псевдорешение соответствует точке, минимизирующей невязку, а нормальное псевдорешение отвечает псевдорешению с наименьшей нормой (см. рис. 3.1)
- Заметим, что псевдообратная матрица $(A^T A)^{-1} A^T$ совпадает с обратной матрицей A^{-1} в случае невырожденных квадратных матриц

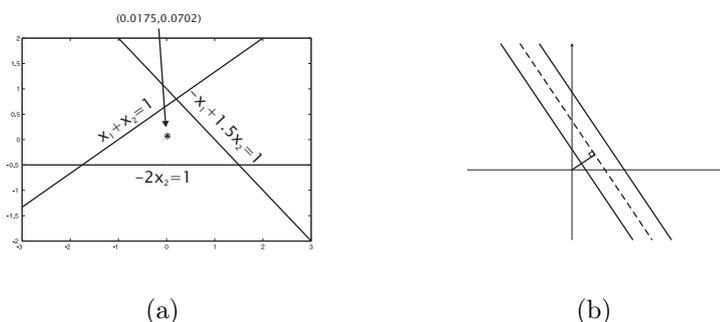


Рис. 3.1. На рисунке (a) показан пример единственного псевдорешения (которое одновременно является и нормальным псевдорешением), на рисунке (b) пунктирная линия обозначает множество всех псевдорешений, а нормальное псевдорешение является основанием перпендикуляра, опущенного из начала координат

3.2 Линейная регрессия

3.2.1 Классическая линейная регрессия

Задача восстановления регрессии

- Задача восстановления регрессии предполагает наличие связи между наблюдаемыми признаками \mathbf{x} и непрерывной переменной t
- В отличие от задачи интерполяции допускаются отклонения решающего правила от правильных ответов на объектах обучающей выборки
- Уравнение регрессии $y(\mathbf{x}, \mathbf{w})$ ищется в некотором параметрическом виде путем нахождения наилучшего значения вектора весов

$$\mathbf{w}_* = \arg \max_{\mathbf{w}} F(X, \mathbf{t}, \mathbf{w})$$

Линейная регрессия

- Наиболее простой и изученной является линейная регрессия
- Главная особенность: настраиваемые параметры входят в решающее правило **линейно**
- Заметим, что линейная регрессия не обязана быть линейной по признакам
- Общее уравнение регрессии имеет вид

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^m w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Особенность выбора базисных функций

- Общего метода выбора базисных функций $\phi_j(\mathbf{x})$ — не существует
- Обычно они подбираются из априорных соображений (например, если мы пытаемся восстановить какой-то периодический сигнал, разумно взять функции тригонометрического ряда) или путем использования некоторых «универсальных» базисных функций
- Наиболее распространенными базисными функциями являются
 - $\phi(\mathbf{x}) = x_k$

- $\phi(\mathbf{x}) = x_{k_1} x_{k_2} \dots x_{k_l}$
- $\phi(\mathbf{x}) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_0\|^p)$, $\gamma, p > 0$.

- Метод построения линейной регрессии (настройки весов \mathbf{w}) **не зависит** от выбора базисных функций

Формализация задачи

- Пусть $S(t, \hat{t})$ — функция потерь от ошибки в определении регрессионной переменной t
- Необходимо минимизировать потери от ошибок на генеральной совокупности

$$\mathbb{E}S(t, y(\mathbf{x}, \mathbf{w})) = \int \int S(t, y(\mathbf{x}, \mathbf{w})) p(\mathbf{x}, t) d\mathbf{x} dt \rightarrow \min_{\mathbf{w}}$$

- Дальнейшие рассуждения зависят от вида функции потерь
- Во многих случаях даже не нужно восстанавливать полностью условное распределение $p(t|\mathbf{x})$

Важная теорема

- Теорема. Пусть функция потерь имеет вид
 - $S(t, \hat{t}) = (t - \hat{t})^2$ — «Потери старушки»;
 - $S(t, \hat{t}) = |t - \hat{t}|$ — «Потери олигарха»;
 - $S(t, \hat{t}) = \delta^{-1}(t - \hat{t})$ — «Потери инвалида».

Тогда величиной, минимизирующей функцию $\mathbb{E}S(t, y(\mathbf{x}, \mathbf{w}))$, является следующая

- $y(\mathbf{x}) = \mathbb{E}p(t|\mathbf{x})$;
- $y(\mathbf{x}) = \text{med } p(t|\mathbf{x})$;
- $y(\mathbf{x}) = \text{mod } p(t|\mathbf{x}) = \arg \max_t p(t|\mathbf{x})$.

- В зависимости от выбранной системы предпочтений, мы будем пытаться оценивать тот или иной функционал от апостериорного распределения **вместо того, чтобы оценивать его самого**

3.2.2 Метод наименьших квадратов

Минимизация невязки

- Наиболее часто используемой функцией потерь является квадратичная $S(t, \hat{t}) = (t - \hat{t})^2$
- Значение регрессионной функции на обучающей выборке в матричном виде может быть записано как $\mathbf{y} = \Phi \mathbf{w}$, где $\Phi = (\phi_{ij}) = (\phi_j(\mathbf{x}_i)) \in \mathbb{R}^{n \times m}$
- Таким образом, приходим к следующей задаче

$$\|\mathbf{y} - \mathbf{t}\|^2 = \|\Phi \mathbf{w} - \mathbf{t}\|^2 \rightarrow \min_{\mathbf{w}}$$

Взяв производную по \mathbf{w} и приравняв ее к нулю, получаем

$$\frac{\partial \|\Phi \mathbf{w} - \mathbf{t}\|^2}{\partial \mathbf{w}} = \frac{\partial [\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t}]}{\partial \mathbf{w}} = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{t} = 0$$

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Регуляризация задачи

- Заметим, что формула для весов линейной регрессии представляет собой псевдорешение уравнения $\Phi \mathbf{w} = \mathbf{t}$
- Матрица $\Phi^T \Phi \in \mathbb{R}^{m \times m}$ вырождена при $m > n$
- Регуляризуя вырожденную матрицу, получаем

$$\mathbf{w} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t}$$

- Отсюда формула для прогноза объектов обучающей выборки по их правильным значениям

$$\hat{\mathbf{t}} = \mathbf{y} = \Phi (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t} = H \mathbf{t}$$

С историческим обозначением прогноза — навешиванием шляпки связано неформальное название матрицы H , по-английски звучащее как hat-matrix

Особенности квадратичной функции потерь

- Достоинства
 - Квадратичная функция потерь гладкая (непрерывная и дифференцируемая)
 - Решение может быть получено в явном виде
 - Существует простая вероятностная интерпретация прогноза и функции потерь
- Недостатки
 - Решение неустойчиво (не робастно) относительно даже малого количества выбросов. Это связано с быстрым возрастанием квадратичной функции потерь при больших отклонениях от нуля
 - Квадратичная функция неприменима к задачам классификации

3.2.3 Вероятностная постановка задачи**Нормальное распределение ошибок**

- Рассмотрим вероятностную постановку задачи восстановления регрессии. Регрессионная переменная t — случайная величина с плотностью распределения $p(t|\mathbf{x})$
- В большинстве случаев предполагается, что t распределена нормально относительно некоторого мат. ожидания $y(\mathbf{x})$, определяемого точкой \mathbf{x}

$$t = y(\mathbf{x}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(\varepsilon|0, \sigma^2)$$

- Необходимо найти функцию $y(\mathbf{x})$, которую мы можем отождествить с уравнением регрессии
- Предположение о нормальном распределении отклонений можно обосновать ссылкой на центральную предельную теорему

Метод максимального правдоподобия для регрессии

- Используем ММП для поиска $y(\mathbf{x})$
- Правдоподобие задается следующей формулой

$$p(\mathbf{t}|\mathbf{y}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t_i - y_i)^2}{2\sigma^2}\right) \rightarrow \max$$

- Взяв логарифм и отбросив члены, не влияющие на положение максимума, получим

$$\sum_{i=1}^n (t_i - y_i)^2 = \sum_{i=1}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 \rightarrow \min_{\mathbf{w}}$$

- Таким образом, применение метода максимального правдоподобия в предположении о нормальности отклонений эквивалентно методу наименьших квадратов

Вероятностный смысл регуляризации

- Теперь будем максимизировать не правдоподобие, а апостериорную вероятность
- По формуле условной вероятности

$$p(\mathbf{w}|\mathbf{t}, X) = \frac{p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{t}, X)} \rightarrow \max_{\mathbf{w}},$$

знаменатель не зависит от \mathbf{w} , поэтому им можно пренебречь

- Пусть $p(\mathbf{w}) \sim \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \left(\frac{\sigma^2}{\lambda}\right) I\right)$. Тогда

$$p(\mathbf{w}|\mathbf{t}, X) \propto \frac{\lambda^{m/2}}{(\sqrt{2\pi}\sigma)^{m+n}} \exp\left(-\frac{1}{2}\left(\sigma^{-2}\|\Phi\mathbf{w} - \mathbf{t}\|^2 + \frac{\lambda}{\sigma^2}\|\mathbf{w}\|^2\right)\right)$$

- Логарифмируя и приравнявая производную по \mathbf{w} к нулю, получаем

$$\mathbf{w} = (\Phi^T \Phi + \lambda I)^{-1} \Phi \mathbf{t}$$

- Регуляризация эквивалентна введению априорного распределения, поощряющего небольшие веса

3.3 Применение регрессионных методов для задачи классификации

3.3.1 Логистическая регрессия

Особенности задачи классификации

- Рассмотрим задачу классификации на два класса $t \in \{-1, +1\}$
- Ее можно свести к задаче регрессии, например, следующим образом

$$\hat{t}(\mathbf{x}) = \text{sign}(y(\mathbf{x})) = \text{sign} \sum_{j=1}^m w_j \phi_j(\mathbf{x})$$

- Возникает вопрос: что использовать в качестве значений регрессионной переменной на этапе обучения?
- Наиболее распространенный подход заключается в использовании значения $+\infty$ для $t = +1$ и $-\infty$ для $t = -1$
- Геометрический смысл: чем дальше от нуля значение $y(\mathbf{x})$, тем увереннее мы в классификации объекта \mathbf{x}

Правдоподобие правильной классификации

- Метод наименьших квадратов, очевидно, неприменим при таком подходе
- Воспользуемся вероятностной постановкой для выписывания функционала качества
- Определим правдоподобие классификации следующим образом

$$p(t|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-ty(\mathbf{x}))}$$

- Это логистическая функция (см. рис. 3.2). Легко показать, что $\sum_t p(t|\mathbf{x}, \mathbf{w}) = 1$ и $p(t|\mathbf{x}, \mathbf{w}) > 0$, а, значит, она является функцией правдоподобия

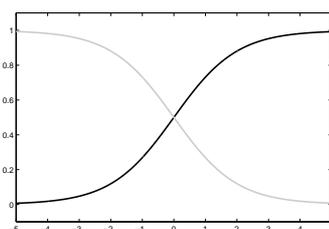


Рис. 3.2. Логистическая функция переводит выход линейной функции в вероятностные значения. Черная кривая показывает правдоподобие для случая $t_i = 1$, а серая кривая — для случая $t_i = -1$

Функционал качества в логистической регрессии

- Правдоподобие правильной классификации всей выборки имеет вид

$$p(\mathbf{t}|X, \mathbf{w}) = \prod_{i=1}^n p(t_i|\mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{1 + \exp\left(-t_i \sum_{j=1}^m w_j \phi_j(\mathbf{x}_i)\right)}$$

3.3.2 Метод IRLS

Особенности функции правдоподобия классификации

- Приравнивание градиента логарифма правдоподобия к нулю приводит к трансцендентным уравнениям, которые неразрешимы аналитически
- Легко показать, что гессиан логарифма правдоподобия неположительно определен

$$\frac{\partial^2 \log p(\mathbf{t}|\mathbf{x}, \mathbf{w})}{\partial \mathbf{w}^2} \leq 0$$

- Это означает, что логарифм функции правдоподобия является вогнутым.
- Логарифм правдоподобия обучающей выборки $L(\mathbf{w}) = \log p(\mathbf{t}|X, \mathbf{w})$, являющийся суммой вогнутых функций, также вогнут, а, значит, имеет **единственный максимум**

Метод оптимизации Ньютона

Основная идея метода Ньютона — это приближение в заданной точке оптимизируемой функции параболой и выбор минимума этой параболы в качестве следующей точки итерационного процесса:

$$f(\mathbf{x}) \rightarrow \min_{\mathbf{w}}$$

$$f(\mathbf{x}) \simeq g(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0))^T(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T(\nabla\nabla f(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0)$$

$$\nabla g(\mathbf{x}_*) = \nabla f(\mathbf{x}_0) + (\nabla\nabla f(\mathbf{x}_0))(\mathbf{x}_* - \mathbf{x}_0) = 0$$

$$\mathbf{x}_* = \mathbf{x}_0 - (\nabla\nabla f(\mathbf{x}_0))^{-1}(\nabla f(\mathbf{x}_0))$$

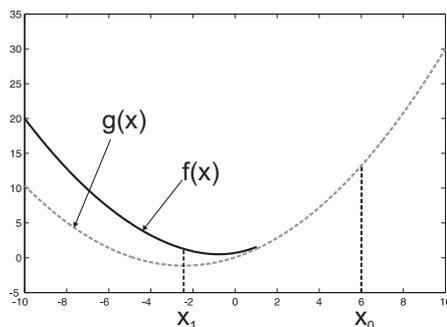


Рис. 3.3. Пример оптимизации с помощью метода Ньютона. Функция $f(x) = \log(1 + \exp(x)) + \frac{x^2}{5}$. В точке $x_0 = 6$ проведено приближение функции $f(x)$ параболой $g(x)$. Точка минимума этой параболы $x_1 = -2.4418$ является следующей точкой итерационного процесса

Итеративная минимизация логарифма правдоподобия

- Так как прямая минимизация правдоподобия невозможна, воспользуемся итерационным методом Ньютона
- Обоснованием корректности использования метода Ньютона является унимодальность оптимизируемой функции $L(\mathbf{w})$ и ее гладкость во всем пространстве весов
- Формула пересчета в методе Ньютона

$$\mathbf{w}^{new} = \mathbf{w}^{old} - H^{-1}\nabla L(\mathbf{w}),$$

где $H = \nabla\nabla L(\mathbf{w})$ — гессиан логарифма правдоподобия обучающей выборки

Формулы пересчета

Обозначим $s_i = \frac{1}{1 + \exp(-t_i y_i)}$, тогда:

$$\nabla L(\mathbf{w}) = \Phi^T \text{diag}(\mathbf{t})\mathbf{s}, \quad \nabla\nabla L(\mathbf{w}) = \Phi^T R \Phi$$

$$R = \begin{pmatrix} s_1(1-s_1) & 0 & \dots & 0 \\ 0 & s_2(1-s_2) & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & s_n(1-s_n) \end{pmatrix}$$

$$\begin{aligned} \mathbf{w}^{new} = \mathbf{w}^{old} - (\Phi^T R \Phi)^{-1} \Phi^T \text{diag}(\mathbf{t}) \mathbf{s} = \\ (\Phi^T R \Phi)^{-1} (\Phi^T R \Phi \mathbf{w}^{old} - \Phi^T R R^{-1} \text{diag}(\mathbf{t}) \mathbf{s}) = (\Phi^T R \Phi)^{-1} \Phi^T R \mathbf{z}, \end{aligned}$$

где $\mathbf{z} = \Phi \mathbf{w}^{old} - R^{-1} \text{diag}(\mathbf{t}) \mathbf{s}$

Название метода (метод наименьших квадратов с итеративно пересчитываемыми весами) связано с тем, что последняя формула является формулой для взвешенного МНК (веса задаются диагональной матрицей R), причем на каждой итерации веса корректируются

Заключительные замечания

- На практике матрица $\Phi^T R \Phi$ часто бывает вырождена (всегда при $m > n$), поэтому обычно прибегают к регуляризации матрицы $(\Phi^T R \Phi + \lambda I)$
- *!! Параметр регуляризации λ является структурным параметром!!*
- *!! Базисные функции $\phi_j(\mathbf{x})$, а значит и матрица Φ являются структурными параметрами!!*
- С поиском методов автоматического выбора базисных функций связана одна из наиболее интригующих проблем современного машинного обучения

Глава 4

Метод опорных векторов и беспризнаковое распознавание образов

В главе подробно рассматривается метод опорных векторов для классификации и восстановления регрессии. Особое внимание уделено формулировке двойственной задачи и использованию правила множителей Лагранжа. Описывается т.н. ядровой переход, представляющий нелинейное обобщение метода опорных векторов, показана связь между этим методом и методом максимального правдоподобия с регуляризацией, а также со статистической теорией обучения Валника-Червоненкиса. В конце главы приведены обобщения метода опорных векторов на задачи, в которых подсчет признаков невозможен или нецелесообразен, но в которых естественным образом можно ввести функцию близости между объектам.

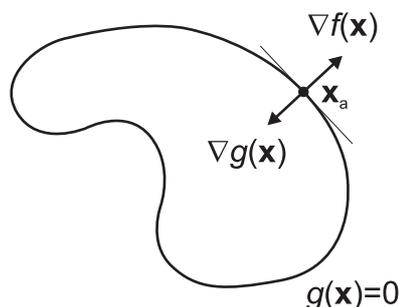


Рис. 4.1. Иллюстрация к задаче оптимизации с ограничениям в виде равенства. В оптимальной точке градиенты ∇f и ∇g должны быть параллельны друг другу

4.1 Ликбез: Условная оптимизация

Задача условной оптимизации

Пусть $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ — гладкая функция. Предположим, что нам необходимо найти ее экстремум:

$$f(\mathbf{x}) \rightarrow \operatorname{extr}_{\mathbf{x}}$$

Для того, чтобы найти экстремум (решить задачу безусловной оптимизации), достаточно проверить условие стационарности:

$$\nabla f(\mathbf{x}) = 0$$

Предположим, что нам необходимо найти экстремум функции при ограничениях:

$$\begin{aligned} f(\mathbf{x}) &\rightarrow \operatorname{extr}_{\mathbf{x}} \\ g(\mathbf{x}) &= 0 \end{aligned}$$

Поверхность ограничения (см. рис. 4.1)

Заметим, что $\nabla g(\mathbf{x})$ ортогонален поверхности ограничения $g(\mathbf{x}) = 0$. Пусть \mathbf{x} и $\mathbf{x} + \boldsymbol{\varepsilon}$ — две близкие точки поверхности. Тогда

$$g(\mathbf{x} + \boldsymbol{\varepsilon}) \simeq g(\mathbf{x}) + \boldsymbol{\varepsilon}^T \nabla g(\mathbf{x})$$

Т.к. $g(\mathbf{x} + \boldsymbol{\varepsilon}) = g(\mathbf{x})$, то $\boldsymbol{\varepsilon}^T \nabla g(\mathbf{x}) \simeq 0$. При стремлении $\|\boldsymbol{\varepsilon}\| \rightarrow 0$ получаем $\boldsymbol{\varepsilon}^T \nabla g(\mathbf{x}) = 0$. Т.к. $\boldsymbol{\varepsilon}$ параллелен поверхности $g(\mathbf{x}) = 0$, то $\nabla g(\mathbf{x})$ является нормалью к этой поверхности.

Функция Лагранжа

Необходимым условием оптимальности является ортогональность $\nabla f(\mathbf{x})$ поверхности ограничения (в противном случае вектор проекции градиента $\nabla f(\mathbf{x})$ на поверхность ограничения имеет ненулевую длину, и можно найти большее значение функции, двигаясь вдоль вектора проекции), т.е.:

$$\nabla f + \lambda \nabla g = 0$$

Здесь $\lambda \neq 0$ — коэффициент Лагранжа. Он может быть любого знака.

Функция Лагранжа

$$L(\mathbf{x}, \lambda) \triangleq f(\mathbf{x}) + \lambda g(\mathbf{x})$$

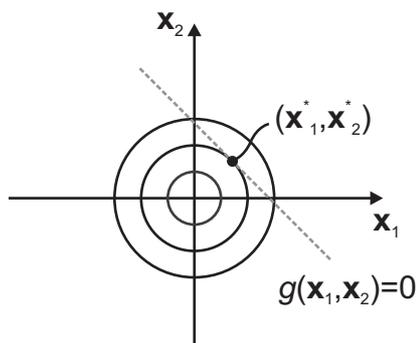


Рис. 4.2. Пример решения задачи условной оптимизации. Черные линии обозначают линии уровня оптимизируемой функции. Серая пунктирная прямая представляет собой область поиска решения. Точка $(x_1^*, x_2^*) = (1/2, 1/2)$ является решением задачи.

Тогда

$$\begin{aligned} \nabla_{\mathbf{x}} L &= 0 & \Rightarrow \text{условие ортогональности } \nabla f + \lambda \nabla g = 0 \\ \frac{\partial}{\partial \lambda} L &= 0 & \Rightarrow g(\mathbf{x}) = 0 \end{aligned}$$

Функция Лагранжа. Пример (см. рис. 4.2)

$$\begin{aligned} f(x_1, x_2) &= 1 - x_1^2 - x_2^2 \rightarrow \max_{x_1, x_2} \\ g(x_1, x_2) &= x_1 + x_2 - 1 = 0 \end{aligned}$$

Функция Лагранжа:

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

Условия стационарности:

$$\begin{aligned} -2x_1 + \lambda &= 0 \\ -2x_2 + \lambda &= 0 \\ x_1 + x_2 - 1 &= 0 \end{aligned}$$

Решение: $(x_1^*, x_2^*) = (\frac{1}{2}, \frac{1}{2})$, $\lambda = 1$.

Ограничение в виде неравенства (см. рис. 4.3)

Задача условной оптимизации

$$\begin{aligned} f(\mathbf{x}) &\rightarrow \max_{\mathbf{x}} \\ g(\mathbf{x}) &\geq 0 \end{aligned}$$

Решение	Ограничение	Условие стационарности
Внутри области $g(\mathbf{x}) > 0$	неактивно	$\nabla f(\mathbf{x}) = 0, \nabla_{\mathbf{x}} L = 0, \lambda = 0$
На границе $g(\mathbf{x}) = 0$	активно	$\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x}), \nabla_{\mathbf{x}, \lambda} L = 0, \lambda > 0$

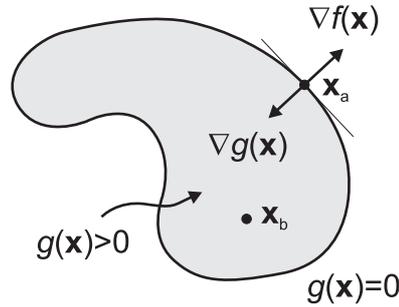


Рис. 4.3. Иллюстрация к задаче оптимизации с ограничениями в виде неравенства. В случае, если оптимальная точка лежит вне области $g(\mathbf{x}) \geq 0$, то в точке оптимума градиенты ∇f и ∇g должны быть коллинеарны и направлены в разные стороны

Условие дополняющей нежесткости:

$$\lambda g(\mathbf{x}) = 0$$

Теорема Каруша-Куна-Таккера

Пусть $f_i : X \rightarrow \mathbb{R}$, $i = 0, 1, \dots, m$ — выпуклые функции, отображающие нормированное пространство X в прямую, $A \in X$ — выпуклое множество. Рассмотрим следующую задачу оптимизации:

$$f_0(\mathbf{x}) \rightarrow \min; \quad f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \quad \mathbf{x} \in A \quad (P)$$

Теорема 1.

1. Если $\hat{\mathbf{x}} \in \text{absmin}(P)$ — решение задачи, то найдется ненулевой вектор множителей Лагранжа $\lambda \in \mathbb{R}^{m+1}$ такой, что для функции Лагранжа $L(\mathbf{x}) = \sum_{i=0}^m \lambda_i f_i(\mathbf{x})$ выполняются условия:

- стационарности $\min_{\mathbf{x} \in A} L(\mathbf{x}) = L(\hat{\mathbf{x}})$
- дополняющей нежесткости $\lambda_i f_i(\hat{\mathbf{x}}) = 0$, $i = 1, \dots, m$
- неотрицательности $\lambda_i \geq 0$

2. Если для допустимой точки $\hat{\mathbf{x}}$ выполняются условия а)-с) и $\lambda_0 \neq 0$, то $\hat{\mathbf{x}} \in \text{absmin}(P)$

3. Если для допустимой точки $\hat{\mathbf{x}}$ выполняются условия а)-с) и $\exists \tilde{\mathbf{x}} \in A : f_i(\tilde{\mathbf{x}}) < 0, i = 0, \dots, m$ (условие Слейтера), то $\hat{\mathbf{x}} \in \text{absmin}(P)$

4.2 Метод опорных векторов для задачи классификации

Задача классификации

- Рассматривается задача классификации на два класса. Имеется обучающая выборка $(X, \mathbf{t}) = \{\mathbf{x}_i, t_i\}_{i=1}^n$, где $\mathbf{x} \in \mathbb{R}^d$, а метка класса $t \in \mathcal{T} = \{-1, 1\}$.
- Необходимо с использованием обучающей выборки построить отображение $A : \mathbb{R}^d \rightarrow \mathcal{T}$, которое для каждого нового входного объекта \mathbf{x}_* выдает его метку класса t_* .

4.2.1 Метод потенциальных функций

Метод потенциальных функций [Айзерман и др., 1964]

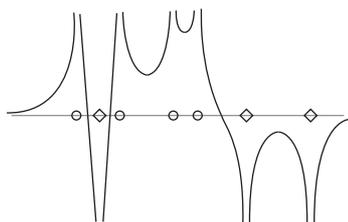


Рис. 4.4. Иллюстрация к методу потенциальных функций.

В каждом объекте \mathbf{x}_i помещен электрический заряд $t_i q_i$. В качестве разделяющей функции используется потенциал создаваемого поля:

$$f(\mathbf{x}) = \sum_{i=1}^n t_i q_i K(\mathbf{x}, \mathbf{x}_i)$$

Функция $K(\mathbf{x}, \mathbf{y})$ называется потенциальной функцией и имеет смысл потенциала в точке \mathbf{y} , создаваемого зарядом в точке \mathbf{x} . Предполагается, что

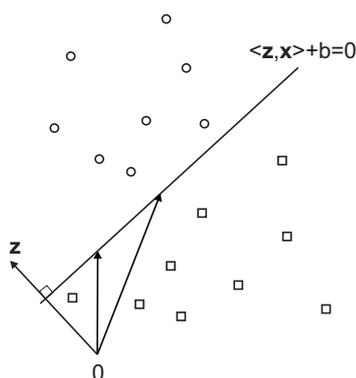
$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &\rightarrow 0 \text{ при } \|\mathbf{x} - \mathbf{y}\| \rightarrow +\infty \\ K(\mathbf{x}, \mathbf{y}) &\rightarrow \max \text{ при } \|\mathbf{x} - \mathbf{y}\| \rightarrow 0 \end{aligned}$$

Алгоритм обучения:

$$f^{new}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) + K(\mathbf{x}, \mathbf{x}_k), & \text{если } t_k = 1 \text{ и } f(\mathbf{x}_k) \leq 0 \\ f(\mathbf{x}) - K(\mathbf{x}, \mathbf{x}_k), & \text{если } t_k = -1 \text{ и } f(\mathbf{x}_k) \geq 0 \\ f(\mathbf{x}), & \text{иначе} \end{cases}$$

4.2.2 Случай линейно разделимых данных

Разделяющая гиперплоскость

Рис. 4.5. Разделяющая гиперплоскость для объектов двух классов. Направление гиперплоскости задается вектором нормали \mathbf{z}

- Гиперплоскость задается направляющим вектором гиперплоскости \mathbf{z} и величиной сдвига b (см. рис. 4.5):

$$\{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{z}, \mathbf{x} \rangle + b = 0\}, \quad \mathbf{z} \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

- Если вектор \mathbf{z} имеет единичную длину, то величина $\langle \mathbf{z}, \mathbf{x} \rangle$ определяет длину проекции вектора \mathbf{x} на направляющий вектор \mathbf{z} . В случае произвольной длины скалярное произведение нормируется на $\|\mathbf{z}\|$.

Каноническая гиперплоскость

- Если величину направляющего вектора \mathbf{z} и величину сдвига b умножить на одно и то же число, то соответствующая им гиперплоскость не изменится.
- Пара (\mathbf{z}, b) задает **каноническую гиперплоскость** для набора объектов $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, если

$$\min_{i=1, \dots, n} |\langle \mathbf{z}, \mathbf{x}_i \rangle + b| = 1 \quad (*)$$

- Условие (*) означает, что ближайший вектор к канонической гиперплоскости находится от нее на расстоянии $1/\|\mathbf{z}\|$:

$$\begin{aligned} \mathbf{x}_i : |\langle \mathbf{z}, \mathbf{x}_i \rangle + b| = 1, \quad \mathbf{x} : \langle \mathbf{z}, \mathbf{x} \rangle + b = 0 \\ |\langle \mathbf{z}, \mathbf{x}_i - \mathbf{x} \rangle| = 1 \Rightarrow \left| \left\langle \frac{\mathbf{z}}{\|\mathbf{z}\|}, \mathbf{x}_i - \mathbf{x} \right\rangle \right| = \frac{1}{\|\mathbf{z}\|} \end{aligned}$$

- Каноническая гиперплоскость определена однозначно с точностью до знака \mathbf{z} и b

Классификатор

- Будем искать классификатор в виде разделяющей гиперплоскости (см. рис. 4.6):

$$\hat{t}(\mathbf{x}) = \text{sign}(y(\mathbf{x})) = \text{sign}(\langle \mathbf{z}, \mathbf{x} \rangle + b)$$

- Предположим, что исходные данные являются линейно разделимыми, т.е.

$$\exists(\mathbf{z}, b) : \hat{t}(\mathbf{x}_i) = t_i \quad \forall i = 1, \dots, n$$

Оптимальная разделяющая гиперплоскость

- Зазором гиперплоскости данной точки (\mathbf{x}, t) называется величина:

$$\rho_{(\mathbf{z}, b)}(\mathbf{x}, t) = \frac{t(\langle \mathbf{z}, \mathbf{x} \rangle + b)}{\|\mathbf{z}\|}$$

Зазором гиперплоскости называется величина:

$$\rho_{(\mathbf{z}, b)} = \min_{i=1, \dots, n} \rho_{(\mathbf{z}, b)}(\mathbf{x}_i, t_i)$$

- Для корректно распознаваемого объекта его величина зазора соответствует расстоянию от этого объекта до гиперплоскости. Отрицательная величина зазора соответствует ошибочной классификации объекта.
- **Оптимальная разделяющая гиперплоскость** — гиперплоскость с максимальной величиной зазора

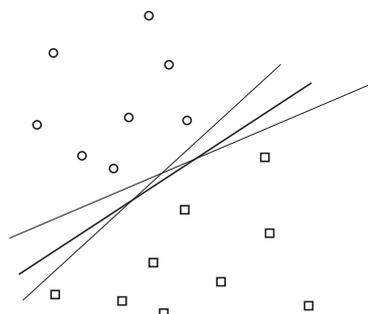


Рис. 4.6. Для линейно разделяемой выборки существует бесконечно большое число гиперплоскостей, корректно разделяющих данные. Жирной линией показана оптимальная разделяющая гиперплоскость.

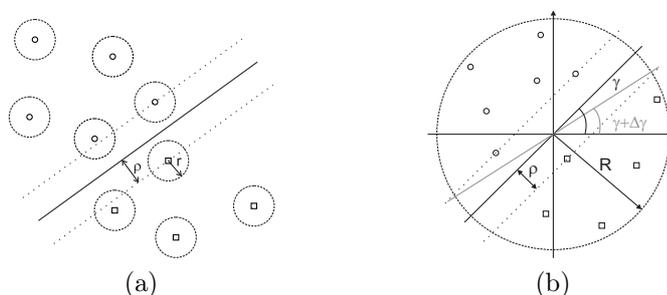


Рис. 4.7. Иллюстрация к оптимальной разделяющей гиперплоскости.

Оптимальная разделяющая гиперплоскость. Примеры

- Предположим, что все объекты тестовой выборки получены путем небольших смещений относительно обучающих объектов, т.е. для объекта обучения (\mathbf{x}, t) тестовый объект может быть представлен как $(\mathbf{x} + \Delta\mathbf{x}, t)$, причем $\|\Delta\mathbf{x}\| \leq r$.
- Гиперплоскость с величиной зазора $\rho > r$ корректно классифицирует выборку (см. рис. 4.7, а).
- Если объекты располагаются на достаточном расстоянии от гиперплоскости, то небольшое изменение параметров (\mathbf{z}, b) не меняет корректного разделения данных (см. рис. 4.7, б).

Оптимальная разделяющая гиперплоскость. Стат. теория Вапника-Червоненкиса

- Пусть имеется некоторое объективное распределение $p(\mathbf{x}, t)$. Обучающая и тестовая совокупности являются н.о.р. выборками из этого распределения.
- Пусть имеется семейство классификаторов $\{f(\mathbf{x}, \mathbf{w}) : \mathbb{R}^d \rightarrow \mathcal{T} | \mathbf{w} \in \Omega\}$ и функция потерь $L : \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}_+$.
- Средним риском называется мат.ожидание функции потерь:

$$R(\mathbf{w}) = \mathbb{E}_{p(\mathbf{x}, t)} L(\cdot) = \int L(t, f(\mathbf{x}, \mathbf{w})) p(\mathbf{x}, t) d\mathbf{x} dt$$

- Эмпирическим риском называется следующая величина:

$$R_{emp}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n L(t_i, f(\mathbf{x}_i, \mathbf{w}))$$

Теорема 2 (Vapnik, 1995). С вероятностью $0 < \eta \leq 1$ выполняется следующее неравенство:

$$R(\mathbf{w}) \leq R_{emp}(\mathbf{w}) + \sqrt{\frac{h(\ln(2n/h) + 1) - \ln(\eta/4)}{n}} \quad (*)$$

Здесь h — положительная величина, называемая ВЧ-размерностью.

Теорема 3 (Vapnik, 1995). Допустим, что вектора $\mathbf{x} \in \mathbb{R}^d$ принадлежат сфере радиуса R . Тогда для семейства разделяющих гиперплоскостей с величиной зазора ρ верно следующее:

$$h \leq \min \left(\left\lceil \frac{R^2}{\rho^2} \right\rceil, d \right) + 1 \quad (**)$$

Для гиперплоскости, корректно разделяющей данные, эмпирический риск равен нулю. Корень в оценке на средний риск (*) является монотонно возрастающей функцией по h . Поэтому минимизация оценки эквивалентна минимизации ВЧ-размерности, что согласно формуле (**) эквивалентно максимизации зазора ρ .

Постановка задачи оптимизации

- Предположим, что в выборке присутствуют объекты двух классов, т.е. $\exists i, j : t_i = 1, t_j = -1$.
- Для построения канонической гиперплоскости с максимальной величиной зазора, корректно разделяющей данные, необходимо решить следующую задачу оптимизации:

$$\begin{aligned} \frac{1}{2} \|z\|^2 &\rightarrow \min_{z, b} \\ t_i(\langle z, \mathbf{x}_i \rangle + b) &\geq 1, \quad \forall i = 1, \dots, n \end{aligned}$$

Функция Лагранжа

Выпишем функцию Лагранжа

$$L(z, b, \mathbf{w}) = \frac{1}{2} \|z\|^2 - \sum_{i=1}^n w_i (t_i(\langle z, \mathbf{x}_i \rangle + b) - 1) \rightarrow \min_{z, b} \max_{\mathbf{w}}$$

Коэффициенты Лагранжа $w_i \geq 0, \forall i = 1, \dots, n$.

$$\frac{\partial}{\partial z} L(z, b, \mathbf{w}) = 0 \Rightarrow z_* = \sum_{i=1}^n w_i t_i \mathbf{x}_i$$

$$\frac{\partial}{\partial b} L(z, b, \mathbf{w}) = 0 \Rightarrow \sum_{i=1}^n w_i t_i = 0$$

$$\begin{aligned} L(z_*, b_*, \mathbf{w}) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j t_i t_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^n \sum_{j=1}^n w_i w_j t_i t_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \\ &+ \sum_{i=1}^n w_i = \sum_{i=1}^n w_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j t_i t_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \max_{\mathbf{w}} \end{aligned}$$

Двойственная задача оптимизации

$$\sum_{i=1}^n w_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j t_i t_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \max_w$$

$$\sum_{i=1}^n t_i w_i = 0$$

$$w_i \geq 0, \forall i = 1, \dots, n$$

Решение

$$\hat{t}(\mathbf{x}) = \text{sign}(\langle \mathbf{z}_*, \mathbf{x} \rangle + b_*) = \text{sign}\left(\sum_{i=1}^n w_i^* t_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b_*\right)$$

Опорные вектора

Условие дополняющей нежесткости:

$$w_i^* (t_i (\langle \mathbf{z}_*, \mathbf{x}_i \rangle + b_*) - 1) = 0$$

Из этого условия следует, что объекты обучающей выборки, для которых $w_i^* > 0$, лежат точно на границе гиперплоскости. Они называются **опорными векторами** (см. рис. 4.8). Значение b_* может быть получено

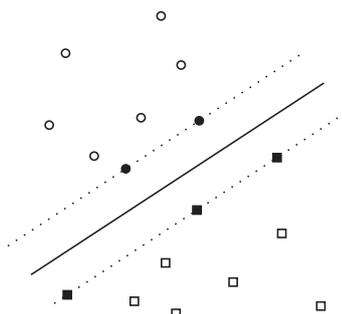


Рис. 4.8. Оптимальная разделяющая гиперплоскость. Объекты, определяющие положение гиперплоскости, называются опорными (обозначены черным цветом). Они располагаются ближе всего к гиперплоскости.

из условия дополняющей нежесткости для любого опорного вектора. При этом с вычислительной точки зрения более устойчивой процедурой является усреднение по всем таким объектам.

Разреженность решения

Решающее правило

$$\hat{t}(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^n w_i t_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b\right)$$

В решающее правило входят только те объекты обучения, для которых $w_i > 0$ (опорные векторы). Такое правило называется **разреженным (sparse model)**. Разреженные модели обладают высокой скоростью распознавания в больших объемах данных, а также «проливают свет» на структуру обучающей совокупности, выделяя наиболее релевантные с точки зрения классификации объекты.

4.2.3 Случай линейно неразделимых данных

Ослабляющие коэффициенты

На практике данные не являются, как правило, линейно разделимыми. Кроме того, даже линейно разделимые выборки могут содержать помехи, ошибочные метки классов и проч. Практический метод распознавания должен учитывать подобные ситуации.

Предположим, что $\forall(\mathbf{z}, b) \exists \mathbf{x}_i : \rho_{(\mathbf{z}, b)}(\mathbf{x}_i, t_i) < 0$. Позволим некоторым из ограничений не выполняться путем введения ослабляющих коэффициентов:

$$t_i(\langle \mathbf{z}, \mathbf{x}_i \rangle + b) \geq 1 \quad \forall i = 1, \dots, n \quad \longrightarrow \quad \begin{cases} t_i(\langle \mathbf{z}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \\ \xi_i \geq 0, \quad \forall i = 1, \dots, n \end{cases}$$

При этом потребуем, чтобы количество нарушений (количество ошибок на обучении) было бы как можно меньшим:

$$\frac{1}{2} \|\mathbf{z}\|^2 + C \sum_{i=1}^n \xi_i \rightarrow \min_{\mathbf{z}, b, \xi}$$

Постановка задачи оптимизации

$$\begin{aligned} \frac{1}{2} \|\mathbf{z}\|^2 + C \sum_{i=1}^n \xi_i &\rightarrow \min_{\mathbf{z}, b} \\ t_i(\langle \mathbf{z}, \mathbf{x}_i \rangle + b) &\geq 1 - \xi_i \quad \forall i = 1, \dots, n \\ \xi_i &\geq 0 \end{aligned}$$

Здесь $C \geq 0$ — некоторый действительный параметр, играющий роль параметра регуляризации

Функция Лагранжа

Выпишем функцию Лагранжа

$$L(\mathbf{z}, b, \xi, \mathbf{w}, \mathbf{v}) = \frac{1}{2} \|\mathbf{z}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n w_i [t_i(\langle \mathbf{z}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^n v_i \xi_i \rightarrow \min_{\mathbf{z}, b, \xi} \max_{\mathbf{w}, \mathbf{v}}$$

Коэффициенты Лагранжа $w_i \geq 0$, $v_i \geq 0$.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{z}} L(\mathbf{z}, b, \xi, \mathbf{w}, \mathbf{v}) &= 0 & \Rightarrow \mathbf{z}_* &= \sum_{i=1}^n w_i t_i \mathbf{x}_i \\ \frac{\partial}{\partial b} L(\mathbf{z}, b, \xi, \mathbf{w}, \mathbf{v}) &= 0 & \Rightarrow \sum_{i=1}^n w_i t_i &= 0 \\ \frac{\partial}{\partial \xi_i} L(\mathbf{z}, b, \xi, \mathbf{w}, \mathbf{v}) &= 0 & \Rightarrow w_i + v_i &= C \end{aligned}$$

Двойственная задача оптимизации

$$\begin{aligned} \sum_{i=1}^n w_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j t_i t_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle &\rightarrow \max_w \\ \sum_{i=1}^n w_i t_i &= 0 \\ 0 \leq w_i \leq C \end{aligned}$$

Решающее правило остается без изменений:

$$\hat{t}(\mathbf{x}) = \text{sign}(\langle \mathbf{z}_*, \mathbf{x} \rangle + b_*) = \text{sign}\left(\sum_{i=1}^n w_i^* t_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b^*\right)$$

4.2.4 Ядровой переход

Ядровой переход

- На практике часто встречается ситуация, когда данные порождаются нелинейной разделяющей поверхностью.
- Для обобщения метода на нелинейный случай заметим, что объекты обучающей выборки входят в двойственную задачу оптимизации только в виде попарных скалярных произведений $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$.
- Предположим, что исходное признаковое пространство было подвергнуто некоторому нелинейному преобразованию (см. рис. 4.9):

$$\Phi: \mathbb{R}^d \rightarrow \mathcal{H}$$

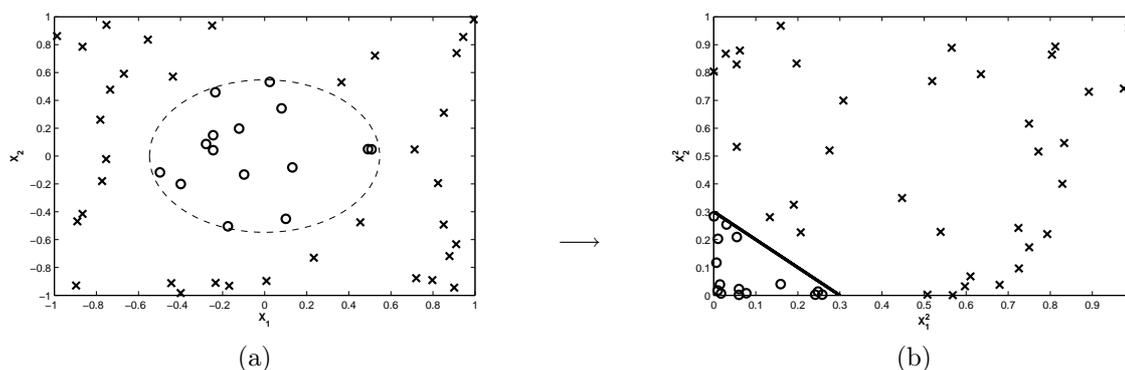


Рис. 4.9. Пример преобразования признакового пространства. На рис. (а) данные разделяются нелинейной поверхностью — эллипсом. Переход из пространства (x_1, x_2) к новому пространству (x_1^2, x_2^2) делает данные линейно разделяемыми.

Ядровой переход

- Для того, чтобы построить гиперплоскость с максимальным зазором в новом пространстве \mathcal{H} необходимо знать лишь $\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle_{\mathcal{H}}$.
- Допустим, что существует некоторая «ядровая функция» $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, такая что

$$K(\mathbf{x}, \mathbf{y}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle_{\mathcal{H}}$$

- Для построения гиперплоскости с максимальным зазором в пространстве \mathcal{H} нет необходимости задавать преобразование Φ в явном виде, достаточно лишь знать K !
- Задача оптимизации зависит только от попарных скалярных произведений, а решающее правило может быть представлено как

$$\hat{t}(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n w_i t_i \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}) \rangle_{\mathcal{H}} + b \right) = \text{sign} \left(\sum_{i=1}^n w_i t_i K(\mathbf{x}_i, \mathbf{x}) + b \right)$$

Требования к ядровой функции

Очевидно, что не для любой функции двух переменных K найдутся такие (\mathcal{H}, Φ) , для которых K будет определять скалярное произведение. Необходимыми и достаточными требованиями являются:

- Симметричность

$$K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$$

- Неотрицательная определенность (условие Мерсера)

$$\forall g(\mathbf{x}) : \int g^2(\mathbf{x}) d\mathbf{x} < \infty$$

$$\int K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0$$

Для фиксированной функции K евклидово пространство \mathcal{H} и преобразование Φ определено не однозначно.

Примеры ядровых функций

- Линейная ядровая функция

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle + \theta, \quad \theta \geq 0$$

- Полиномиальная ядровая функция

$$K(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \theta)^d, \quad \theta \geq 0, d \in \mathbb{N}$$

- Гауссиана

$$K(\mathbf{x}, \mathbf{y}) = \exp \left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2} \right), \quad \sigma > 0$$

- Сигмоидная ядровая функция

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\langle \mathbf{x}, \mathbf{y} \rangle + r), \quad r \in \mathbb{R}$$

Это семейство не удовлетворяет условию Мерсера!

4.2.5 Заключительные замечания

Пример использования метода опорных векторов

Зависимость от ширины гауссианы (см. рис. 4.10)

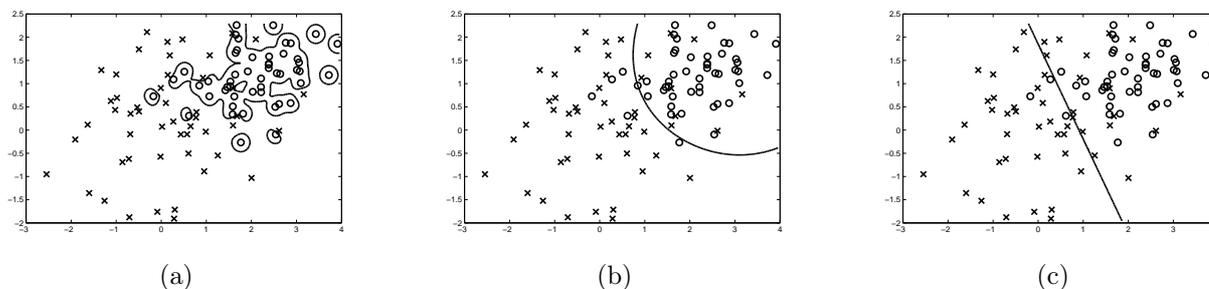


Рис. 4.10. Примеры решения двухклассовой задачи классификации с помощью метода опорных векторов. На рисунке (a) показана разделяющая поверхность для случая использования в качестве ядровой функции гауссианы с параметрами $C = 1$, $\sigma^2 = 0.1$, (b) соответствует $C = 1$, $\sigma^2 = 2$, (c) — $C = 1$, $\sigma^2 = 1000$

Зависимость от штрафного коэффициента (см. рис. 4.11)

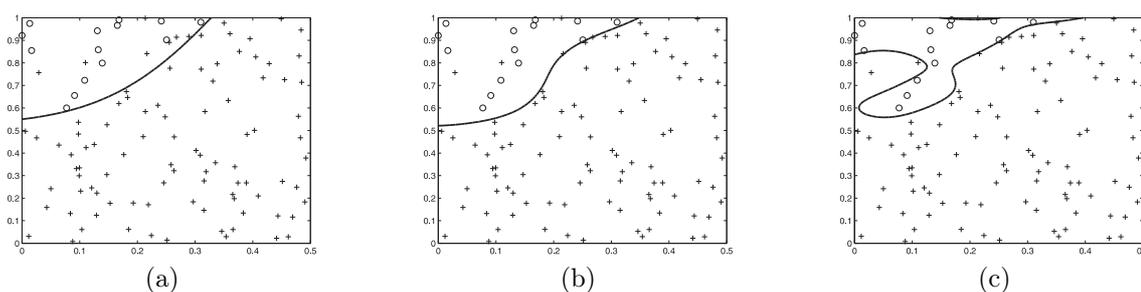


Рис. 4.11. Примеры решения двухклассовой задачи классификации с помощью метода опорных векторов. На рисунке (a) показана разделяющая поверхность для случая $C = 10^{-2}$, (b) соответствует $C = 1$, (c) — $C = 10^5$

Глобальность и единственность решения

- Задача обучения SVM — задача квадратичного программирования
- Известно, что для любой задачи выпуклого программирования (в частности, квадратичного) любой локальный максимум является и глобальным. Кроме того, решение будет единственным, если целевая функция строго вогнута (гессиян отрицательно определен).
- Для обучения SVM можно воспользоваться любым стандартным методом решения задачи квадратичного программирования, однако лучше использовать специальные алгоритмы, учитывающие особенности задачи квадратичного программирования в SVM (например, *SMO* или *SVM^{light}*).
- Подробнее см. <http://www.kernel-machines.org>

Задача обучения SVM как задача максимума регуляризованного правдоподобия

$$\begin{aligned} \frac{1}{2}\|\mathbf{z}\|^2 + C \sum_{i=1}^n \xi_i &\rightarrow \min_{\mathbf{z}, b, \xi} \\ t_i(\langle \mathbf{z}, \mathbf{x}_i \rangle + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}$$

Если $t_i y(\mathbf{x}_i) \geq 1$, то $\xi_i = 0$. Для остальных точек $\xi_i = 1 - t_i y(\mathbf{x}_i)$. Следовательно, оптимизируемую функцию можно переписать в виде

$$\sum_{i=1}^n E_{SV}(t_i y(\mathbf{x}_i)) + \lambda \|\mathbf{z}\|^2$$

Здесь $\lambda = (2C)^{-1}$, а $E_{SV}(\cdot)$ — функция потерь, определенная как

$$E_{SV}(s) = [1 - s]_+$$

SVM vs. Логистическая регрессия

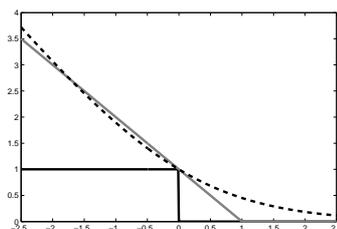


Рис. 4.12. Различные функции ошибок. Черная кривая соответствует индикаторной функции ошибки, серая кривая — функция ошибки в методе опорных векторов, пунктирная линия — функция ошибки в логистической регрессии

Задача оптимизации в логистической регрессии:

$$\sum_{i=1}^n E_{LR}(t_i y(\mathbf{x}_i)) + \lambda \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}$$

Здесь $E_{LR}(s) = \log(1 + \exp(-s))$.

Задача оптимизации в SVM:

$$\sum_{i=1}^n E_{SV}(t_i y(\mathbf{x}_i)) + \lambda \|\mathbf{z}\|^2 \rightarrow \min_{\mathbf{z}}$$

Достоинства и недостатки SVM

- + Высокое качество распознавания за счет построения нелинейных разделяющих поверхностей, максимизирующих зазор
- + Глобальность и в ряде случаев единственность получаемого решения
- Низкая скорость обучения и большие требования к памяти для задач больших размерностей
- Необходимость грамотного выбора штрафного коэффициента C и параметров ядровой функции

4.3 Метод опорных векторов для задачи регрессии

Линейная регрессия vs. SVR

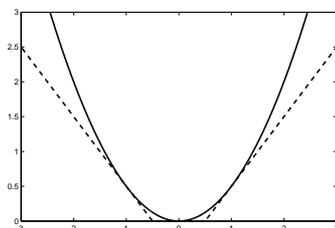


Рис. 4.13. Различные функции ошибок. Черная кривая — квадратичная функция ошибок линейной регрессии, пунктирная линия — функция ошибок в методе опорных векторов для регрессии

Задача оптимизации в линейной регрессии

$$\frac{1}{2} \sum_{i=1}^n (t_i - y(\mathbf{x}_i))^2 + \frac{1}{2} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}$$

Для того, чтобы добиться разреженного решения, заменим квадратичную функцию потерь на ε -нечувствительную:

$$E_{\varepsilon}(t - y(\mathbf{x})) = \begin{cases} 0, & \text{если } |t - y(\mathbf{x})| < \varepsilon \\ |t - y(\mathbf{x})| - \varepsilon, & \text{иначе} \end{cases}$$

Тогда мы приходим к следующей оптимизационной задаче:

$$C \sum_{i=1}^n E_{\varepsilon}(y(\mathbf{x}_i) - t_i) + \frac{1}{2} \|\mathbf{z}\|^2 \rightarrow \min_{\mathbf{z}}$$

Ослабляющие коэффициенты

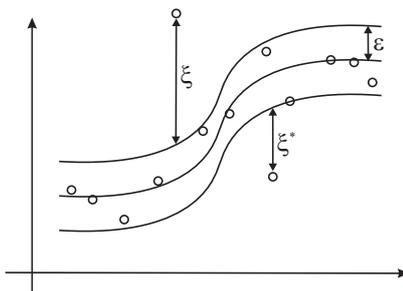


Рис. 4.14. Иллюстрация к введению ослабляющих коэффициентов. Объекты, лежащие выше ε -трубки, имеют положительное значение ξ_i , а объекты, лежащие ниже ε -трубки, имеют положительное значение ξ_i^*

$$C \sum_{i=1}^n (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{z}\|^2 \rightarrow \min_{\mathbf{z}, b, \xi, \xi^*}$$

$$t_i \leq y(\mathbf{x}_i) + \varepsilon + \xi_i$$

$$t_i \geq y(\mathbf{x}_i) - \varepsilon - \xi_i^*$$

$$\xi_i, \xi_i^* \geq 0$$

✓ Упр. Двойственная задача

$$-\frac{1}{2} \sum_{i,j=1}^n (w_i - w_i^*)(w_j - w_j^*)K(\mathbf{x}_i, \mathbf{x}_j) - \varepsilon \sum_{i=1}^n (w_i + w_i^*) + \sum_{i=1}^n t_i (w_i - w_i^*) \rightarrow \max_{\mathbf{w}, \mathbf{w}^*}$$

$$\sum_{i=1}^n (w_i - w_i^*) = 0$$

$$0 \leq w_i, w_i^* \leq C$$

Функция регрессии

$$y(\mathbf{x}) = \sum_{i=1}^n (w_i - w_i^*)K(\mathbf{x}, \mathbf{x}_i) + b$$

Условия дополняющей нежесткости

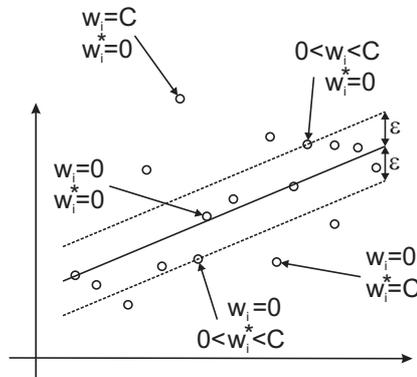


Рис. 4.15. Иллюстрация к опорным объектам в задаче регрессии

Запишем условия дополняющей нежесткости

$$w_i(\varepsilon + \xi_i - t_i + \langle \mathbf{z}, \mathbf{x}_i \rangle + b) = 0$$

$$w_i(\varepsilon + \xi_i^* + t_i - \langle \mathbf{z}, \mathbf{x}_i \rangle - b) = 0$$

$$(C - w_i)\xi_i = 0, \quad (C - w_i^*)\xi_i^* = 0$$

Из них следует, что опорные объекты лежат за пределами или на границе ε -трубки, определяемой функцией $\langle \mathbf{z}, \mathbf{x} \rangle + b$

4.4 Беспризнаковое распознавание образов

4.4.1 Основная методика беспризнакового распознавания образов

Задачи беспризнакового распознавания образов

Существует ряд задач распознавания образов, в которых трудно выбрать признаковое пространство, однако, относительно легко ввести меру сходства или несходства между парами объектов. Примеры:

- Задача распознавания личности по фотопортрету
- Задача идентификации личности по подписи в процессе ее формирования
- Задача распознавания классов пространственной структуры белков по последовательностям составляющих их аминокислот

Беспризнаковое распознавание образов: основная идея

- Предположим, что объекты выборки $\omega_1, \dots, \omega_n \in \Omega$
- Пространство Ω является гильбертовым, т.е. на нем определены операции суммы и произведения на число

$$\begin{aligned} \forall \alpha_1, \alpha_2 \in \Omega \exists \alpha = \alpha_1 + \alpha_2 \in \Omega \\ \forall \alpha_1 \in \Omega, c \in \mathbb{R} \exists \alpha = c\alpha_1 \in \Omega, \end{aligned}$$

удовлетворяющие аксиомам линейного пространства:

1. $\alpha_1 + \alpha_2 = \alpha_2 + \alpha_1$
2. $\alpha_1 + (\alpha_2 + \alpha_3) = (\alpha_1 + \alpha_2) + \alpha_3$
3. $\exists \phi \in \Omega : \alpha + \phi = \phi + \alpha = \alpha$
4. $\forall \alpha \exists (-\alpha) : \alpha + (-\alpha) = \phi$
5. $c(\alpha_1 + \alpha_2) = c\alpha_1 + c\alpha_2$
6. $(c + d)\alpha = c\alpha + d\alpha$
7. $(cd)\alpha = c(d\alpha)$
8. $1\alpha = \alpha$

- Существует функция $K : \Omega \times \Omega \rightarrow \mathbb{R}$, определяющая скалярное произведение в пространстве Ω :

1. $K(\alpha_1, \alpha_2) = K(\alpha_2, \alpha_1)$
2. $K(\alpha, \alpha) \geq 0, = 0 \Leftrightarrow \alpha = \phi$
3. $K(\alpha_1 + \alpha_2, \alpha) = K(\alpha_1, \alpha) + K(\alpha_2, \alpha)$
4. $K(c\alpha_1, \alpha_2) = cK(\alpha_1, \alpha_2)$

Решающая функция $\hat{t}(\omega) = \text{sign}(K(\vartheta, \omega) + b)$

Задача оптимизации $\frac{1}{2}K(\vartheta, \vartheta) + C \sum_{i=1}^n \xi_i \rightarrow \min_{\vartheta, b, \xi}$
 $t_i(K(\vartheta, \omega_i) + b) \geq 1 - \xi_i$
 $\xi_i \geq 0$

Двойственная задача оптимизации $\sum_{i=1}^n w_i - \frac{1}{2} \sum_{i,j=1}^n t_i t_j w_i w_j K(\omega_i, \omega_j) \rightarrow \max_{\mathbf{w}}$
 $\sum_{i=1}^n t_i w_i = 0$
 $0 \leq w_i \leq C$

Решение $\hat{t}(\omega) = \text{sign}(\sum_{i=1}^n t_i w_i K(\omega_i, \omega) + b)$

- Для решения задачи нет необходимости явно задавать линейные операции в пространстве Ω , достаточно лишь потребовать их существование в пространстве, где функция K задает скалярное произведение
- Для применения данного подхода достаточно задать функцию K . Эта функция может быть построена непосредственно, либо с использованием меры сходства

4.4.2 Построение функции, задающей скалярное произведение

Явное задание функции K . Пример. [Середин, 2001]

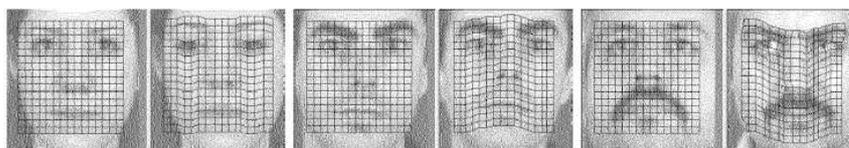


Рис. 4.16. Примеры фотографий лиц в задаче идентификации личности по фотопортрету. На фотографиях также представлено эластичное преобразование раstra при сравнении двух фотопортретов, получающееся в результате решения задачи минимизации искажений при условии минимального отклонения точек раstra от исходного положения

- Рассмотрим задачу распознавания личности по фотопортрету.
- Пусть два изображения заданы векторами яркости в точках раstra: $\mathbf{y}' = \mathbf{y}(\omega') = (y'_t, t \in T)$ и $\mathbf{y}'' = \mathbf{y}(\omega'') = (y''_t, t \in T)$, $T = \{t = (t_1, t_2), t_1 = \overline{1, n_1}, t_2 = \overline{1, n_2}\}$
- Потенциальная функция может быть задана как:

$$- K(\mathbf{y}', \mathbf{y}'') = \langle \mathbf{y}', \mathbf{y}'' \rangle = \sum_{t \in T} y'_t y''_t$$

$$- K(\mathbf{y}', \mathbf{y}'') = [\langle \mathbf{y}', \mathbf{y}'' \rangle + 1]^\alpha$$

$$- K(\mathbf{y}', \mathbf{y}'') = \exp(-\alpha \|\mathbf{y}' - \mathbf{y}''\|^2) = \exp(-\alpha [\langle \mathbf{y}', \mathbf{y}' \rangle + \langle \mathbf{y}'', \mathbf{y}'' \rangle - 2\langle \mathbf{y}', \mathbf{y}'' \rangle])$$

- Пусть задана эластичная деформация раstra $\mathbf{t} \rightarrow \mathbf{t} + \mathbf{x}_t$ (см. рис. 4.16). Эластичная потенциальная функция:

$$K(\mathbf{y}', \mathbf{y}'') = \sum_{t \in T} y'_t y''_{t+\mathbf{x}_t}$$

Построение функции K с использованием метрики

- Пусть задано метрическое пространство A с метрикой $\rho : A \times A \rightarrow \mathbb{R}$:

1. $\rho(\alpha_1, \alpha_2) \geq 0, = 0 \Leftrightarrow \alpha_1 = \alpha_2$
2. $\rho(\alpha_1, \alpha_2) = \rho(\alpha_2, \alpha_1)$
3. $\rho(\alpha_1, \alpha_3) \leq \rho(\alpha_1, \alpha_2) + \rho(\alpha_2, \alpha_3)$

- Общностью двух элементов из A относительно некоторого центра $\phi \in A$ назовем следующую величину:

$$\mu_\phi(\alpha_1, \alpha_2) = \frac{1}{2} [\rho^2(\alpha_1, \phi) + \rho^2(\alpha_2, \phi) - \rho^2(\alpha_1, \alpha_2)]$$

- Свойства общности

1. $\mu_\phi(\alpha_1, \alpha_2) = \mu_\phi(\alpha_2, \alpha_1)$
2. $\mu_\phi(\alpha, \alpha) \geq 0 = 0 \Leftrightarrow \alpha = \phi$
3. $\forall \alpha \in A \mu_\phi(\alpha, \phi) = 0$
4. $\mu_\phi(\alpha, \alpha) = \rho^2(\alpha, \phi)$
5. $\rho(\alpha_1, \alpha_2) = (\mu_\phi(\alpha_1, \alpha_1) + \mu_\phi(\alpha_2, \alpha_2) - 2\mu_\phi(\alpha_1, \alpha_2))^{1/2}$
6. $\mu_\phi(\alpha_1, \alpha_2) \leq \frac{1}{2} [\mu_\phi(\alpha_1, \alpha_1) + \mu_\phi(\alpha_2, \alpha_2)]$
7. $|\mu_\phi(\alpha_1, \alpha_2)| \leq \sqrt{\mu_\phi(\alpha_1, \alpha_1)} \sqrt{\mu_\phi(\alpha_2, \alpha_2)}$
8. $\mu_{\tilde{\phi}}(\alpha_1, \alpha_2) = \mu_\phi(\alpha_1, \alpha_2) - \mu_\phi(\alpha_1, \tilde{\phi}) - \mu_\phi(\alpha_2, \tilde{\phi}) + \mu_\phi(\phi, \tilde{\phi})$

- По своим свойствам общность очень похожа на скалярное произведение!

Матрица общностей

- Пусть $\{\alpha_1, \dots, \alpha_q\} \subset A$ — конечная совокупность элементов метрического пространства с центром $\phi \in A$. Составим матрицу общностей этих элементов:

$$M_\phi = (\mu_\phi(\alpha_i, \alpha_j), i, j = 1, \dots, q)$$

- Матрица общностей является симметричной, диагональные элементы неотрицательны.
- В отличие от матрицы скалярных произведений, которая всегда неотрицательно определена, матрица общностей может иметь собственные значения любого знака

Теорема 4. Если M_ϕ является неотрицательно определенной для любой конечной совокупности элементов, то матрица $M_{\tilde{\phi}}$ относительно другого центра $\tilde{\phi}$ также является неотрицательно определенной.

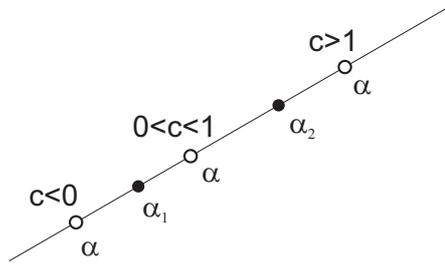


Рис. 4.17. Иллюстрация к понятию соосных элементов

Соосность элементов

Элемент $\alpha \in A$ называется **соосным элементом** для упорядоченной пары $[\alpha_1, \alpha_2]$ с коэффициентом $c \in \mathbb{R}$ и обозначается $\alpha = \text{соах}([\alpha_1, \alpha_2]; c)$, если

$$\rho(\alpha_1, \alpha) = |c|\rho(\alpha_1, \alpha_2), \quad \rho(\alpha_2, \alpha) = |c - 1|\rho(\alpha_1, \alpha_2)$$

Очевидно, что $\alpha_1 = \text{соах}([\alpha_1, \alpha_2]; 0)$, $\alpha_2 = \text{соах}([\alpha_1, \alpha_2]; 1)$. Если $\alpha = \text{соах}([\alpha_1, \alpha_2]; c)$, то $\alpha = \text{соах}([\alpha_2, \alpha_1]; 1 - c)$. Для элементов $[\alpha_1, \alpha_2]$ и $\alpha = \text{соах}([\alpha_1, \alpha_2]; c)$ неравенство треугольника переходит в равенство:

$$\begin{aligned} \rho(\alpha_1, \alpha_2) + \rho(\alpha_2, \alpha) &= \rho(\alpha_1, \alpha), & \text{если } c > 1 \\ \rho(\alpha_1, \alpha) + \rho(\alpha, \alpha_2) &= \rho(\alpha_1, \alpha_2), & \text{если } 0 < c \leq 1 \\ \rho(\alpha, \alpha_1) + \rho(\alpha_1, \alpha_2) &= \rho(\alpha, \alpha_2), & \text{если } c \leq 0 \end{aligned}$$

Введение линейных операций

Метрическое пространство A называется **евклидовым метрическим пространством**, если $\forall [\alpha_1, \alpha_2], \alpha_1, \alpha_2 \in A$ и $\forall c \in \mathbb{R} \exists \alpha \in A : \alpha = \text{соах}([\alpha_1, \alpha_2]; c)$ и матрица общности всякой конечной совокупности элементов из A является неотрицательно определенной.

Введем операции суммы и умножения на число:

$$c\alpha \triangleq \text{соах}([\phi, \alpha]; c) \quad \alpha_1 + \alpha_2 = 2 \text{соах} \left([\alpha_1, \alpha_2]; \frac{1}{2} \right)$$

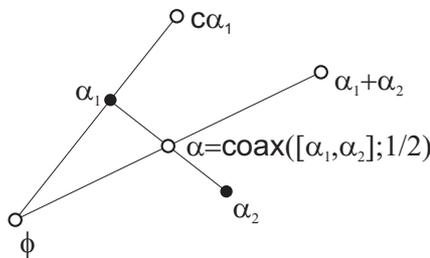


Рис. 4.18. Введение линейных операций в евклидовом метрическом пространстве

Свойства введенных операций

Теорема 5. В евклидовом метрическом пространстве введенные операции сложения и умножения на число, а также скалярное произведение понимаемое как общность элементов $\langle \alpha_1, \alpha_2 \rangle = \mu_\phi(\alpha_1, \alpha_2)$ удовлетворяют всем требованиям гильбертова пространства:

- $\alpha_1 + \alpha_2 = \alpha_2 + \alpha_1, (\alpha_1 + \alpha_2) + \alpha_3 = \alpha_1 + (\alpha_2 + \alpha_3)$
- $\alpha + \phi = \alpha, c\phi = \phi$
- $\forall \alpha \exists (-\alpha) : (-\alpha) + \alpha = \phi$
- $c_1(c_2\alpha) = (c_1c_2)\alpha$
- $1\alpha = \alpha$
- $(c_1 + c_2)\alpha = c_1\alpha + c_2\alpha, c(\alpha_1 + \alpha_2) = c\alpha_1 + c\alpha_2$
- $\langle \alpha_1, \alpha_2 \rangle = \langle \alpha_2, \alpha_1 \rangle, \langle c_1\alpha_1 + c_2\alpha_2, \alpha_3 \rangle = c_1\langle \alpha_1, \alpha_3 \rangle + c_2\langle \alpha_2, \alpha_3 \rangle$
- $\langle \alpha, \alpha \rangle \geq 0, = 0 \Leftrightarrow \alpha = \phi$
- $\|\alpha_1 - \alpha_2\| = \sqrt{\langle \alpha_1 - \alpha_2, \alpha_1 - \alpha_2 \rangle} = \rho(\alpha_1, \alpha_2)$

Глава 5

Задачи выбора модели

В главе рассматриваются вопросы выбора модели при обучении ЭВМ. Изложена суть проблемы и ее методологический характер. Приведены многочисленные примеры задач выбора модели, с которыми приходится сталкиваться при решении конкретных прикладных проблем. Рассмотрено несколько общих (не зависящих от эвристик) методов выбора модели, проанализированы их преимущества и недостатки.

5.1 Ликбез: Оптимальное кодирование

Оптимальное кодирование

- Рассматривается задача кодирования алфавита \mathcal{A} словами из алфавита \mathcal{B} , как правило содержащего меньше символов
- Пусть каждый символ $a \in \mathcal{A}$ встречается в текстах с вероятностью $p(a)$. Обозначим $l(a)$ длину его кодирования в \mathcal{B}
- Задача построить схему кодирования, обеспечивающую минимальную среднюю длину кодированных сообщений, т.е.

$$\mathbb{E}_{\mathcal{A}} l(a) = \sum_{a \in \mathcal{A}} p(a) l(a) \rightarrow \min$$

Теорема Шеннона

- Теорема Шеннона. Если кодируемые элементы a могут встречаться с разными вероятностями $p(a)$, то существует оптимальное кодирование данного множества элементов и длина описания элемента a равна $l(a) = -\log_B p(a)$
- Основание логарифма B — это мощность кодирующего алфавита. Если алфавит \mathcal{B} состоит из двух символов (например, ноль и единица), то логарифм двоичный, а длина описания измеряется в битах. Если логарифм натуральный, то длина описания измеряется в натах, хотя довольно трудно представить себе алфавит из 2.7182... символов :)
- Теорема Шеннона согласуется со здравым смыслом: чем чаще встречается символ, тем короче должна быть длина его описания, и наоборот. Лень Морзе дорого обошлась (и обходится) человечеству из-за того, что пропускная способность каналов связи оказалась ниже оптимальной

5.2 Постановка задачи выбора модели

5.2.1 Общий характер проблемы выбора модели

Определение модели

- Пусть $\mathcal{A}(w)$ — решающее правило, полученное в результате настройки весов $w \in \Omega$ в ходе обучения
- Модель — это совокупность всех решающих правил, которые получаются путем присваивания весам всех возможных допустимых значений $\mathcal{A}(\Omega)$
- Модель определяется множеством допустимых весов Ω , априорными распределениями весов на этом множестве $P(w)$ и структурой решающего правила $\mathcal{A}(\cdot)$

Задача выбора модели

- Выбрать модель — значит определить множество Ω , указать на нем априорные вероятности весов $P(w)$ и определить структуру (схему, по которой настраиваемые веса будут образовывать композицию с входными данными) решающего правила \mathcal{A}
- Поскольку рассмотреть всевозможные множества, распределения и структуры невозможно, их обычно ограничивают некоторым параметрическим семейством, зависящим от **структурных параметров**
- Под задачей выбора модели будем понимать проблему автоматической настройки всех структурных параметров для данного алгоритма машинного обучения

Смысл проблемы выбора модели

- Нетрудно придумать алгоритм, блестяще работающий на обучающей выборке (например, запомнить для нее правильные ответы). Вот только на новых объектах такой алгоритм, скорее всего, будет работать плохо
- Возникает идея ограничить множество допустимых решающих правил, чтобы в процессе обучения мы не могли бы получить «плохие» решения
- Ценой неизбежно становится ухудшение работы алгоритма на обучающей выборке
- **Процесс выбора модели в машинном обучении — это поиск компромисса между точностью на обучении и его «надежностью» на произвольных объектах генеральной совокупности**

Философский смысл проблемы выбора модели

- Данная проблема проявляется в различных областях математики в частности, и человеческой деятельности вообще
- Не будет преувеличением сказать, что она носит общенаучный и даже философский характер
- Это проблема выбора средств: гайку можно выточить на старом добром станке (в котором из механизмов только электродвигатель времен Александра II), а можно использовать новейшее оборудование с ЧПУ (три микропроцессора, система калибровки, связанная со спутником и 20 кг технической документации)
- Гайка, сделанная на новейшем станке, будет иметь более высокое качество. Разумеется, если не сломается система калибровки, не сгорит хотя бы один из трех процессоров, и за станок не сядет какой-нибудь «естествоиспытатель»...

5.2.2 Примеры задач выбора модели

Задача классификации методом опорных векторов

- Решающее правило имеет вид

$$y(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n w_i K(\mathbf{x}, \mathbf{x}_i) + b \right)$$

- Оптимизационная задача для поиска весов

$$\sum_{i=1}^n w_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n t_i t_j w_i w_j K(\mathbf{x}_i, \mathbf{x}_j) \rightarrow \max$$

$$\sum_{i=1}^n t_i w_i = 0 \quad 0 \leq w_i \leq C$$

- Коэффициент регуляризации C и ядровая функция $K(\mathbf{x}', \mathbf{x}'')$ определяют модель метода опорных векторов (см. рис. 5.1, 4.10, 4.11)

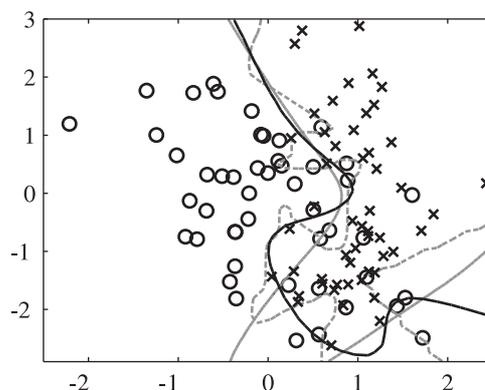


Рис. 5.1. Результаты классификации SVM с различными структурными параметрами

Задача регрессии

- Обобщенная линейная регрессия

$$y(\mathbf{x}) = \sum_{j=1}^m w_j \phi_j(\mathbf{x})$$

- Веса регрессии вычисляются по следующей формуле

$$\mathbf{w} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t}$$

- Параметр регуляризации $\lambda \geq 0$, система базисных функций $\{\phi_j(\mathbf{x})\}_{j=1}^m$ и их количество m определяют модель

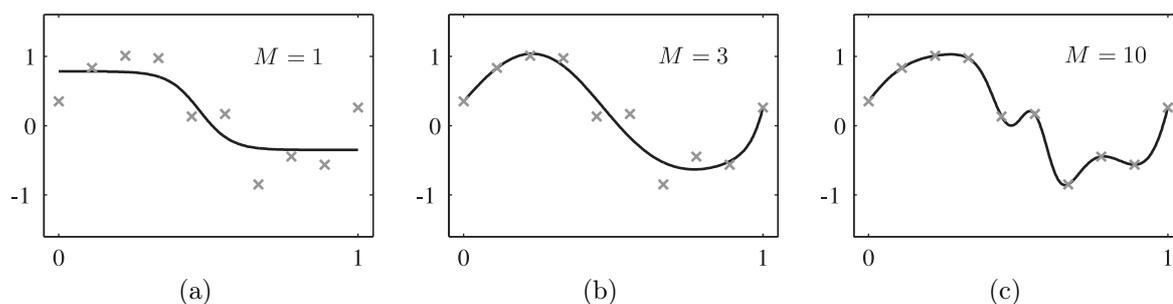


Рис. 5.2. Результаты восстановления регрессии с различным числом базисных функций

Задача кластеризации

- Большинство методов кластеризации предполагают задание пользователем количества кластеров, на которые будут разбиваться входные данные (см. рис. 5.3)

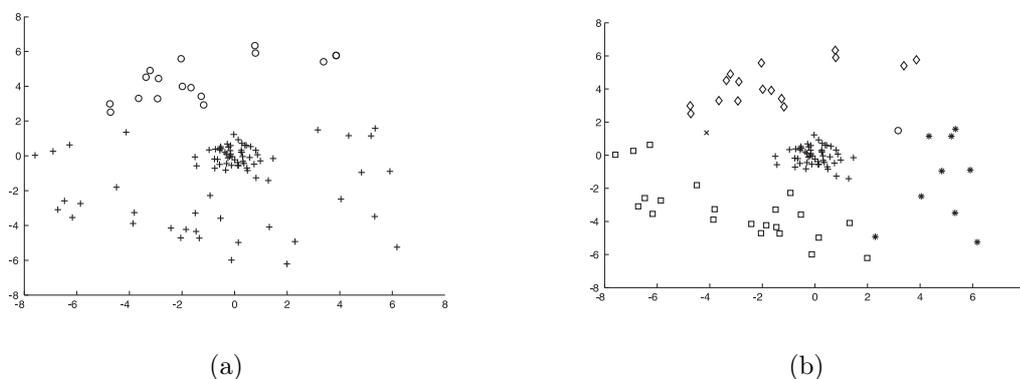


Рис. 5.3. Решение задачи кластеризации с помощью метода К средних с двумя кластерами (а) и с шестью кластерами (б)

Нейронные сети

- Выбор архитектуры нейронной сети (количество нейронов на каждом уровне) и функции активации определяют модель нейронной сети

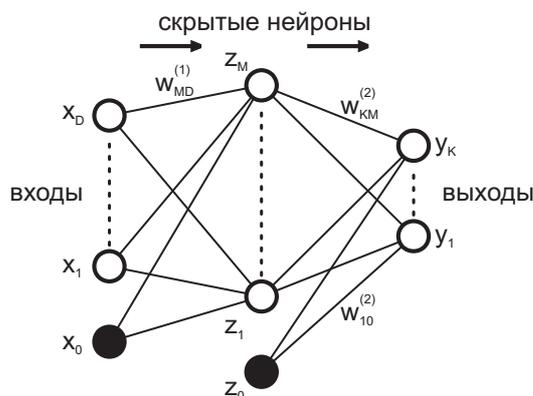


Рис. 5.4. Архитектура нейронной сети. Качества классификации нейронной сети существенно зависит от выбора ее структуры

5.3 Общие методы выбора модели

5.3.1 Кросс-валидация

Что такое общие методы выбора модели?

- Под общими методами выбора модели будем понимать алгоритмы, позволяющие проводить автоматическую настройку **любых** структурных параметров для широкого множества задач машинного обучения (например, для всех задач классификации)

- На практике такой метод может использоваться для настройки лишь некоторых структурных параметров, но гипотетически он должен позволять определять любые характеристики наилучшей модели

Скольльзящий контроль

- Процедура скользящего контроля (кросс-валидации) заключается в последовательном исключении части объектов из обучающей выборки, обучении на оставшихся объектах и распознавании исключенных объектов
- Тем самым эмулируется наличие тестовой выборки, которая не участвует в обучении, но для которой известны правильные ответы
- Структурные параметры настраиваются путем минимизации ошибки на скользящем контроле
- Процедура скользящего контроля является на сегодняшний день **самым лучшим средством настройки структурных параметров**

Схема k-fold cross validation

- Выборка разбивается на k непересекающихся (одинаковых по объему) частей. На каждой итерации обучение проводится по $k - 1$ части, а тестирование на исключенных объектах



Рис. 5.5. Процедура 4-fold cross validation. Каждый раз из выборки исключается четверть объектов, на остальной части проводится обучение, а затем исключенные объекты распознаются

- На рисунке 5.5 приведена процедура 4-fold cross validation
- При $k = n$ процедура называется leave-one-out
- Наилучшим режимом скользящего контроля считается 5×2 -fold cross validation.

Особенности скользящего контроля

- Ошибка на скользящем контроле является довольно точной оценкой ошибки на генеральной совокупности (обобщающей способности)
- Проведение скользящего контроля требует значительного времени на многократное повторное обучение алгоритмов и применимо лишь для «быстрых» методов машинного обучения
- С помощью скользящего контроля можно настраивать не более двух-трех структурных параметров, т.к. настройка производится путем **полного перебора** всевозможных сочетаний параметров
- При его использовании для выбора модели ошибку на скользящем контроле **нельзя** рассматривать как оценку ошибки на генеральной совокупности, т.к. она получается заниженной
- Скользящий контроль неприменим в задачах кластерного анализа и прогнозирования временных рядов

5.3.2 Теория Вапника-Червоненкиса

Идея теории структурной минимизации риска

- Теория Вапника-Червоненкиса использует косвенные характеристики для оценки обобщающей способности (среднего риска)
- Ключевым понятием является т.н. емкость (размерность Вапника-Червоненкиса, VC-dimension) модели
- Идея данного подхода к выбору модели (структурной минимизации риска) заключается в следующем: **Чем более «гибкой» является модель, тем хуже ее обобщающая способность**
- В самом деле, «гибкое» решающее правило способно настроиться на малейшие шумы, содержащиеся в обучающей выборке

Понятие емкости

- Рассмотрим задачу классификации на два класса
- (Несколько упрощая,) емкостью данной модели будем называть максимальное число объектов обучающей выборки, для которых при **любой** их разметке на классы найдется хотя бы один алгоритм из модели, безошибочно их классифицирующий
- По аналогии вводятся определения емкости для других задач машинного обучения
- Важный пример модели, для которой известна емкость — классификатор, строящий линейную гиперплоскость. Емкость линейного классификатора равна $h = d + 1$, где d — размерность пространства признаков
- Следствие: $n \leq d + 1$ объектов **всегда** можно безошибочно разделить гиперплоскостью

Формула Вапника

- Очевидно, что чем больше емкость, тем хуже. Значит нужно добиваться минимально возможного количества ошибок на обучении при минимальной возможной емкости
- Ошибку на обучении (эмпирический риск) $P_{train}(\mathbf{w})$, емкость $h(\Omega)$ и ошибку на генеральной совокупности (средний риск) $P_{test}(\mathbf{w})$ связывает известная формула Вапника

$$P_{test}(\mathbf{w}) \leq P_{train}(\mathbf{w}) + \sqrt{\frac{h(\Omega)(\log(2n/h(\Omega)) + 1) - \log(\eta/4)}{n}}$$

Неравенство верно с вероятностью $1 - \eta$ для $\forall \mathbf{w} \in \Omega$

- Последовательно анализируя модели с увеличивающейся емкостью, согласно теории ВЧ, необходимо выбирать модель с наименьшей верхней оценкой тестовой ошибки

Достоинства и недостатки теории ВЧ

- Достоинства
 - Серьезное теоретическое обоснование, связь с ошибкой на генеральной совокупности
 - Теория продолжает развиваться и в наши дни (эффективная емкость, локальная емкость, комбинаторный подход и т.д.)
- Недостатки
 - Оценки сильно завышены
 - Для большинства моделей емкость не поддается оценке
 - Многие модели с бесконечной емкостью показывают хорошие результаты на практике

Пути развития теории ВЧ (и не только ее)

- Емкость вводится для всевозможных положений объектов выборки, в то время как реально приходится иметь дело с одной конкретной выборкой и опять-таки имеет смысл рассматривать степень адаптируемости модели под эту конкретную выборку, а не под абстрактно возможную
- В процессе обучения поиск алгоритма в модели ведется не по всем ее представителям, а лишь по конечному числу, которое и имеет смысл рассматривать при интерпретации емкости, как степени адаптируемости модели под данные
- Искомая закономерность может обладать рядом дополнительных свойств, которые сокращают объем допустимых алгоритмов модели

5.3.3 Принцип минимальной длины описания

Предпосылка метода

- Из пункта А в пункт В передается закодированное сообщение о классификации обучающей выборки. Нужно добиться минимально возможного размера сообщения
- Стратегия 1: передаем каждый объект и его метку класса $\{(\mathbf{x}_i, t_i)\}_{i=1}^n$
- Стратегия 2: передаем длинное описание сложного алгоритма, который можно использовать для правильной классификации всей обучающей выборки $Descr(A)$
- Стратегия 3: передаем короткое описание простого алгоритма $Descr(A')$, который правильно классифицирует большинство объектов обучающей выборки, а классификацию неправильно распознанных объектов передаем отдельным списком $\{\mathbf{x}_{i_k}, t_{i_k}\}_{k=1}^p, p < n$

Смысл метода

- Чем точнее на обучающей выборке алгоритм, тем он сложнее, а значит тем длиннее будет его описание...
- ... но тем меньше будет список неправильно распознанных объектов (см. рис. 5.6)
- Принцип минимальной длины описания (minimum description length MDL, Rissanen, 1978) штрафует излишнюю алгоритмическую сложность решающего правила



Рис. 5.6. Иллюстрация метода минимальной длины описания

Особенности подсчета длины описания

- Существует множество подходов к оценке длины описания алгоритма вплоть до длины кода реализующей его программы
- Необходимо отметить, что кодирование должно быть эффективным, т.к. даже самый простой алгоритм можно закодировать в очень длинное сообщение
- Согласно теореме Шеннона, при оптимальном кодировании длина описания структуры пропорциональна логарифму ее вероятности, взятому с противоположным знаком

Эквивалентность MDL и максимизации апостериорной вероятности

- Пусть на множестве алгоритмов задано априорное распределение $p(\mathbf{w})$

$$l(\mathbf{w}) = -\log p(\mathbf{w})$$

- Длина описания данных тем меньше, чем выше вероятность данной классификации при использовании данного алгоритма, т.е. чем выше правдоподобие $p(\mathbf{t}|X, \mathbf{w})$

$$l(\mathbf{t}|\mathbf{w}) = -\log p(\mathbf{t}|X, \mathbf{w})$$

- Отсюда получаем выражение, объединяющее точность на обучении и сложность алгоритма **в единое выражение**

$$l(\mathbf{t}, \mathbf{w}) = -\log p(\mathbf{t}|X, \mathbf{w}) - \log p(\mathbf{w})$$

$$\arg \min_{\mathbf{w}} l(\mathbf{t}, \mathbf{w}) = \arg \max_{\mathbf{w}} p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})$$

- Таким образом, MDL обосновывает идею максимизации регуляризованного правдоподобия

Отличительные особенности MDL

- MDL позволяет обосновать корректность регуляризации правдоподобия
- Область применения MDL шире, чем у статистических методов обучения, т.е. MDL можно применять и там, где вводить вероятности некорректно или бессмысленно
- При использовании MDL предполагается, **что чем сложнее алгоритм, тем хуже его обобщающая способность**. Современные исследования (в частности, boosting) показывают, что это далеко не всегда так

5.3.4 Информационные критерии

Информационный критерий Акаике

- В 1973г. Акаике установил связь между правдоподобием (ключевое понятие статистики) и дивергенцией Кульбака-Лейблера (ключевое понятие в теории информации)
- Ему удалось получить приблизительное соотношение между правдоподобием генеральной совокупности и правдоподобием обучающей выборки (т.е. данных, по которым с помощью ММП производится настройка параметров решающего правила)

$$AIC = \log p(\mathbf{t}|X, \mathbf{w}_{ML}) - M,$$

где M — число настраиваемых параметров

- Пример использования: задача восстановления регрессии с известным гауссовским шумом в одномерном пространстве при помощи полинома степени k

$$k = \arg \min \left(\frac{\sum_{i=1}^n (t_i - y_k(\mathbf{x}_i))^2}{2\sigma^2} + k + 1 \right)$$

Информационный критерий Шварца

- Критерий Шварца (часто именуемый Байесовским информационным критерием) представляет собой простейшее приближение обоснованности, широко используемой в байесовском обучении

$$BIC \approx \int p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})d\mathbf{w}$$

- Используя приближение интеграла гауссианой и сильно огрубляя, получаем

$$BIC = \log p(\mathbf{t}|X, \mathbf{w}_{MP}) - \frac{1}{2}M \log n$$

- Пример использования: Задача восстановления регрессии с известным гауссовским шумом в одномерном пространстве при помощи полинома степени k

$$k = \arg \min \left(\frac{\sum_{i=1}^n (t_i - y_k(\mathbf{x}_i))^2}{2\sigma^2} + (k + 1) \log \sqrt{n} \right)$$

Особенности информационных критериев

- Оба критерия являются (весьма грубыми) приближениями более сложных выражений, часто не поддающихся аналитическому вычислению. Значение критерия может расцениваться лишь как приближительная характеристика обобщающей способности полученного решающего правила
- Критерии разумно использовать когда **все M настраиваемых параметров оказывают примерно одинаковое влияние** на вид решающего правила, например, входят в него линейно.
- Байесовский критерий сильнее штрафует сложные (с точки зрения дополнительных параметров) модели
- В байесовский критерий входит значение правдоподобия в точке \mathbf{w}_{MP} , а в критерий Акаике — в точке \mathbf{w}_{ML}