

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Лекция 1. Различные задачи машинного обучения

Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Концепция машинного обучения

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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.

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Классификация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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 $\{p(s|\mathbf{x})\}_{s=1}^l$

Классификация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

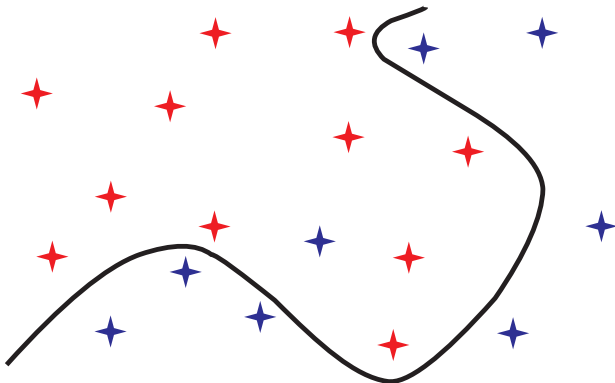
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач классификации

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Регрессия

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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 построить 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 $p(t|\mathbf{x})$

Регрессия

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

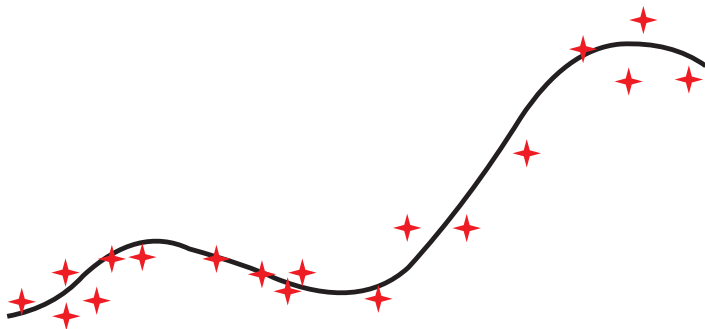
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач восстановления регрессии

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации
Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации
Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Кластеризация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации
Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации
Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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) $X = \bigcup_{j=1}^k C_j$, $C_j \subset \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, $C_i \cap C_j = \emptyset$
- В каждый класс должны попасть объекты в некотором смысле похожие друг на друга

Кластеризация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации
Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

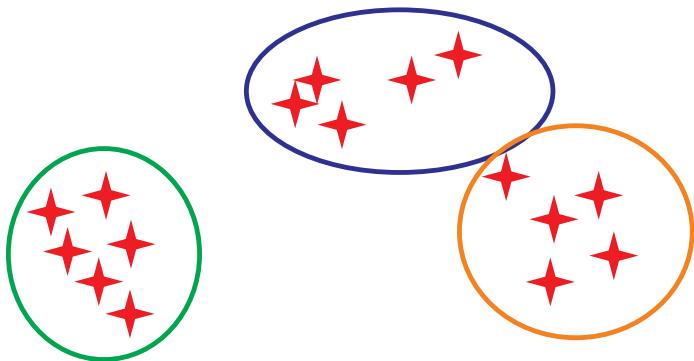
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач кластерного анализа

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации
Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации
Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Идентификация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

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

Идентификация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

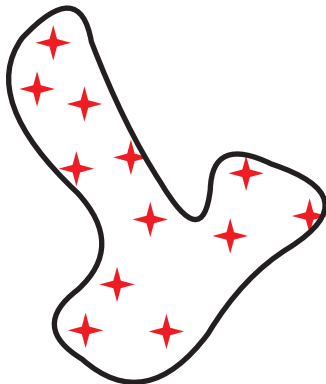
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач идентификации

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации
Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Прогнозирование

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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
- В отличие от задачи восстановления регрессии, здесь осуществляется прогноз **по времени**, а не по признакам

Прогнозирование

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

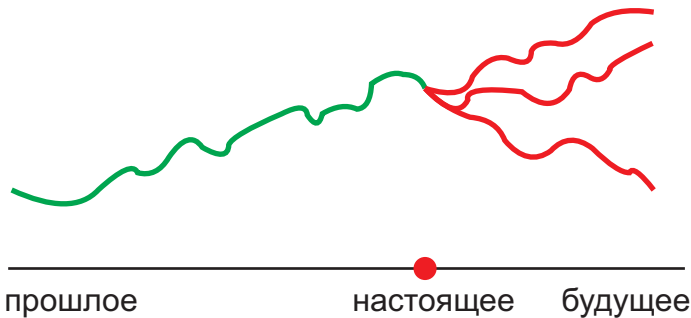
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач прогнозирования

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Извлечение знаний

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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 построить алгоритм, генерирующий набор объективных закономерностей между признаками, имеющих место в генеральной совокупности
- 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»)

Извлечение знаний

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

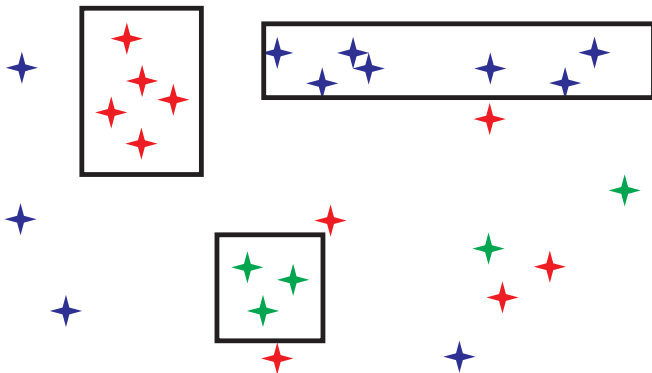
Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание



Примеры задач извлечения знаний

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Задача
классификации

Задача
восстановления
регрессии

Задача
кластеризации
(обучения без
учителя)

Задача
идентификации

Задача прогно-
зирования

Задача
извлечения
знаний

Основные
проблемы
машинного
обучения

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных
Переобучение

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Объем выборки I

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных
Переобучение

Напоминание

- 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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных
Переобучение

Напоминание

- Одна и та же выборка может являться большой для простых моделей алгоритмов и малой для сложных моделей.
- 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)

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Неполнота признакового описания

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

- 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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

- 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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных
Переобучение

Напоминание

- 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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки
Некорректность
входных данных
Переобучение

Напоминание

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Идея машинного обучения

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

- Задача машинного обучения заключается в восстановлении зависимостей по конечным выборкам данных (прецедентов)
- Пусть $(X, t) = (\mathbf{x}_i, t_i)_{i=1}^n$ – обучающая выборка, где $\mathbf{x}_i \in \mathbb{R}^d$ – признаковое описание объекта, а $t \in \mathcal{T}$ – значение скрытой компоненты (классовая принадлежность (не по Марксу!), 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

Проблема переобучения

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

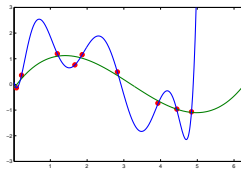
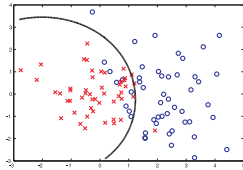
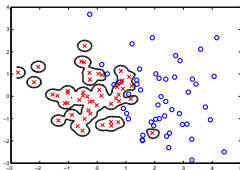
Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

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



Способы оценки и увеличения обобщающей способности

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

- На сегодняшний день единственным универсальным способом оценивания обобщающей способности является скользящий контроль
- 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 (Х. Акаике, 1974, Г. Шварц, 1978)
 - Максимизация обоснованности (Д. МакКай, 1992)
- Последний принцип позволяет надеяться на конструктивное решение задачи выбора модели

Примеры задач выбора модели

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Малый объем
обучающей
выборки

Некорректность
входных данных

Переобучение

Напоминание

- Определение числа кластеров в данных
- Выбор коэффициента регуляризации в задаче машинного обучения
- Установка степени полинома при интерполяции сплайнами
- Выбор наилучшей базисной функции в обобщенных линейных моделях
- Определение количества ветвей в решающем дереве
- и многое другое...

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Матричная нотация

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

- При работе с многомерными величинами очень удобна матричная нотация, т.е. представление многих операций над векторами и числами в виде операций над матрицами
- Скалярное произведение двух векторов $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ принимает вид

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^d x_i y_i = \mathbf{x}^T \mathbf{y},$$

т.е. вектора трактуются как частные случаи матриц

- Квадратичная форма

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^d \sum_{j=1}^d x_i a_{ij} y_j = \mathbf{x}^T \mathbf{A} \mathbf{y}$$

- Матричная нотация облегчает математические выкладки и позволяет реализовать вычисления на ЭВМ более эффективно

Пример использования

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

- Предположим нам надо решить несовместную систему линейных уравнений $Ax \approx b$, $A \in \mathbb{R}^{m \times n}$
- Для этого будем минимизировать квадрат нормы невязки (система-то нерешаемая) $\|Ax - b\|^2 \rightarrow \min_x$
- Представляя норму в матричной виде, дифференцируя по вектору и приравнивая производную к нулю получаем известную формулу для псевдорешения СЛАУ

$$\begin{aligned}\|Ax - b\|^2 &= \langle Ax - b, Ax - b \rangle = (Ax - b)^T (Ax - b) = \\ &= (Ax)^T Ax - b^T Ax - (Ax)^T b + b^T b = \\ &= x^T A^T Ax - 2x^T A^T b + b^T b\end{aligned}$$

$$\frac{\partial}{\partial x} (x^T A^T Ax - 2x^T A^T b + b^T b) = 2A^T Ax - 2A^T b = 0$$

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

- Заметим, что если матрица A квадратная (число уравнений равно числу неизвестных) и невырожденная, то последняя формула переходит в формулу обычного решения СЛАУ $x = A^{-1}b$

Задача условной оптимизации

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

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

$$f(\mathbf{x}) \rightarrow \underset{\mathbf{x}}{\text{extr}}$$

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:

$$f(\mathbf{x}) \rightarrow \underset{\mathbf{x}}{\text{extr}}$$

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

Поверхность ограничения

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

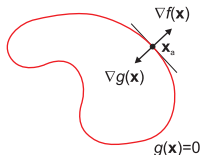
Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики



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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

Необходимым условием оптимальности является ортогональность $\nabla f(\mathbf{x})$ поверхности ограничения, т.е.:

$$\nabla f + \lambda \nabla g = 0 \quad (1)$$

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

Функция Лагранжа

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

Тогда

$$\nabla_{\mathbf{x}} L = 0 \quad \Rightarrow \quad \text{условие (1)}$$

$$\frac{\partial}{\partial \lambda} L = 0 \quad \Rightarrow \quad g(\mathbf{x}) = 0$$

Функция Лагранжа. Пример.

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

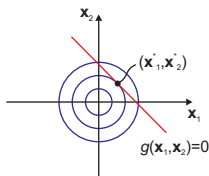
Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики



$$f(x_1, x_2) = 1 - x_1^2 - x_2^2 \rightarrow \max_{x_1, x_2}$$

$$g(x_1, x_2) = x_1 + x_2 - 1 = 0$$

Функция Лагранжа:

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

Условия стационарности:

$$-2x_1 + \lambda = 0$$

$$-2x_2 + \lambda = 0$$

$$x_1 + x_2 - 1 = 0$$

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

План лекции

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

1 Некоторые задачи машинного обучения

Задача классификации

Задача восстановления регрессии

Задача кластеризации (обучения без учителя)

Задача идентификации

Задача прогнозирования

Задача извлечения знаний

2 Основные проблемы машинного обучения

Малый объем обучающей выборки

Некорректность входных данных

Переобучение

3 Напоминание

Полезные сведения из линейной алгебры и теории оптимизации

Основные понятия мат. статистики

Краткое напоминание основных вероятностных понятий

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

- $X : \Omega \rightarrow \mathbb{R}$ – 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

Нормальное распределение

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

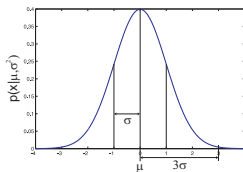
Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

- Нормальное распределение играет важнейшую роль в математической статистике

$$X \sim \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$
$$\mu = \mathbb{E}X, \quad \sigma^2 = \mathbb{D}X \triangleq \mathbb{E}(X - \mathbb{E}X)^2$$



- 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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

- Многомерное нормальное распределение имеет вид

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

где $\mu = \mathbb{E}X$, $\Sigma = \mathbb{E}(X - \mu)(X - \mu)^T$ — 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

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

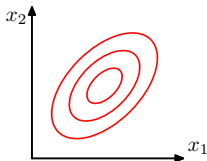
Основные
проблемы
машинного
обучения

Напоминание

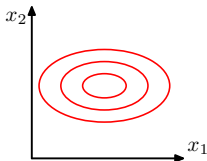
Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

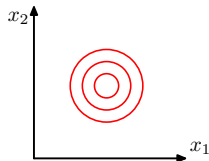
- Нормальное распределение **полностью задается** первыми двумя моментами (мат. ожидание и матрица ковариаций/дисперсия)
- Матрица ковариаций неотрицательно определена, причем на диагоналях стоят дисперсии соответствующих компонент
- Нормальное распределение имеет очень легкие хвосты: большие отклонения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



(a)



(b)



(c)

Основная задача мат. статистики

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия 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 (не путать с одноименной кафедрой) является асимптотически оптимальным при $n \rightarrow \infty$
- Обычно максимизируют не само правдоподобие, а его логарифм, т.к. это вычислительно проще (произведение плотностей по всем объектам переходит в сумму логарифмов плотностей)

Пример использования

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

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

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

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

$$\mu_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\begin{aligned} \sigma_{ML}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \\ &= \frac{1}{n} \sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{j=1}^n x_j \right)^2 \end{aligned}$$

Выводы

Лекция 1.
Различные
задачи
машинного
обучения

Журавлев,
Ветров

Некоторые
задачи
машинного
обучения

Основные
проблемы
машинного
обучения

Напоминание

Полезные
сведения из
линейной
алгебры и
теории
оптимизации

Основные
понятия мат.
статистики

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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Обобщенные линейные модели

Ю. И. Журавлев¹, Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

Условная вероятность

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

- Пусть X и Y — случайные величины с плотностями $p(x)$ и $p(y)$ соответственно
- В общем случае их совместная плотность $p(x, y) \neq p(x)p(y)$. Если это равенство выполняется, величины называют **независимыми**
- Условной плотностью называется величина

$$p(x|y) = \frac{p(x, y)}{p(y)}$$

- Смысл: как факт $Y = y$ влияет на распределение X . Заметим, что $\int p(x|y)dx \equiv 1$, но $\int p(x|y)dy$ не обязан равняться единице, т.к. относительно y это не плотность, а **функция правдоподобия**
- Очевидная система тождеств $p(x|y)p(y) = p(x, y) = p(y|x)p(x)$ позволяет легко переходить от $p(x|y)$ к $p(y|x)$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Правило суммирования вероятностей

- Все операции над вероятностями базируются на применении всего двух правил
- Правило суммирования: Пусть A_1, \dots, A_k взаимоисключающие события, одно из которых **всегда происходит**. Тогда

$$P(A_i \cup A_j) = P(A_i) + P(A_j) \quad \sum_{i=1}^k P(A_i) = 1$$

- Очевидное следствие (формула полной вероятности): $\forall B$ верно $\sum_{i=1}^k P(A_i|B) = 1$, откуда

$$\sum_{i=1}^k \frac{P(B|A_i)P(A_i)}{P(B)} = 1 \quad P(B) = \sum_{i=1}^k P(B|A_i)P(A_i)$$

- В интегральной форме

$$p(b) = \int p(b, a) da = \int p(b|a)p(a) da$$

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Правило произведения вероятностей

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

- Правило произведения гласит, что любую совместную плотность всегда можно разбить на множители

$$p(a, b) = p(a|b)p(b) \quad P(A, B) = P(A|B)P(B)$$

- Аналогично для многомерных совместных распределений

$$p(a_1, \dots, a_n) =$$

$$p(a_1|a_2, \dots, a_n)p(a_2|a_3, \dots, a_n) \dots p(a_{n-1}|a_n)p(a_n)$$

- Можно показать (Jaynes, 1995), что правила суммирования и произведения вероятностей являются единственными возможными операциями, позволяющими рассматривать вероятности как промежуточную ступень между истиной и ложью

Априорные и апостериорные суждения

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

- Предположим, мы пытаемся изучить некоторое явление
- У нас имеются некоторые знания, полученные до (лат. a priori) наблюдений/эксперимента. Это может быть опыт прошлых наблюдений, какие-то модельные гипотезы, ожидания
- В процессе наблюдений эти знания подвергаются постепенному уточнению. После (лат. a posteriori) наблюдений/эксперимента у нас формируются новые знания о явлении
- Будем считать, что мы пытаемся оценить неизвестное значение величины θ посредством наблюдений некоторых ее косвенных характеристик $x|\theta$

Формула Байеса

- Знаменитая формула Байеса (1763 г.) устанавливает правила, по которым происходит преобразование знаний в процессе наблюдений
- Обозначим априорные знания о величине θ за $p(\theta)$
- В процессе наблюдений мы получаем серию значений $\mathbf{x} = (x_1, \dots, x_n)$. При разных θ наблюдение выборки \mathbf{x} более или менее вероятно и определяется значением правдоподобия $p(\mathbf{x}|\theta)$
- За счет наблюдений наши представления о значении θ меняются согласно формуле Байеса

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta)d\theta}$$

- Заметим, что знаменатель не зависит от θ и нужен исключительно для нормировки апостериорной плотности

Обобщенные
линейные модели

Ветров

Напоминание

Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

План лекции

Обобщенные
линейные модели

Ветров

Напоминание
Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

Псевдообращение матриц

Обобщенные
линейные модели

Ветров

Напоминание
Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

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

$$A^T A \mathbf{x} = A^T \mathbf{b}$$

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

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

- Операция $(A^T A)^{-1} A^T$ называется псевдообращением матрицы A , а \mathbf{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

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

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

Обобщенные
линейные модели

Ветров

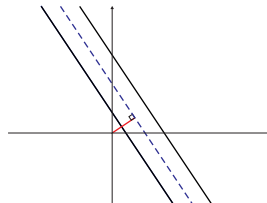
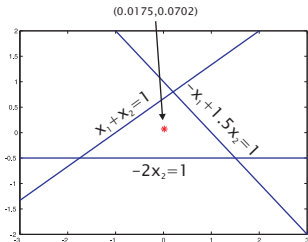
Напоминание
Формула Байеса
Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Графическая иллюстрация

- Псевдорешение соответствует точке, минимизирующей невязку, а нормальное псевдорешение отвечает псевдорешению с наименьшей нормой



- Заметим, что псевдообратная матрица $(A^T A)^{-1} A^T$ совпадает с обратной матрицей A^{-1} в случае невырожденных квадратных матриц

Обобщенные
линейные модели

Ветров

Напоминание
Формула Байеса

Решение
нерешаемых
систем
уравнений

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

Задача восстановления регрессии

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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

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

Линейная регрессия

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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**

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

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

$$\begin{aligned} \frac{\partial \|\Phi \mathbf{w} - \mathbf{t}\|^2}{\partial \mathbf{w}} &= \frac{\partial [\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t}]}{\partial \mathbf{w}} = \\ &= 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{t} = 0 \\ \mathbf{w} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned}$$

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

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 (Упр.) 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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Особенности квадратичной функции потерь

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

Нормальное распределение ошибок

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

- 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 (не путать с одноименной кафедрой) 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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

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

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

- Регуляризация эквивалентна введению априорного распределения, поощряющего небольшие веса

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

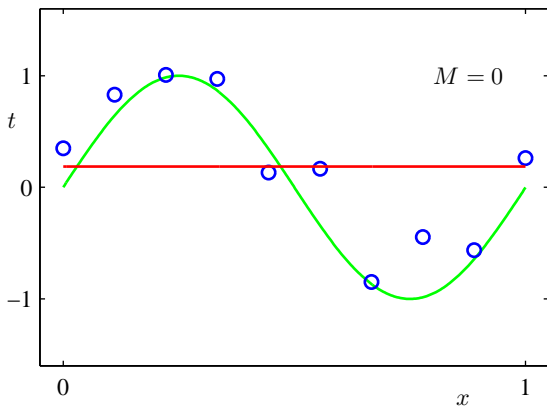
Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями: $x \in \mathbb{R}$, $\phi_j(x) = x^j$, $j = 0, \dots, M$



Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

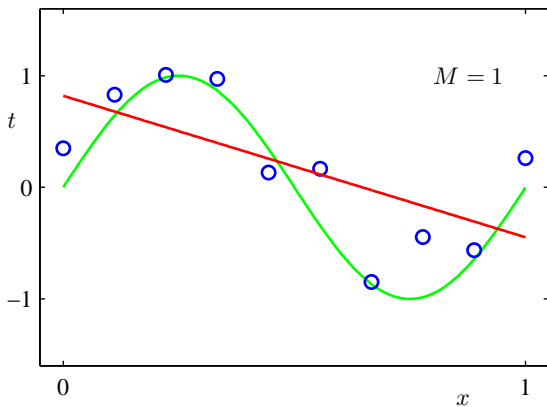
Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями: $x \in \mathbb{R}$, $\phi_j(x) = x^j$, $j = 0, \dots, M$



Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

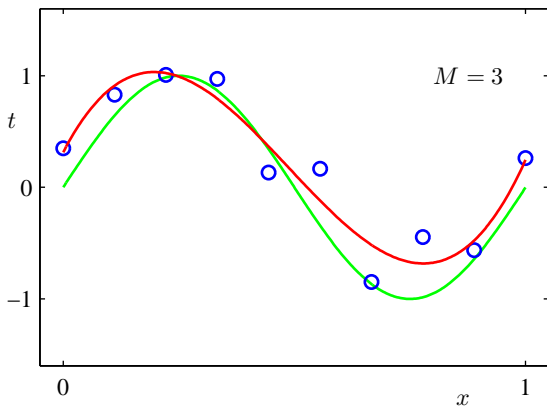
Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями: $x \in \mathbb{R}$, $\phi_j(x) = x^j$, $j = 0, \dots, M$



Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

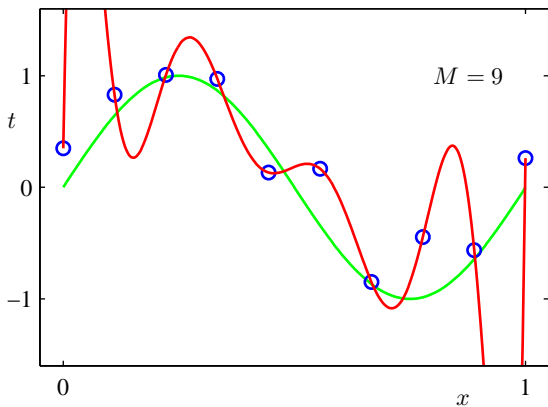
Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями: $x \in \mathbb{R}$, $\phi_j(x) = x^j$, $j = 0, \dots, M$



Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия
Классическая
линейная
регрессия
Метод
наименьших
квадратов
Вероятностная
постановка
задачи
Применение
регрессионных
методов для
задачи
классификации

Значения наиболее правдоподобных весов

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Классическая
линейная
регрессия

Метод
наименьших
квадратов

Вероятностная
постановка
задачи

Применение
регрессионных
методов для
задачи
классификации

weight	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0	0.19	0.82	0.31	0.35
w_1		-1.27	7.99	232.37
w_2			-25.43	-5321.83
w_3			17.37	48568.31
w_4				-231639.30
w_5				640042.26
w_6				-1061800.52
w_7				1042400.18
w_8				-557682.99
w_9				125201.43

Таблица: Значения наиболее правдоподобных весов в зависимости от степени полинома. С увеличением степени, абсолютные значения весов быстро растут

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

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}

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS

Правдоподобие правильной классификации

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS

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

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

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

Функционал качества в логистической регрессии

Обобщенные
линейные модели

Ветров

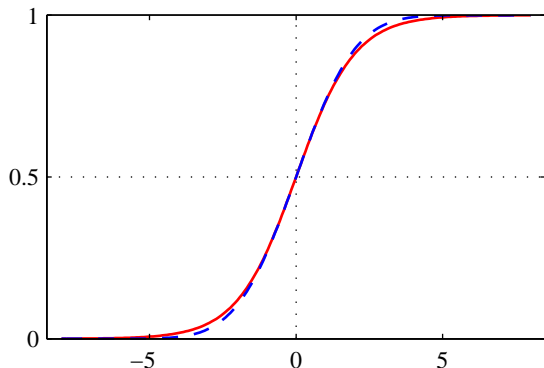
Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS



- Правдоподобие правильной классификации всей выборки имеет вид

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

План лекции

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи

классификации
Логистическая
регрессия

Метод IRLS

Напоминание

Формула Байеса

Решение нерешаемых систем уравнений

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Применение регрессионных методов для задачи классификации

Логистическая регрессия

Метод IRLS

Особенности функции правдоподобия классификации

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

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

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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS

Метод оптимизации Ньютона

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия
Метод IRLS

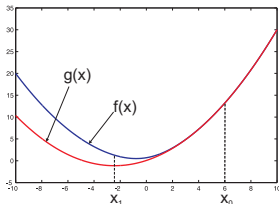
Основная идея метода Ньютона — это приближение в заданной точке оптимизируемой функции параболой и выбор минимума этой параболы в качестве следующей точки итерационного процесса:

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

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

$$\nabla g(\mathbf{x}_*) = \nabla f(\mathbf{x}_0) + (\nabla \nabla f(\mathbf{x}_0))(\mathbf{x}_* - \mathbf{x}_0) = 0 \Rightarrow \mathbf{x}_* = \mathbf{x}_0 - (\nabla \nabla f(\mathbf{x}_0))^{-1}(\nabla f(\mathbf{x}_0))$$

Пример. Функция $f(x) = \log(1 + \exp(x)) + \frac{x^2}{5}$.
 $x_0 = 6$, $x_1 = -2.4418$.



Итеративная минимизация логарифма правдоподобия

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации
Логистическая
регрессия

Метод IRLS

- 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:

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

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

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

где $\mathbf{z} = \Phi \mathbf{w}^{old} - R^{-1} (\mathbf{s} - \mathbf{t}^*)$, $t_i^* = \mathcal{H}(t_i)$

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

Обобщенные
линейные модели

Ветров

Напоминание

Линейная
регрессия

Применение
регрессионных
методов для
задачи
классификации

Логистическая
регрессия

Метод IRLS

- 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

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Лекция 2. Графические модели. Общее представление

Ю. И. Журавлев¹, Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Ликбез

Графические
модели

Байесовские сети

Марковские сети

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Условная вероятность

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Пусть X и Y — случайные величины с плотностями $p(x)$ и $p(y)$ соответственно
- В общем случае их совместная плотность $p(x, y) \neq p(x)p(y)$. Если это равенство выполняется, величины называют **независимыми**
- Условной плотностью называется величина

$$p(x|y) = \frac{p(x, y)}{p(y)}$$

- Смысл: как факт $Y = y$ влияет на распределение X . Заметим, что $\int p(x|y)dx \equiv 1$, но $\int p(x|y)dy$ не обязан равняться единице, т.к. относительно y это не плотность, а **функция правдоподобия**
- Очевидная система тождеств $p(x|y)p(y) = p(x, y) = p(y|x)p(x)$ позволяет легко переходить от $p(x|y)$ к $p(y|x)$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Правило суммирования вероятностей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Все операции над вероятностями базируются на применении всего двух правил
- Правило суммирования: Пусть A_1, \dots, A_k взаимоисключающие события, одно из которых **всегда происходит**. Тогда

$$P(A_i \cup A_j) = P(A_i) + P(A_j) \quad \sum_{i=1}^k P(A_i) = 1$$

- Очевидное следствие (формула полной вероятности): $\forall B$ верно $\sum_{i=1}^k P(A_i|B) = 1$, откуда

$$\sum_{i=1}^k \frac{P(B|A_i)P(A_i)}{P(B)} = 1 \quad P(B) = \sum_{i=1}^k P(B|A_i)P(A_i)$$

- В интегральной форме

$$p(b) = \int p(b, a) da = \int p(b|a)p(a) da$$

Правило произведения вероятностей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Правило произведения гласит, что любую совместную плотность всегда можно разбить на множители

$$p(a, b) = p(a|b)p(b) \quad P(A, B) = P(A|B)P(B)$$

- Аналогично для многомерных совместных распределений

$$p(a_1, \dots, a_n) =$$

$$p(a_1|a_2, \dots, a_n)p(a_2|a_3, \dots, a_n) \dots p(a_{n-1}|a_n)p(a_n)$$

- Можно показать (Jaunes, 1995), что правила суммирования и произведения вероятностей являются единственными возможными операциями, позволяющими рассматривать вероятности как промежуточную ступень между истиной и ложью

Априорные и апостериорные суждения

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Предположим, мы пытаемся изучить некоторое явление
- У нас имеются некоторые знания, полученные до (лат. a priori) наблюдений/эксперимента. Это может быть опыт прошлых наблюдений, какие-то модельные гипотезы, ожидания
- В процессе наблюдений эти знания подвергаются постепенному уточнению. После (лат. a posteriori) наблюдений/эксперимента у нас формируются новые знания о явлении
- Будем считать, что мы пытаемся оценить неизвестное значение величины θ посредством наблюдений некоторых ее косвенных характеристик $x|\theta$

Формула Байеса

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Знаменитая формула Байеса (1763 г.) устанавливает правила, по которым происходит преобразование знаний в процессе наблюдений
- Обозначим априорные знания о величине θ за $p(\theta)$
- В процессе наблюдений мы получаем серию значений $\mathbf{x} = (x_1, \dots, x_n)$. При разных θ наблюдение выборки \mathbf{x} более или менее вероятно и определяется значением правдоподобия $p(\mathbf{x}|\theta)$
- За счет наблюдений наши представления о значении θ меняются согласно формуле Байеса

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta)d\theta}$$

- Заметим, что знаменатель не зависит от θ и нужен исключительно для нормировки апостериорной плотности

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Условная независимость случайных величин

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез
Формула Байеса
Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Случайные величины x и y называются условно независимыми от z , если

$$p(x, y|z) = p(x|z)p(y|z)$$

- Другими словами вся информация о взаимозависимостях между x и y содержится в z
- Заметим, что из безусловной независимости не следует условная и наоборот
- Основное свойство условно независимых случайных величин

$$p(z|x, y) = \frac{p(x, y|z)p(z)}{p(x, y)} = \frac{p(x|z)p(y|z)p(z)}{p(x, y)} =$$

$$\frac{p(x|z)p(z)p(y|z)p(z)}{p(x, y)p(z)} = \frac{p(z|x)p(z|y)}{p(z)p(x)p(y)p(x, y)} = \frac{1}{Z} \frac{p(z|x)p(z|y)}{p(z)}$$

Пример

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Формула Байеса

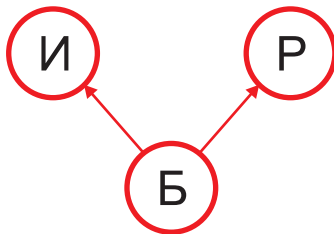
Условная
независимость
случайных
величин

Графические
модели

Байесовские сети

Марковские сети

- Рассмотрим следующую гипотетическую ситуацию: римские легионы во главе с императором атакуют вторгшихся варваров
- События «гибель императора» и «уничтожение Рима» не являются независимыми
- Однако, если нам дополнительно известен исход битвы с варварами, эти два события становятся независимыми
- В самом деле, если легионы битву проиграли, то судьба Рима мало зависит от того, был ли император убит в сражении



План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Классическая задача машинного обучения

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

- Задачу машинного обучения можно трактовать как восстановление неизвестных зависимостей между наблюдаемыми переменными X и скрытыми (латентными) переменными T . В случае обучения с учителем такое восстановление производится по обучающей выборке Y
- В классических задачах машинного обучения предполагается, что обучающая выборка сформирована из **однородных и независимых** объектов $Y = \{(\mathbf{x}_i, \mathbf{t}_i)\}_{i=1}^n$
- До недавнего времени вероятностные методы обработки данных ограничивались только таким простейшим случаем, а изложение каждого метода начиналось со слов «Предположим, что нам дана выборка из независимых одинаково распределенных случайных величин...»

Задачи со структурными ограничениями

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

- Во многих задачах взаимосвязи между наблюдаемыми и скрытыми переменными носят сложный характер
- В частности, между отдельными переменными существуют вероятностные зависимости
- Факт зависимости переменных друг от друга удобно отображать с помощью неориентированного графа (марковской сети)
- Если связи между переменными причинно-следственные, то их удобно отображать в виде ориентированных графов (байесовских сетей)
- Основным средством работы с графическими моделями служит аппарат теории вероятностей, в ее байесовской интерпретации

Простой пример

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

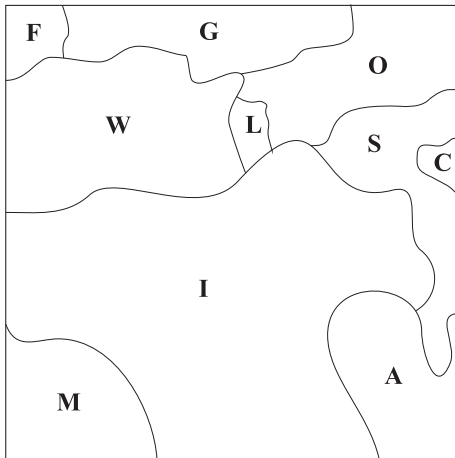
Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Задача о раскраске областей на плоскости так, чтобы никакие соседние не были окрашены в одинаковый цвет



Простой пример

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

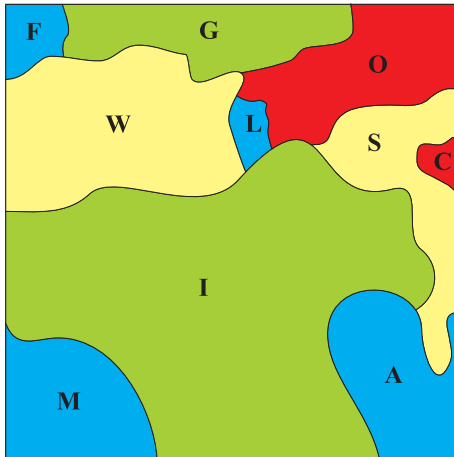
Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Задача о раскраске областей на плоскости так, чтобы никакие соседние не были окрашены в одинаковый цвет



Простой пример

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

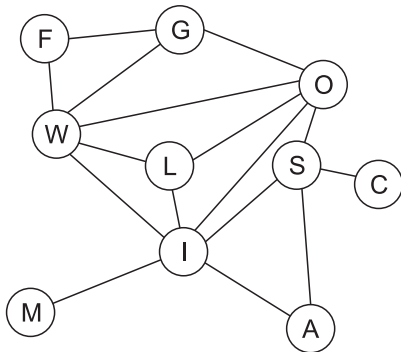
Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Такая задача легко формулируется в терминах графической модели, в которой каждая вершина графа может находиться в одном из четырех состояний



Вопрос залу: почему четырех?

Примеры задач со структурными связями

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

- Обработка изображений, сигналов
- Анализ социальных сетей
- Поиск залежей полезных ископаемых
- Анализ естественных языков
- Биомедицина и биоинформатика
- Веб-поиск
- и др.

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Графические модели

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

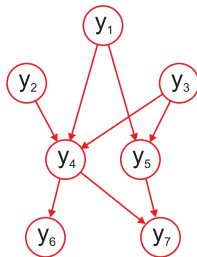
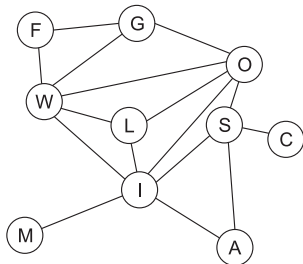
Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

- Графическая модель представляет собой ориентированный или неориентированный граф
- Вершины графа соответствуют переменным
- Ребра графа соответствуют вероятностным отношениям, определяющим непосредственные зависимости



Главные задачи в анализе графических моделей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

Обозначим совокупность наблюдаемых переменных X , а ненаблюдаемых переменных T . Основными задачами в анализе графических моделей являются

- Подсчет условного распределения на значения отдельной скрытой переменной $p(t_i|X)$ —?
- Нахождение наиболее вероятной конфигурации скрытых переменных $p(T|X) \rightarrow \max_T$
- Оценка адекватности выбранной графической модели данным $p(X)$ —?

Трудности, возникающие при использовании графических моделей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Задачи со
структурными
ограничениями

Основные
проблемы в
анализе
графических
моделей

Байесовские сети

Марковские сети

- Не во всех случаях существуют строгие алгоритмы вывода и обучения графических моделей
- Даже там, где они существуют, их применение может оказаться невозможно из-за высоких вычислительных требований и требований к памяти
- В настоящее время в мире активно разрабатываются приближенные эффективные методы обучения и принятия решения в графических моделях (Monte Carlo Markov chains, Variational bounds, Expectation propagation, Belief propagation, и др.)

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа
Пример
использования

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клик

Пример использования

Связь с байесовскими сетями

Совместное распределение переменных

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

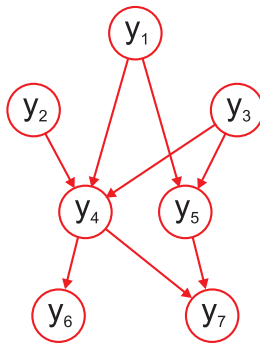
Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети



Совместное распределение системы переменных задается выражением

$$p(Y) = p(y_1, y_2, y_3, y_4, y_5, y_6, y_7) = \\ p(y_1)p(y_2)p(y_3)p(y_4|y_1, y_2, y_3)p(y_5|y_1, y_3)p(y_6|y_4)p(y_7|y_4, y_5).$$

Совместное и условные распределения

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

- В общем случае совместное распределение для ориентированного графа с n вершинами

$$p(Y) = \prod_{i=1}^n p(y_i | \text{pa}_i),$$

где pa_i — множество вершин-родителей y_i

- Обычно предполагается, что атомарные условные распределения $p(y_i | \text{pa}_i)$ известны
- Зная атомарные распределения, мы можем рассчитать (хотя бы теоретически) любые условные вероятности одних подмножеств переменных по другим подмножествам переменных

Вычисление условных распределений I

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

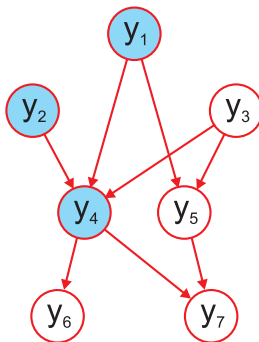
Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

- Вернемся к иллюстрации графической модели из семи переменных
- Пусть нам необходимо найти распределение (y_5, y_7) при заданных значениях y_1, y_2, y_4 и неизвестных y_3, y_6



Вычисление условных распределений II

Лекция 2.
Графические
модели. Общее
представление

Ветров

- По определению условной вероятности

$$p(y_5, y_7 | y_1, y_2, y_4) = \frac{p(y_1, y_2, y_4, y_5, y_7)}{p(y_1, y_2, y_4)}$$

- Расписываем знаменатель

$$p(y_1, y_2, y_4) = p(y_1)p(y_2)p(y_4 | y_1, y_2) = \{Sum\ rule\}$$

$$p(y_1)p(y_2) \int p(y_4 | y_1, y_2, y_3)p(y_3)dy_3$$

- Аналогично числитель

$$p(y_1, y_2, y_4, y_5, y_7) = p(y_1)p(y_2)p(y_4 | y_1, y_2)p(y_5 | y_1)p(y_7 | y_5, y_4) = p(y_1) \times p(y_2) \left(\int p(y_4 | y_1, y_2, y_3)p(y_3)dy_3 \right) \left(\int p(y_5 | y_1, y_3)p(y_3)dy_3 \right) p(y_7 | y_5, y_4)$$

- Для взятия возникающих интегралов обычно пользуются различными аппроксимационными методами
- Таким образом, условное распределение выражено через известные атомарные распределения вида $p(y_i | pa_i)$

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети
Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Особенности использования байесовских сетей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети
Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

- По смыслу построения байесовские сети не могут содержать ориентированные циклы, т.к. это будет нарушать правило умножения вероятностей
- Главным достоинством графических моделей является относительно простое выделение условно-независимых величин, которое облегчает дальнейший анализ, позволяя значительно уменьшить количество факторов, влияющих на данную переменную
- В байесовских сетях сделать это несколько сложнее, чем в марковских

Граф 1

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

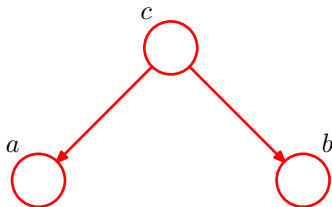
Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети



- Аналогия: Рим (a), император (b) и варвары (c)
- Переменные a и b независимы при заданном c
- Возможна маргинализация (исключение переменной)

$$p(a, b) = \int p(a|c)p(b|c)p(c)dc \neq p(a)p(b)$$

Граф 2

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

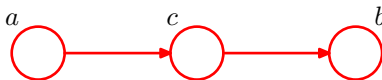
Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети



- Аналогия: хорошая работа (a), премия (c), яхта (b)
- Переменные a и b независимы при заданном c
- Возможна маргинализация (исключение переменной)

$$p(a, b) = p(a) \int p(b|c)p(c|a)dc \neq p(a)p(b)$$

Граф 3

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

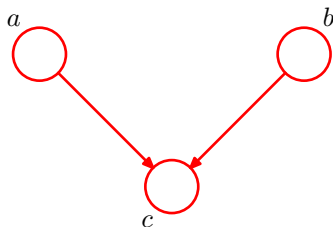
Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети



- Аналогия: вор (a), землетрясение (b) и сигнализация (c)
- Переменные a и b независимы, т.е. $p(a, b) = p(a)p(b)$, но не условно независимы!
- Зависимость $p(c|a, b)$ не может быть выражена через $p(c|a)$ и $p(c|b)$, хотя обратное верно

$$p(c|a) = \int p(c|a, b)p(b)db$$

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клика

Пример использования

Связь с байесовскими сетями

Пример нестрогих вероятностных рассуждений I

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

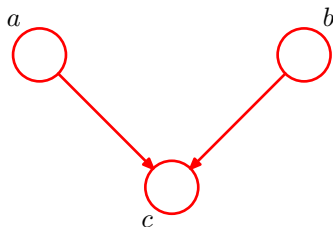
Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети



- Рассмотрим последний граф подробнее. Введем обозначения событий: «сигнализация сработала/не сработала» ($s/\neg s$), «вор есть/вора нет» ($v/\neg v$) и «землетрясение произошло/не произошло» ($z/\neg z$)
- Пусть $p(s|v, \neg z) = p(s|v, z) = 1$, $p(s|\neg v, z) = 0.1$, $p(s|\neg v, \neg z) = 0$, $p(v) = 2 \times 10^{-4}$, $p(z) = 10^{-2}$.
Графическая модель полностью определена

Пример нестрогих вероятностных рассуждений II

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети
Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

Допустим, мы получили сигнал тревоги. Необходимо оценить вероятность того, что в квартире вор $p(v|s)$

$$p(v|s) = \frac{1}{Z} p(s|v)p(v) = \frac{p(s|v)p(v)}{p(s|v)p(v) + p(s|\neg v)p(\neg v)}$$

$$p(s|\neg v) = p(s|\neg v, z)p(z) + p(s|\neg v, \neg z)p(\neg z) = 10^{-3}$$

$$p(s|v) = 1$$

$$p(v|s) \approx \frac{1}{6}, \quad p(\neg v|s) \approx \frac{5}{6}, \quad Z \approx 1.2 \times 10^{-3}$$

Пример нестрогих вероятностных рассуждений III

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

- Пусть теперь дополнительно стало известно, что произошло землетрясение. Как изменится вероятность того, что в квартире вор $p(v|s, z)$?

$$p(v|s, z) = \frac{1}{Z} p(s|v, z) p(v|z), \quad p(v|z) = p(v)$$

$$Z = p(s|v, z) p(v|z) + p(s|\neg v, z) p(\neg v|z) =$$

$$1 \times 2 \times 10^{-4} + 0.1 \times (1 - 2 \times 10^{-4}) = 0.1002$$

$$p(v|s, z) = 0.002, \quad p(\neg v|s, z) = 0.998$$

- Заметим, что события z и v перестали быть независимыми, и добавление сведений о значении z меняет знания о значении v . Это называется эффектом оправдания (explaining away)

Примеры байесовских сетей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Факторизация
байесовских
сетей

Три
элементарных
графа

Пример
использования

Марковские сети

- Скрытые марковские модели
- Фильтр Калмана
- Экспертные системы
- Вероятностный РСА
- Смеси экспертов
- Факторный анализ
- и др.

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования
Связь с
байесовскими
сетями

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клик

Пример использования

Связь с байесовскими сетями

Неориентированные графические модели

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

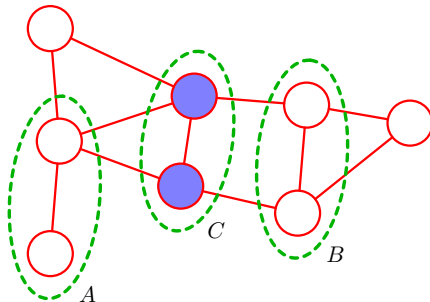
Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования
Связь с
байесовскими
сетями

- При использовании ориентированных графов определение условной независимости не очень просто
- В марковских сетях это проще. На рисунке A и B независимы при условии C
- Ребра графа связывают переменные, между которыми существуют непосредственные (а не опосредованные) зависимости



Факторизация в марковских сетях

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования
Связь с
байесовскими
сетями

- Пусть y_i и y_j независимы при условии, что все остальные переменные нам известны, т.е.
$$p(y_i, y_j | Y_{\{i,j\}}) = p(y_i | Y_{\{i,j\}})p(y_j | Y_{\{i,j\}})$$
- Это означает, что y_i и y_j не соединены ребром (иначе не было бы условной независимости)

- Запишем совместное распределение и применим правило умножения вероятностей

$$p(Y) = p(y_i, y_j | Y_{\{i,j\}})p(Y_{\{i,j\}}) = p(y_i | Y_{\{i,j\}})p(y_j | Y_{\{i,j\}})p(Y_{\{i,j\}})$$

- Таким образом, переменные, не соединенные ребрами, **входят в разные множители** совместного распределения

Потенциалы марковской сети

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

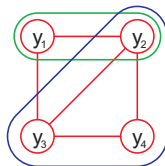
Пример
использования
Связь с
байесовскими
сетями

- В общем виде совместное распределение значений элементов сети записывается с помощью неотрицательных потенциальных функций, определенных на максимальных кликах

$$p(Y) = \frac{1}{Z} \prod_C \psi_C(Y_C), \quad Z = \sum_Y \prod_C \psi_C(Y_C), \quad \psi_C(Y_C) \geq 0$$

- На рисунке синяя клика является максимальной, а зеленая — нет. Совместное распределение имеет вид

$$p(Y) = \frac{1}{Z} \psi_1(y_1, y_2, y_3) \psi_2(y_2, y_3, y_4)$$



Потенциальная функция не обязана иметь вероятностную природу, но чем она больше, тем более вероятны соответствующие значения переменных. Обычно потенциальные функции задаются пользователем исходя из априорных предпочтений тех или иных конфигураций переменных. Реже — настраиваются по данным

Энергетическая нотация

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования
Связь с
байесовскими
сетями

- Иногда удобно ввести обозначение $\psi_C(Y_C) = \exp(-E_C(Y_C))$, где $E_C(Y_C)$ имеет смысл энергии
- Тогда задача нахождения наиболее вероятного состояния системы сводится к задаче минимизации полной энергии системы

$$\arg \max p(Y) = \arg \max \frac{1}{Z} \prod_C \psi_C(Y_C) =$$

$$\arg \max \exp \left(- \sum_C E_C(Y_C) \right) = \arg \min \sum_C E_C(Y_C)$$

- Заметим, что в отличие от байесовских сетей для полного задания графической модели необходимо знать (или конструктивно уметь подсчитывать) нормировочную константу Z

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клик

Пример использования

Связь с байесовскими сетями

Фильтрация изображений

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

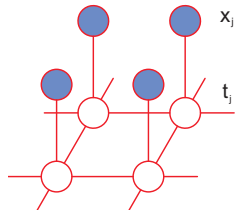
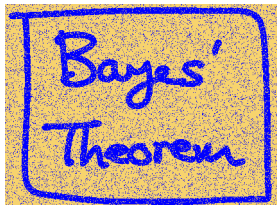
Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями



- Рассмотрим задачу фильтрации изображения. Пусть $x_i \in \{-1, 1\}$ — наблюдаемые пиксели бинарного изображения, а $t_i \in \{-1, 1\}$ — истинные значения пикселей
- Введем энергию системы

$$E(X, T) = h \sum_i t_i - \beta \sum_{i,j} t_i t_j - \eta \sum_i t_i x_i,$$

где $h \in \mathbb{R}$ позволяет отразить априорные предпочтения в пользу того или иного цвета (например, указать, что желтый цвет встречается чаще, чем синий), $\beta > 0$ выражает степень зависимости между соседними пикселями, а $\eta > 0$ показывает интенсивность шума

Разметка областей I

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

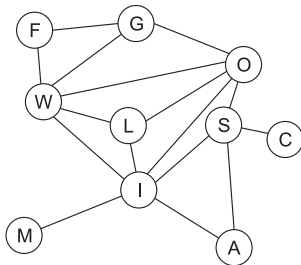
Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

Вернемся к примеру со странами



Совместная плотность задается формулой

$$p(X) = \frac{1}{Z} \psi_1(F, G, W) \psi_2(G, O, W) \psi_3(W, O, L, I) \psi_4(I, S, O) \times \\ \times \psi_5(S, I, A) \psi_6(S, C) \psi_7(M, I)$$

Разметка областей II

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

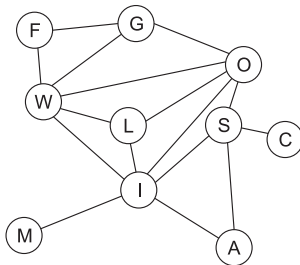
Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями



- Предположим, что переменные могут принимать одно из четырех значений $\{red, yellow, blue, green\}$
- Требование несовпадающих цветов областей эквивалентно условию равенства нулю потенциала, если хотя бы два его аргумента имеют одинаковое значение, например $\psi_5(red, blue, red) = 0$

Разметка областей III

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

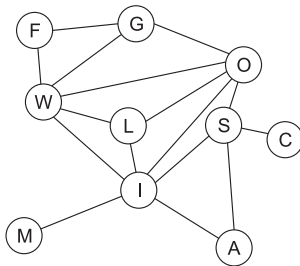
Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями



Мы можем снизить число нежелаемых цветовых переходов (например, из желтого в красный) снизив соответствующие значения потенциалов $\psi_7(\text{red}, \text{yellow})$, $\psi_7(\text{yellow}, \text{red})$, $\psi_6(\text{red}, \text{yellow})$ и $\psi_6(\text{yellow}, \text{red})$

Разметка областей IV

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

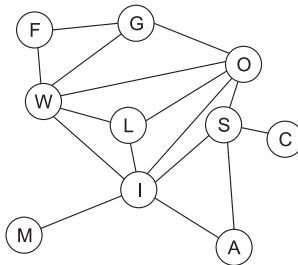
Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями



Мы можем искусственно способствовать окраске отдельных регионов в выбранные цвета, вводя индивидуальные множители, например

$$\psi_1(F, G, W) = \phi_1(F, G, W)\phi_2(F)\phi_3(W).$$

Теперь можно увеличить значение $\phi_2(\text{blue})$ и $\phi_3(\text{green})$, чтобы получить политическую карту, привычную российскому глазу

Примеры марковских сетей

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

- Изображения
- Социальные сети
- Случайные поля
- Карты сайтов
- Условные случайные поля
- и др.

План

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

Ликбез

Формула Байеса

Условная независимость случайных величин

Графические модели

Задачи со структурными ограничениями

Основные проблемы в анализе графических моделей

Байесовские сети

Факторизация байесовских сетей

Три элементарных графа

Пример использования

Марковские сети

Потенциалы и энергия клик

Пример использования

Связь с байесовскими сетями

Марковские vs. Байесовские сети

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик
Пример
использования

Связь с
байесовскими
сетями

Сходства и различия двух типов графических моделей

Свойство	Марковские сети	Байесовские сети
Форма	Произв. потенциалов	Произв. потенциалов
Потенциалы	Произвольные	Усл. вероятности
Циклы	Разрешены	Запрещены
Нормировка	$Z = ?$	$Z = 1$
Условная нез-ть	Легко проверяема	Сложнее
Полнота	нет	нет
Анализ	MCMC, VR, и т.д..	Сводит к марковским

Существующее положение дел

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

- На сегодняшний день существуют эффективные алгоритмы (sum-product, max-product) анализа ациклических графов (деревьев), решающие все три основные задачи анализа графических моделей
- Частным случаем деревьев являются графы-цепочки, характеризующие, например, сигналы во времени
- В случае наличия циклов сложность точных алгоритмов резко возрастает
- Для анализа графов с циклами в основном используются приближенные методы (loopу BP, EP, MCMC)
- В некоторых частных случаях для анализа циклических сетей существуют эффективные точные алгоритмы, например, разрезы графов

Сведение байесовских сетей к марковским

Лекция 2.
Графические
модели. Общее
представление

Ветров

Ликбез

Графические
модели

Байесовские сети

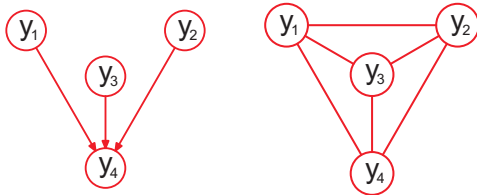
Марковские сети

Потенциалы и
энергия клик

Пример
использования

Связь с
байесовскими
сетями

- Наиболее разработаны в настоящее время методы анализа марковских сетей
- Байесовскую сеть можно легко свести к марковской с потерей информации, переживнив всех родителей (морализация)
- Заметим, что в приведенном примере все полезные свойства оказались потеряны, и мы получили банальную клику, в которой все зависят от всех



Лекция 3. Скрытые марковские модели. Часть 1

Ю. И. Журавлев¹, Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Метод
динамического
программирования

Основы
применения
СММ

EM-алгоритм

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

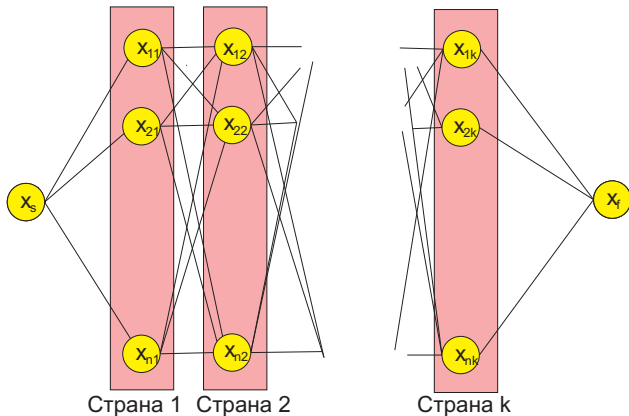
Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

Задача объезда стран



- Рассмотрим такую задачу: необходимо в определенной последовательности объехать k стран, в каждой из которых провести одну ночь в отеле, потратив минимум денег
- Если в каждой стране имеется n городов, то задача сводится к перебору n^k вариантов

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез
Метод
динамического
программирования

Основы
применения
СММ

EM-алгоритм

Функция Беллмана

Лекция 3.
Скрытые
марковские
модели. Часть 1

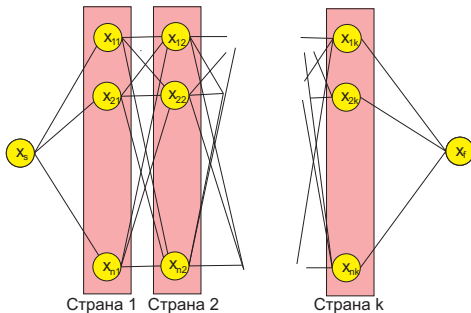
Ветров

Ликбез

Метод
динамического
программирования

Основы
применения
СММ

EM-алгоритм



- Пусть цена билета между городами x_{ij} и $x_{i+1,l}$ задается функцией $f(x_{ij}, x_{i+1,l})$, а цена ночлега в городе — функцией $h(x_{ij})$
- Подсчитаем функцию Беллмана $V(x)$, определяемую рекуррентно: $V(x_s) = 0$

$$V(x_{i+1,l}) = \min_j [V(x_{ij}) + f(x_{ij}, x_{i+1,l}) + h(x_{i+1,l})]$$

- Физический смысл функции Беллмана это наименьшая сумма денег, которую нужно потратить, чтобы добраться из города x_s в город $x_{i+1,l}$ (легко показать по индукции)

Динамическое программирование

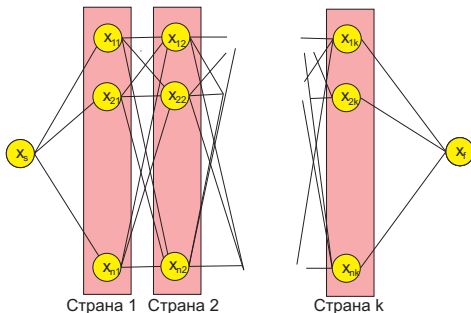
Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез
Метод
динамического
программирования

Основы
применения
СММ

EM-алгоритм



- Определим также функцию $S(x)$, возвращающую город, откуда мы приехали в город x : $S(x_s) = \emptyset$

$$S(x_{i+1,l}) = \arg \min_j [V(x_{ij}) + f(x_{ij}, x_{i+1,l}) + h(x_{i+1,l})]$$

- Тогда оптимальный путь $(x_1^*, x_2^*, \dots, x_k^*)$ может быть получен путем рекуррентного вызова функции $S(x)$: $x_n^* = S(x_f)$

$$x_m^* = S(x_{m+1}^*)$$

Особенности динамического программирования

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

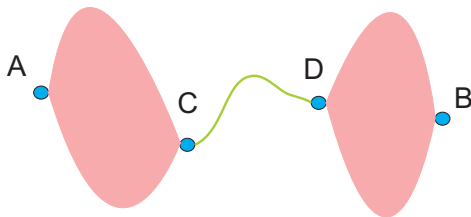
Ликбез
Метод
динамического
программирования

Основы
применения
СММ

EM-алгоритм

Метод динамического программирования эффективен, если имеет место

- Перекрывающиеся подзадачи, которые необходимо решить, чтобы решить исходную задачу
- Оптимальная подструктура (принцип кусочной оптимальности)
- Возможность запомнить решения подзадач



Если известно, что оптимальный путь проходит через точки C и D и известен оптимальный путь между ними, то этот путь станет частью оптимального пути между A и B

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

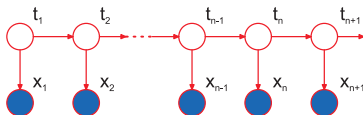
Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

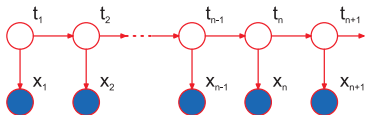
Скрытая Марковская модель (СММ)



Скрытая Марковская модель [первого порядка] — это вероятностная модель последовательности, которая

- Состоит из набора наблюдаемых переменных $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^d$ и латентных (скрытых) переменных $T = \{t_1, \dots, t_N\}$, $t_n \in \{0, 1\}^K$, $\sum_{j=1}^K t_{nj} = 1$
- Латентные переменные T являются **бинарными** и кодируют K состояний, поэтому их иногда называют переменными состояния
- Значение наблюдаемого вектора \mathbf{x}_n , взятого в момент времени n , зависит только от скрытого состояния t_n , которое в свою очередь зависит только от скрытого состояния в предыдущий момент времени t_{n-1}

Примеры использования СММ



Что можно анализировать с помощью СММ

- Речь
- Видео
- Поведение
- Фондовые рынки
- Естественный язык
- ДНК
- и др.

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

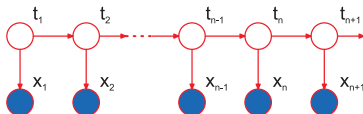
Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

Скрытая Марковская модель (СММ)



- Скрытая марковская модель является частным случаем байесовской сети (графической модели, задаваемой ориентированным графом)
- Граф, задающий СММ, является ациклическим, поэтому для СММ существуют эффективные алгоритмы вывода
- Для полного задания модели достаточно задать все условные распределения вида $p(x_n | t_n)$, $p(t_n | t_{n-1})$ и априорное распределение $p(t_1)$

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем
Алгоритм
Витерби

EM-алгоритм

Спецификация вероятностной модели

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем
Алгоритм
Витерби

EM-алгоритм

Пусть имеется K возможных состояний. Закодируем состояние в каждый момент времени n бинарным вектором $\mathbf{t}_n = (t_{n1}, \dots, t_{nK})$, где

$$t_{nj} = \begin{cases} 1, & \text{если в момент } n \text{ модель находится в состоянии } j \\ 0, & \text{иначе} \end{cases}$$

Тогда распределение $p(\mathbf{t}_n | \mathbf{t}_{n-1})$ можно задать матрицей перехода A размера $K \times K$, где $A_{ij} = p(t_{nj} = 1 | t_{n-1,i} = 1)$, $\sum_j A_{ij} = 1$, т.е.

$$p(\mathbf{t}_n | \mathbf{t}_{n-1}) = \prod_{i=1}^K \prod_{j=1}^K A_{ij}^{t_{n-1,i} t_{nj}}$$

Пусть в первый момент времени $p(t_{1j} = 1) = \pi_j$. Тогда

$$p(\mathbf{t}_1) = \prod_{j=1}^K \pi_j^{t_{1j}}$$

Замечание

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

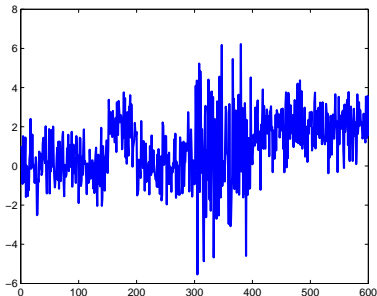
Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Хотя матрица A может быть произвольного вида с учетом ограничений на неотрицательность и сумму элементов строки, с точки зрения СММ представляет интерес диагональное преобладание матрицы перехода
- В этом случае можно ожидать, что процесс находится в некотором состоянии на протяжении какого-то отрезка времени
- Появляется простая физическая интерпретация СММ: имеется процесс, который иногда (относительно редко) скачкообразно меняет свои характеристики



Спецификация вероятностной модели

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Условное распределение $p(\mathbf{x}_n | \mathbf{t}_n)$ определяется текущим состоянием \mathbf{t}_n
- Обычно предполагают, что оно нам известно с точностью до параметров ϕ_k , $k \in \{1, \dots, K\}$, т.е. если $t_{n1} = 1$, то \mathbf{x}_n взят из распределения $p(\mathbf{x}_n | \phi_1)$, если $t_{n2} = 1$, то \mathbf{x}_n взят из распределения $p(\mathbf{x}_n | \phi_2)$, и т.д.
- Таким образом

$$p(\mathbf{x}_n | \mathbf{t}_n) = \prod_{j=1}^K (p(\mathbf{x}_n | \phi_j))^{t_{nj}}$$

Задачи в СММ

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем
Алгоритм
Витерби

EM-алгоритм

Обозначим полный набор параметров $\Theta = \{\pi, A, \phi\}$. Тогда основные задачи, возникающие в СММ, можно сформулировать следующим образом:

- **Обучение с учителем.** Известна некоторая последовательность X , для которой заданы T . Задача состоит в оценке по обучающей выборке набора параметров Θ .
- **Сегментация.** Известна некоторая последовательность X и набор параметров Θ . Задача состоит в получении наиболее правдоподобной последовательности состояний T как $\arg \max_T p(T|X, \Theta)$ (алгоритм Витерби).
- **Обучение без учителя.** Известна некоторая последовательность X и число состояний K . Задача состоит в оценке параметров Θ (EM-алгоритм).
 - **Нахождение маргинального распределения** $p(t_n|X, \Theta)$ компоненты t_n по заданным X и Θ
- **Прогнозирование.** Известна некоторая последовательность X . Задача состоит в оценке наблюдаемого вектора в следующий момент времени $N + 1$ — $p(x_{N+1}|X)$.

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

Совместное распределение переменных СММ

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Предположим, нам задана **обучающая выборка** (X, T) , представляющая собой одну или несколько последовательностей, в которых известны значения скрытых компонент
- Требуется оценить вектор параметров Θ
- По построению байесовской сети совместное распределение переменных задается формулой

$$\begin{aligned} p(X, T | \Theta) &= p_{\pi}(t_1) \prod_{n=1}^N p_{\phi}(\mathbf{x}_n | t_n) \prod_{n=2}^N p_A(t_n | t_{n-1}) = \\ &= \prod_{j=1}^K \pi_j^{t_{1j}} \left(\prod_{n=2}^N \prod_{i=1}^K \prod_{j=1}^K A_{ij}^{t_{n-1}, i t_{nj}} \right) \left(\prod_{n=1}^N \prod_{k=1}^K (p(\mathbf{x}_n | \phi_k))^{t_{nk}} \right) \end{aligned}$$

Метод максимального правдоподобия

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Для оценки параметров Θ воспользуемся методом максимального правдоподобия

$$\Theta_{ML} = \arg \max p(X, T | \Theta) = \arg \max \log p(X, T | \Theta)$$

- Для удобства перейдем к логарифму функции правдоподобия (положение максимума, очевидно, не изменится)

$$\begin{aligned} \log p(X, T | \Theta) = & \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) + \\ & + \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(\mathbf{x}_n | \phi_k) \right) \end{aligned}$$

Функция Лагранжа

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Параметры, входящие в Θ , не могут принимать произвольные значения, следовательно необходима оптимизация при ограничениях

$$\sum_{j=1}^K \pi_j = 1, \quad \sum_{j=1}^K A_{ij} = 1, \quad \forall i = \overline{1, K}$$

- Воспользуемся правилом множителей Лагранжа и выпишем лагранжиан

$$\mathcal{L}(\Theta, \lambda, \boldsymbol{\mu}) = \log p(X, T | \Theta) + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) + \sum_{i=1}^K \mu_i \left(\sum_{j=1}^K A_{ij} - 1 \right) \rightarrow \text{extr}$$

Оценка π

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

$$\begin{aligned}\log p(X, T|\Theta) &= \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) + \\ &+ \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(x_n | \phi_k) \right) \\ \mathcal{L}(\Theta, \lambda, \boldsymbol{\mu}) &= \log p(X, T|\Theta) + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) + \sum_{i=1}^K \mu_i \left(\sum_{j=1}^K A_{ij} - 1 \right) \rightarrow \text{extr} \\ \frac{\partial \mathcal{L}(\Theta, \lambda, \boldsymbol{\mu})}{\partial \pi_j} &= \frac{t_{1j}}{\pi_j} + \lambda = 0 \Rightarrow \pi_j = -\frac{t_{1j}}{\lambda} \\ \sum_{j=1}^K \pi_j &= 1 \Rightarrow \lambda = -\sum_{j=1}^K t_{1j} = -1 \\ \pi_j &= t_{1j}\end{aligned}$$

Оценка матрицы A

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

$$\log p(X, T|\Theta) = \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) +$$
$$+ \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(\mathbf{x}_n | \phi_k) \right)$$
$$\mathcal{L}(\Theta, \lambda, \boldsymbol{\mu}) = \log p(X, T|\Theta) + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) + \sum_{i=1}^K \mu_i \left(\sum_{j=1}^K A_{ij} - 1 \right) \rightarrow \text{extr}$$

$$\frac{\partial \mathcal{L}(\Theta, \lambda, \boldsymbol{\mu})}{\partial A_{ij}} = \sum_{n=2}^N \frac{t_{n-1,i} t_{nj}}{A_{ij}} + \mu_i = 0 \Rightarrow A_{ij} = - \sum_{n=2}^N \frac{t_{n-1,i} t_{nj}}{\mu_i}$$

$$\sum_{j=1}^K A_{ij} = 1 \Rightarrow \mu_i = - \sum_{n=2}^N \sum_{j=1}^K t_{n-1,i} t_{nj} = - \sum_{n=2}^N t_{n-1,i}$$

$$A_{ij} = \frac{\sum_{n=2}^N t_{n-1,i} t_{nj}}{\sum_{n=2}^N t_{n-1,i}}$$

Оценка условной плотности $p(\mathbf{x}|\mathbf{t})$

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

$$\begin{aligned}\log p(X, T|\Theta) &= \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) + \\ &+ \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(\mathbf{x}_n|\phi_k) \right) \\ \mathcal{L}(\Theta, \lambda, \boldsymbol{\mu}) &= \log p(X, T|\Theta) + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) + \sum_{i=1}^K \mu_i \left(\sum_{j=1}^K A_{ij} - 1 \right) \rightarrow \text{extr} \\ \frac{\partial \mathcal{L}(\Theta, \lambda, \boldsymbol{\mu})}{\partial \phi_k} &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{\partial \log p(\mathbf{x}_n|\phi_k)}{\partial \phi_k} = \sum_{\{n: t_{nk}=1\}} \frac{\partial \log p(\mathbf{x}_n|\phi_k)}{\partial \phi_k} = 0\end{aligned}$$

Получили стандартную задачу максимизации правдоподобия по выборке независимых одинаково-распределенных объектов

$$\phi_k = \arg \max \sum_{\{n: t_{nk}=1\}} \log p(\mathbf{x}_n|\phi_k)$$

Для оценки параметров ϕ_k можно воспользоваться методами восстановления плотностей, например, EM-алгоритмом

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

Сегментация сигнала

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- Пусть известна некоторая последовательность наблюдений X и набор параметров СММ Θ . Требуется определить наиболее вероятную последовательность состояний T , т.е. найти $\arg \max_T p(T|X, \Theta)$
- Заметим, что $p(X|\Theta)$ не зависит от T , поэтому

$$\begin{aligned}\arg \max_T p(T|X, \Theta) &= \arg \max_T \frac{p(X, T|\Theta)}{p(X|\Theta)} = \\ &= \arg \max_T p(X, T|\Theta) = \arg \max_T \log p(X, T|\Theta)\end{aligned}$$

- Но это же классическая задача динамического программирования!

Аналогия с задачей объезда стран

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

- В самом деле логарифм совместной плотности по определению

$$\begin{aligned} \log p(X, T | \Theta) = & \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) + \\ & + \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(\mathbf{x}_n | \phi_k) \right) \end{aligned}$$

- Первое слагаемое определяет «пункт отбытия», второе слагаемое — стоимость переезда из города страны $n - 1$ в город в стране n , а третье слагаемое отражает «стоимость ночлега» в выбранном городе страны n

Алгоритм Витерби

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

$$\log p(X, T | \Theta) = \left(\sum_{j=1}^K t_{1j} \log \pi_j \right) + \\ + \left(\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} \right) + \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p(\mathbf{x}_n | \phi_k) \right)$$

Функция Беллмана $V_{1j} = \log \pi_j$

$$V_{nj} = \max_i [V_{n-1,i} + \log A_{ij} + \log p(\mathbf{x}_n | \phi_j)]$$

Функция S_{nj} определяется аналогично: $S_{1j} = \emptyset$

$$S_{nj} = \arg \max_i [V_{n-1,i} + \log A_{ij} + \log p(\mathbf{x}_n | \phi_j)]$$

Алгоритм Витерби

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

$$V_{nj} = \max_i [V_{n-1,i} + \log A_{ij} + \log p(\mathbf{x}_n | \phi_j)]$$

$$S_{nj} = \arg \max_i [V_{n-1,i} + \log A_{ij} + \log p(\mathbf{x}_n | \phi_j)]$$

Выполнив прямой проход по сигналу, мы оцениваем V_{nj} и S_{nj} , а выполнив обратный проход, мы получаем оптимальные номера оптимальных состояний ($i^*(1), \dots, i^*(N)$): $i^*(N) = \arg \max_i V_{Ni}$

$$i^*(n) = S_{n+1, i^*(n+1)}$$

Легко видеть, что значения переменных t_n определяются так:
 $t_{n, i^*(n)} = 1, t_{ni} = 0, \forall i \neq i^*(n)$

- Алгоритм Витерби позволяет быстро проводить сегментацию очень длинных сигналов
- Существует версия алгоритма Витерби, позволяющая осуществлять сегментацию в реальном времени (точнее, с небольшой задержкой)

Пример использования

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

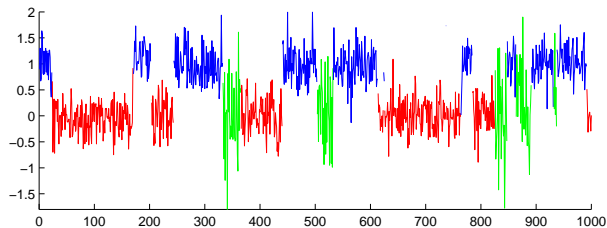
Определение
СММ

Обучение СММ
с учителем

Алгоритм
Витерби

EM-алгоритм

Пример разметки сигнала на три состояния с помощью алгоритма Витерби



План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

Алгоритм Витерби

EM-алгоритм

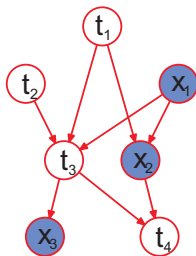
Графические модели с неполными данными

Разделение гауссовской смеси

Максимум неполного правдоподобия

- Предположим имеется графическая модель, в которой известна только часть значений переменных
- Атомарные распределения известны с точностью до вектора параметров θ
- Требуется оценить параметры по наблюдаемым величинам с помощью метода максимального правдоподобия, т.е. найти

$$\theta_{ML} = \arg \max p(X|\theta)$$



Трудности оптимизации неполного правдоподобия

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- По правилу суммирования вероятностей неполное правдоподобие может быть получено в виде суммирования по скрытым переменным полного правдоподобия

$$p(X|\theta) = \sum_T p(X, T|\theta)$$

- Во многих случаях (в частности, в байесовских сетях) подсчет полного правдоподобия тривиален
- При оптимизации правдоподобия удобно переходить к логарифму, в частности выше мы получили явные формулы для $\arg \max_{\theta} p(X, T|\theta) = \arg \max_{\theta} \log p(X, T|\theta)$ в СММ
- Прямая оптимизация логарифма неполного правдоподобия очень затруднительна даже в итерационной форме, т.к. функционал имеет вид «логарифм суммы», в то время как удобно оптимизировать «сумму логарифмов»

$$\log p(X|\theta) = \log \sum_T p(X, T|\theta)$$

Схема EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- На входе: выборка X , зависящая от набора параметров θ
- Инициализируем θ некоторыми начальным приближением
- **E-шаг:** Оцениваем распределение скрытой компоненты при фиксированном значении параметров θ_{old}

$$p(T|X, \theta_{old}) = \frac{p(X, T|\theta_{old})}{\sum_T p(X, T|\theta_{old})}$$

- **M-шаг:** Оптимизируем

$$\mathbb{E}_{T|X, \theta_{old}} \log p(X, T|\theta) = \sum_T p(T|X, \theta_{old}) \log p(X, T|\theta) \rightarrow \max_{\theta}$$

Если бы мы **точно знали значение** $T = T_0$, то вместо мат. ожидания по всевозможным (с учетом наблюдаемых данных) $T|X, \theta_{old}$, мы бы оптимизировали $\log p(X, T_0|\theta)$

- Переход к E-шагу, пока процесс не сойдется

Замечания

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- Оптимизация проводится итерационно методом покоординатного спуска: на каждой итерации последовательно уточняются возможные значения T (E-шаг), а потом пересчитываются значения θ (M-шаг)
- Во многих случаях на M-шаге можно получить явные формулы, т.к. там происходит оптимизация выпуклой комбинации логарифмов полных правдоподобий

$$\sum_T p(T|X, \theta_{old}) \log p(X, T|\theta) \rightarrow \max_{\theta},$$

имеющей вид взвешенной «суммы логарифмов»

План

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Ликбез

Метод динамического программирования

Основы применения СММ

Определение СММ

Обучение СММ с учителем

Алгоритм Витерби

EM-алгоритм

Графические модели с неполными данными

Разделение гауссовской смеси

Смесь гауссовских распределений

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

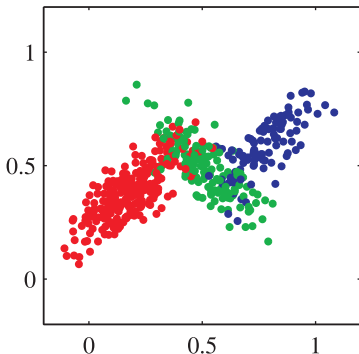
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- Имеется выборка $X \sim \sum_{j=1}^l w_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \in \mathbb{R}^d$, $w_j \geq 0$,
 $\sum_{j=1}^l w_j = 1$
- Требуется восстановить плотность генеральной совокупности



EM-алгоритм

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- Выбираем начальное приближение μ_j, w_j, Σ_j
- E-шаг: Вычисляем распределение скрытых переменных $z_i \in \{0, 1\}$, $\sum_j z_{ij} = 1$, которые определяют, к какой компоненте смеси принадлежит объект \mathbf{x}_i

$$\gamma(z_{ij}) = \frac{w_j \mathcal{N}(\mathbf{x}_i | \mu_j, \Sigma_j)}{\sum_{k=1}^l w_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)}$$

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

$$\mu_j^{new} = \frac{1}{N_j} \sum_{i=1}^n \gamma(z_{ij}) \mathbf{x}_i \quad w_j^{new} = \frac{N_j}{n} \quad N_j = \sum_{i=1}^n \gamma(z_{ij})$$

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

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

Вывод формул на M-шаге

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Выражение для мат. ожидания логарифма правдоподобия по апостериорному распределению имеет вид

$$\mathbb{E}_{Z|X, \theta} \log p(X, Z | \theta) = \sum_{i=1}^n \sum_{j=1}^l \gamma(z_{ij}) (\log w_j + \log \mathcal{N}(x_i | \mu_j, \Sigma_j))$$

Дифференцирование по μ_j и Σ_j (точнее по Σ_j^{-1}) выполняется аналогично случаю одной многомерной гауссианы, только объекты теперь берутся с весом $\gamma(z_{ij})$.

Оптимизация по w_j проводится с учетом ограничения $\sum_{j=1}^l w_j = 1$ по правилу множителей Лагранжа

$$\sum_{i=1}^n \sum_{j=1}^l \gamma(z_{ij}) \log w_j + \lambda \left(\sum_{j=1}^l w_j - 1 \right)$$

Дифференцируя лагранжиан по w_j получаем уравнения

$$\frac{\sum_{i=1}^n \gamma(z_{ij})}{w_j} + \lambda = 0, \quad w_j = -\frac{\sum_{i=1}^n \gamma(z_{ij})}{\lambda}$$

Учитывая, что $\sum_{j=1}^l w_j = 1$ окончательно получаем

$$\lambda = -n, \quad w_j = \frac{\sum_{i=1}^n \gamma(z_{ij})}{n} = \frac{N_j}{n}$$

Пример работы EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

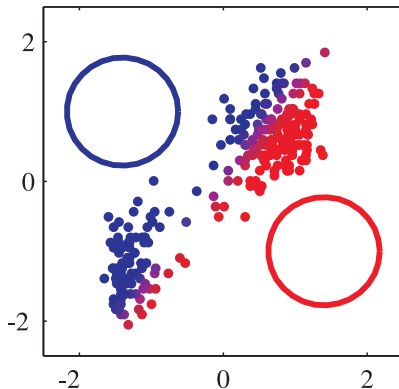
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Применение EM-алгоритма для разделения смеси двух гауссиан. Итерация 0.



Пример работы EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

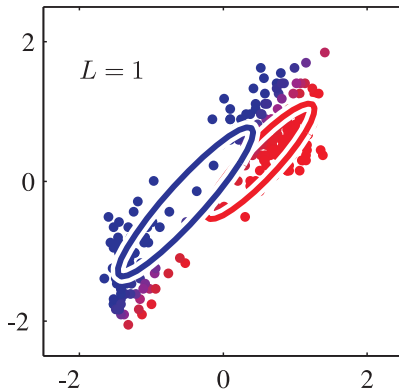
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Применение EM-алгоритма для разделения смеси двух гауссиан. Итерация 1.



Пример работы EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

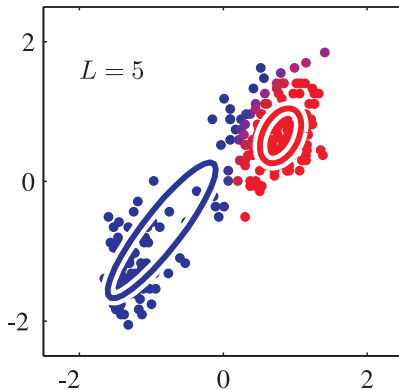
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Применение EM-алгоритма для разделения смеси двух гауссиан. Итерация 5.



Пример работы EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

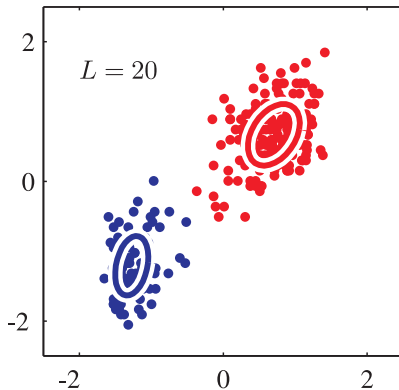
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Применение EM-алгоритма для разделения смеси двух гауссиан. Итерация 20.



Недостатки EM-алгоритма

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

- 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!!**

Применение EM-алгоритма для описания сложных плотностей

Лекция 3.
Скрытые
марковские
модели. Часть 1

Ветров

Ликбез

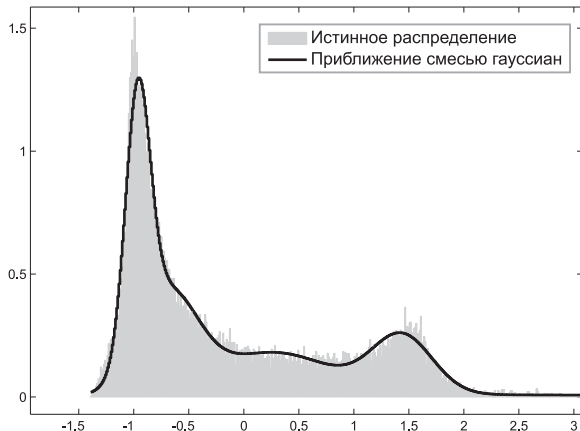
Основы
применения
СММ

EM-алгоритм

Графические
модели с
неполными
данными

Разделение
гауссовской
смеси

Смесью гауссиан можно эффективно приближать сложные, плохо-параметризуемые плотности. При этом получаются аналитически заданные приближения



Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Лекция 4. Скрытые марковские модели. Часть 2.

Ю. И. Журавлев¹, Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

Скрытая Марковская модель (СММ)

Лекция 4.
Скрытые
марковские
модели. Часть 2.

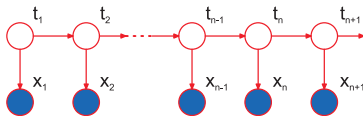
Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример



Скрытая Марковская модель [первого порядка] — это вероятностная модель последовательности, которая

- Состоит из набора наблюдаемых переменных $X = \{x_1, \dots, x_N\}$ и латентных (скрытых) переменных $T = \{t_1, \dots, t_N\}$.
- Латентные переменные T являются **дискретными**, поэтому их иногда называют переменными состояния.
- Значение наблюдаемого вектора в момент времени n x_n зависит только от скрытого состояния t_n , которое в свою очередь зависит только от скрытого состояния в предыдущий момент времени t_{n-1} .

Спецификация вероятностной модели I

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Пусть имеется K возможных состояний. Закодируем состояние в каждый момент времени n бинарным вектором $\mathbf{t}_n = (t_{n1}, \dots, t_{nK})$, где

$$t_{nj} = \begin{cases} 1, & \text{если в момент } n \text{ модель находится в состоянии } j \\ 0, & \text{иначе} \end{cases}$$

Тогда распределение $p(\mathbf{t}_n | \mathbf{t}_{n-1})$ можно задать матрицей перехода A размера $K \times K$, где $A_{ij} = p(t_{nj} = 1 | t_{n-1,i} = 1)$, $\sum_j A_{ij} = 1$, т.е.

$$p(\mathbf{t}_n | \mathbf{t}_{n-1}) = \prod_{i=1}^K \prod_{j=1}^K A_{ij}^{t_{n-1,i} t_{nj}}$$

Пусть в первый момент времени $p(t_{1j} = 1) = \pi_j$. Тогда

$$p(\mathbf{t}_1) = \prod_{j=1}^K \pi_j^{t_{1j}}$$

Спецификация вероятностной модели II

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Пусть для каждого состояния $k \in \{1, \dots, K\}$ в момент времени n известна модель генерации наблюдаемых данных $\mathbf{x}_n p(\mathbf{x}_n | \phi_k)$, задаваемая с помощью набора параметров ϕ_k . Тогда

$$p(\mathbf{x}_n | \mathbf{t}_n) = \prod_{j=1}^K (p(\mathbf{x}_n | \phi_k))^{t_{nk}}$$

Обозначим полный набор параметров СММ через $\Theta = \{\pi, A, \phi\}$. Тогда правдоподобие в СММ вычисляется как

$$\begin{aligned} p(X, T | \Theta) &= p(\mathbf{t}_1, \pi) \prod_{n=2}^N p(\mathbf{t}_n | \mathbf{t}_{n-1}, A) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{t}_n, \phi) = \\ &= \prod_{j=1}^K \pi_j^{t_{1j}} \left(\prod_{n=2}^N \prod_{i=1}^K \prod_{j=1}^K A_{ij}^{t_{n-1,i} t_{nj}} \right) \left(\prod_{n=1}^N \prod_{k=1}^K (p(\mathbf{x}_n | \phi_k))^{t_{nk}} \right) \end{aligned}$$

Задачи в СММ

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

- **Распознавание** (предыдущая лекция). Известна некоторая последовательность X и набор параметров Θ . Задача состоит в получении наиболее правдоподобной последовательности состояний T как $\arg \max_T p(T|X, \Theta)$ (алгоритм Витерби).
- **Обучение с учителем** (предыдущая лекция). Известна некоторая последовательность X , для которой заданы T . Задача состоит в оценке по обучающей выборке набора параметров Θ .
- **Обучение без учителя**. Известна некоторая последовательность X и число состояний K . Задача состоит в оценке параметров Θ (EM-алгоритм).
 - **Нахождение маргинального распределения** $p(t_n|X, \Theta)$ компоненты t_n по заданным X и Θ
- **Прогнозирование**. Известна некоторая последовательность X . Задача состоит в оценке наблюдаемого вектора в следующий момент времени $N + 1 - p(\mathbf{x}_{N+1}|X)$.

План лекции

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Метод максимального правдоподобия для СММ

Алгоритм «вперед-назад»

Устойчивые формулы для алгоритма «вперед-назад»

Модельный пример

План лекции

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Метод максимального правдоподобия для СММ

Алгоритм «вперед-назад»

Устойчивые формулы для алгоритма «вперед-назад»

Модельный пример

Метод максимального правдоподобия

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Для оценки параметров СММ Θ воспользуемся методом максимального правдоподобия:

$$\Theta_* = \arg \max_{\Theta} p(X|\Theta) = \arg \max_{\Theta} \sum_T p(X, T|\Theta)$$

Прямая максимизация правдоподобия затруднительна, т.к. оптимизируемая функция не является выпуклой и, кроме того, для вычисления функции требуется суммирование N^K слагаемых. Можно воспользоваться итерационным EM-алгоритмом.

EM-алгоритм в общем виде

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Требуется найти максимум правдоподобия в вероятностной модели со скрытыми переменными:

$$p(X|\Theta) = \sum_T p(X, T|\Theta) \rightarrow \max_{\Theta} \Leftrightarrow \log \left(\sum_T p(X, T|\Theta) \right) \rightarrow \max_{\Theta}$$

- **Е-шаг.** Фиксируется значение параметров Θ_{old} . Оценивается апостериорное распределение на скрытые переменные $p(T|X, \Theta_{old})$, и полное правдоподобие усредняется по полученному распределению:

$$\mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta) = \sum_T \log p(X, T|\Theta) p(T|X, \Theta_{old})$$

- **М-шаг.** Фиксируется апостериорное распределение $p(T|X, \Theta_{old})$, и производится поиск новых значений параметров Θ_{new} :

$$\Theta_{new} = \arg \max_{\Theta} \mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta)$$

- Шаги Е и М повторяются до сходимости.

E-шаг

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Обозначим

$$\gamma(t_{nj}) = \mathbb{E}_{T|X, \Theta} t_{nj} = p(t_{nj} = 1 | X, \Theta_{old}),$$

$$\xi(t_{n-1,i}, t_{nj}) = \mathbb{E}_{T|X, \Theta}(t_{n-1,i} t_{nj}) = p(t_{n-1,i} = 1, t_{nj} = 1 | X, \Theta_{old}).$$

Тогда

$$\begin{aligned} \log p(X, T | \Theta) &= \sum_{j=1}^K t_{1j} \log \pi_j + \\ &\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K t_{n-1,i} t_{nj} \log A_{ij} + \sum_{n=1}^N \sum_{j=1}^K t_{nj} \log p(\mathbf{x}_n | \phi_k) \end{aligned}$$

$$\mathbb{E}_{T|X, \Theta_{old}} \log p(X, T | \Theta) = \sum_{j=1}^K \gamma(t_{1j}) \log \pi_j +$$

$$\sum_{n=2}^N \sum_{i=1}^K \sum_{j=1}^K \xi(t_{n-1,i}, t_{nj}) \log A_{ij} + \sum_{n=1}^N \sum_{j=1}^K \gamma(t_{nj}) \log p(\mathbf{x}_n | \phi_k)$$

M-шаг

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

$$\pi : \sum_{j=1}^K \gamma(t_{1j}) \log \pi_j + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) \rightarrow \text{extr}_{\pi, \lambda}$$

$$\frac{\gamma(t_{1j})}{\pi_j} + \lambda = 0 \Rightarrow \pi_j = -\frac{\gamma(t_{1j})}{\lambda}$$

$$\sum_{j=1}^K \pi_j = 1 \Rightarrow \lambda = -\sum_{i=1}^K \gamma(t_{1i})$$

$$\pi_j^{\text{new}} = \frac{\gamma(t_{1j})}{\sum_{i=1}^K \gamma(t_{1i})}$$

Действуя аналогично для A , принимая во внимание, что $\sum_{j=1}^K A_{ij} = 1 \forall i$, получаем:

$$A_{ij}^{\text{new}} = \frac{\sum_{n=2}^N \xi(t_{n-1, i} t_{nj})}{\sum_{k=1}^K \sum_{n=2}^N \xi(t_{n-1, i} t_{nk})}$$

M-шаг для компонент ϕ

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

M-шаг для компонент генерации данных $p(\mathbf{x}_n|\phi_k)$ абсолютно аналогичен M-шагу для оценки параметров при восстановлении смесей распределений. В частности, если в качестве компонент выступают многомерные нормальные распределения

$$p(\mathbf{x}_n|\phi_k) = \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

то задача оптимизации для параметров $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ может быть решена в явном виде:

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(t_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(t_{nk})}$$
$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma(t_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(t_{nk})}$$

Инициализация параметров Θ

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Для начала работы EM-алгоритма необходимо задать начальные значения параметров $\Theta = (\boldsymbol{\pi}, A, \phi)$. Заметим, что если какой-нибудь параметр инициализирован нулем, то в процессе итераций EM его значение не изменится.

- Значения параметров $\boldsymbol{\pi}$ и A обычно выбираются случайными при соблюдении ограничений $\sum_j \pi_j = 1$ и $\sum_j A_{ij} = 1 \forall i$.
- Инициализация ϕ зависит от формы распределений $p(\mathbf{x}|\phi)$. В случае нормальных распределений можно провести кластеризацию данных на K кластеров и выбрать в качестве $\boldsymbol{\mu}_k$ и Σ_k центр и разброс соответствующего кластера.

План лекции

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Метод максимального правдоподобия для СММ

Алгоритм «вперед-назад»

Устойчивые формулы для алгоритма «вперед-назад»

Модельный пример

Вычисление апостериорных распределений на скрытые компоненты

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

На E-шаге алгоритма обучения СММ требуется вычисление апостериорных распределений на скрытые компоненты

$$\gamma(t_{nj}) = p(t_{nj} = 1 | X, \Theta), \quad \xi(t_{n-1,i}, t_{nj}) = p(t_{n-1,i} = 1, t_{nj} = 1 | X, \Theta)$$

Алгоритм «вперед-назад» (Баума-Уэлша) позволяет эффективно вычислять эти величины для всех n, i, j за линейное по N время.

В дальнейшем для удобства будем опускать Θ во всех формулах, считая набор параметров фиксированным.

Свойства условной независимости для СММ

Лекция 4.
Скрытые
марковские
модели. Часть 2.

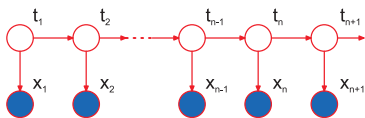
Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример



$$p(X|t_n) = p(x_1, \dots, x_n | t_n) p(x_{n+1}, \dots, x_N | t_n)$$

$$p(x_1, \dots, x_{n-1} | x_n, t_n) = p(x_1, \dots, x_{n-1} | t_n)$$

$$p(x_1, \dots, x_{n-1} | t_{n-1}, t_n) = p(x_1, \dots, x_{n-1} | t_{n-1})$$

$$p(x_{n+1}, \dots, x_N | t_n, t_{n+1}) = p(x_{n+1}, \dots, x_N | t_{n+1})$$

$$p(x_{n+2}, \dots, x_N | t_{n+1}, x_{n+1}) = p(x_{n+2}, \dots, x_N | t_{n+1})$$

$$p(X | t_{n-1}, t_n) = p(x_1, \dots, x_{n-1} | t_{n-1}) p(x_n | t_n) \times \\ \times p(x_{n+1}, \dots, x_N | t_n)$$

$$p(x_{N+1} | X, t_{N+1}) = p(x_{N+1} | t_{N+1})$$

$$p(t_{N+1} | t_N, X) = p(t_{N+1} | t_N)$$

Вычисление $\gamma(t_{nj})$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

По формуле Байеса

$$\gamma(t_n) = p(t_n|X) = \frac{p(X|t_n)p(t_n)}{p(X)}$$

Пользуясь свойствами условной независимости для СММ, получаем

$$\gamma(t_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, t_n)p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N|t_n)}{p(X)} = \frac{\alpha(t_n)\beta(t_n)}{p(X)}$$

Здесь

$$\alpha(t_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, t_n)$$

$$\beta(t_n) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N|t_n)$$

Алгоритм «вперед-назад» позволяет быстро рекуррентно вычислять значение $\alpha(t_n)$ через $\alpha(t_{n-1})$ (проход вперед) и $\beta(t_n)$ через $\beta(t_{n+1})$ (проход назад).

Рекуррентная формула для $\alpha(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

$$\begin{aligned}\alpha(t_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, t_n) = \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | t_n) p(t_n) = \\ &= p(\mathbf{x}_n | t_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | t_n) p(t_n) = \\ &= p(\mathbf{x}_n | t_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, t_n) = \\ &= p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, t_{n-1}, t_n) = \\ &= p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, t_n | t_{n-1}) p(t_{n-1}) = \\ &= p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | t_{n-1}) p(t_n | t_{n-1}) p(t_{n-1}) = \\ &= p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, t_{n-1}) p(t_n | t_{n-1}) = \\ &= p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} \alpha(t_{n-1}) p(t_n | t_{n-1})\end{aligned}$$

Вычисление $\alpha(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Для запуска рекуррентного процесса необходимо
вычислить

$$\alpha(t_1) = p(x_1, t_1) = p(t_1)p(x_1|t_1) = \prod_{j=1}^K (\pi_j p(x_1|\phi_j))^{t_{1j}}$$

На каждом шаге рекурсии вектор $\alpha(t_n)$ длины K умножается на матрицу $p(t_n|t_{n-1})$ размера $K \times K$. Поэтому сложность всего рекуррентного процесса составляет $O(NK^2)$.

Рекуррентная формула для $\beta(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

$$\begin{aligned}\beta(t_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | t_n) = \\ &= \sum_{t_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, t_{n+1} | t_n) = \\ &= \sum_{t_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | t_{n+1}) p(t_{n+1} | t_n) = \\ &= \sum_{t_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | t_{n+1}) p(\mathbf{x}_{n+1} | t_{n+1}) p(t_{n+1} | t_n) = \\ &= \sum_{t_{n+1}} \beta(t_{n+1}) p(\mathbf{x}_{n+1} | t_{n+1}) p(t_{n+1} | t_n)\end{aligned}$$

Инициализация рекуррентного процесса для $\beta(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Вспомним формулу, связывающую значение $\gamma(t_n)$ с $\alpha(t_n)$ и $\beta(t_n)$:

$$\gamma(t_n) = \frac{p(x_1, \dots, x_n, t_n)p(x_{n+1}, \dots, x_N | t_n)}{p(X)} = \frac{\alpha(t_n)\beta(t_n)}{p(X)}$$

Подставляя в формулу значение $n = N$, получаем

$$\gamma(t_N) = p(t_N | X) = \frac{p(X, t_N)\beta(t_N)}{p(X)}$$

Таким образом, $\beta(t_N) = 1$.

Нормировочная константа $p(X)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Значение $\gamma(t_n)$ определено с точностью до нормировочной константы $p(X)$. Однако, на M-шаге значение $\gamma(t_n)$, как правило, входит в числитель и в знаменатель. Таким образом, нормировочная константа сокращается. Например, при вычислении центра нормального распределения:

$$\mu_k = \frac{\sum_{n=1}^N \gamma(t_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(t_{nk})} = \frac{\sum_{n=1}^N \alpha(t_{nk}) \beta(t_{nk}) \mathbf{x}_k}{\sum_{n=1}^N \alpha(t_{nk}) \beta(t_{nk})}$$

Тем не менее, сама нормировочная константа $p(X)$ — это значение правдоподобия, которое может представлять отдельный интерес (например, можно отслеживать возрастание правдоподобия при итерациях EM-алгоритма). Зная $\alpha(t_n)$ и $\beta(t_n)$, правдоподобие может быть легко вычислено как

$$p(X) = \sum_{t_n} \alpha(t_n) \beta(t_n), \quad \forall n \Rightarrow p(X) = \sum_{t_N} \alpha(t_N)$$

Формула для $\xi(\mathbf{t}_{n-1}, \mathbf{t}_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

$$\begin{aligned}\xi(\mathbf{t}_{n-1}, \mathbf{t}_n) &= p(\mathbf{t}_{n-1}, \mathbf{t}_n | X) = \\ &= \frac{p(X | \mathbf{t}_{n-1}, \mathbf{t}_n) p(\mathbf{t}_{n-1}, \mathbf{t}_n)}{p(X)} = \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{t}_{n-1}) p(\mathbf{x}_n | \mathbf{t}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{t}_n) p(\mathbf{t}_n | \mathbf{t}_{n-1}) p(\mathbf{t}_{n-1})}{p(X)} = \\ &= \frac{\alpha(\mathbf{t}_{n-1}) p(\mathbf{x}_n | \mathbf{t}_n) p(\mathbf{t}_n | \mathbf{t}_{n-1}) \beta(\mathbf{t}_n)}{p(X)}\end{aligned}$$

Итоговый EM-алгоритм для случая, когда $p(\mathbf{x}_n | t_n)$ — нормальные распределения

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

- Начальная инициализация параметров $\Theta = (\boldsymbol{\pi}, A, \boldsymbol{\phi})$. Параметры $\boldsymbol{\pi}$ и A можно инициализировать случайно, а $\boldsymbol{\phi}$ с помощью кластеризации данных.
- E-шаг. Вычисление $\boldsymbol{\alpha}(t_n)$ и $\boldsymbol{\beta}(t_n)$ с помощью рекуррентного алгоритма «вперед-назад». Вычисление величин $\gamma(t_{nj})$ и $\xi(t_{n-1,i}, t_{nj})$ и, возможно, правдоподобия $p(X)$.
- M-шаг. Вычисление новых значений параметров:

$$\pi_j^{new} = \frac{\gamma(t_{1j})}{\sum_{i=1}^K \gamma(t_{1i})}, \quad A_{ij}^{new} = \frac{\sum_{n=2}^N \xi(t_{n-1,i}, t_{nj})}{\sum_{k=1}^K \sum_{n=2}^N \xi(t_{n-1,i}, t_{nk})}$$
$$\boldsymbol{\mu}_k^{new} = \frac{\sum_{n=1}^N \gamma(t_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(t_{nk})}, \quad \Sigma_k^{new} = \frac{\sum_{n=1}^N \gamma(t_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(t_{nk})}$$

- Повторять шаги E и M до сходимости (пока значение $p(X)$ или Θ не стабилизируется).

Задача прогнозирования

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

$$\begin{aligned} p(\mathbf{x}_{N+1}|X) &= \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}, \mathbf{t}_{N+1}|X) = \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1})p(\mathbf{t}_{N+1}|X) = \\ &= \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1}) \left(\sum_{\mathbf{t}_N} p(\mathbf{t}_{N+1}, \mathbf{t}_N|X) \right) = \\ &= \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1}) \left(\sum_{\mathbf{t}_N} p(\mathbf{t}_N|X)p(\mathbf{t}_{N+1}|\mathbf{t}_N) \right) = \\ &= \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1}) \left(\sum_{\mathbf{t}_N} p(\mathbf{t}_{N+1}|\mathbf{t}_N) \frac{p(\mathbf{t}_N, X)}{p(X)} \right) = \\ &= \frac{1}{p(X)} \sum_{\mathbf{t}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1}) \left(\sum_{\mathbf{t}_N} p(\mathbf{t}_{N+1}|\mathbf{t}_N) \alpha(\mathbf{t}_N) \right) \end{aligned}$$

Это фактически смесь распределений с компонентами $p(\mathbf{x}_{N+1}|\mathbf{t}_{N+1})$ и весами $w_k = \frac{1}{p(X)} \sum_{\mathbf{t}_N} p(\mathbf{t}_{N+1}|\mathbf{t}_N) \alpha(\mathbf{t}_N)$. Для получения точечного прогноза \mathbf{x}_{N+1}^* можно воспользоваться МСМС.

План лекции

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Метод максимального правдоподобия для СММ

Алгоритм «вперед-назад»

Устойчивые формулы для алгоритма «вперед-назад»

Модельный пример

Необходимость устойчивых вычислений для $\alpha(t_n)$ и $\beta(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Формулы пересчета для $\alpha(t_n)$ и $\beta(t_n)$:

$$\alpha(t_n) = p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} \alpha(t_{n-1}) p(t_n | t_{n-1})$$

$$\beta(t_n) = \sum_{t_{n+1}} \beta(t_{n+1}) p(t_{n+1} | t_n) p(\mathbf{x}_{n+1} | t_{n+1})$$

На практике значения вероятностей $p(t_n | t_{n-1})$ и $p(\mathbf{x}_n | t_n)$ могут быть существенно меньше единицы. В процессе пересчета эти вероятности умножаются друг на друга, и получающиеся значения перестают укладываться в машинную точность.

Устойчивые формулы для $\alpha(t_n)$ I

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Предлагается вместо $\alpha(t_n)$ рассмотреть следующую величину:

$$\hat{\alpha}(t_n) = p(t_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{\alpha(t_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n)}$$

Можно надеяться, что значения $\hat{\alpha}(t_n)$ будут существенно отличны от нуля, т.к. $\sum_{t_n} \hat{\alpha}(t_n) = 1$.

Рассмотрим также величины

$$c_n = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \dots, \mathbf{x}_1).$$

Вычисление этих величин также будет устойчивым, т.к. они имеют смысл одномерных вероятностных распределений.

Устойчивые формулы для $\alpha(t_n)$ II

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Очевидно, что

$$p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n c_i$$

$$\alpha(t_n) = p(t_n | \mathbf{x}_1, \dots, \mathbf{x}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \hat{\alpha}(t_n) \left(\prod_{i=1}^n c_i \right)$$

Подставляя это выражение в формулу пересчета для $\alpha(t_n)$

$$\alpha(t_n) = p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} \alpha(t_{n-1}) p(t_n | t_{n-1}),$$

получаем

$$c_n \hat{\alpha}(t_n) = p(\mathbf{x}_n | t_n) \sum_{t_{n-1}} \hat{\alpha}(t_{n-1}) p(t_n | t_{n-1})$$

Значение c_n определяется из условия нормировки для $\hat{\alpha}(t_n)$.

Устойчивые формулы для $\beta(t_n)$

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Рассмотрим следующую величину

$$\hat{\beta}(t_n) = \frac{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | t_n)}{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)} = \frac{\beta(t_n)}{\prod_{i=n+1}^N c_i}$$

Вычисление данной величины будет устойчивым, т.к. $\hat{\beta}(t_n)$ является отношением двух распределений на $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N)$.
Подставляя выражение для $\hat{\beta}(t_n)$ в формулу пересчета

$$\beta(t_n) = \sum_{t_{n+1}} \beta(t_{n+1}) p(t_{n+1} | t_n) p(\mathbf{x}_{n+1} | t_{n+1}),$$

получаем

$$c_{n+1} \hat{\beta}(t_n) = \sum_{t_{n+1}} \hat{\beta}(t_{n+1}) p(t_{n+1} | t_n) p(\mathbf{x}_{n+1} | t_{n+1})$$

Значения c_n определяются из формул пересчета для $\hat{\alpha}(t_n)$.

Окончательные выражения для E-шага

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Правдоподобие $p(X)$ вычисляется как

$$p(X) = \prod_{n=1}^N c_n$$

Другие необходимые величины вычисляются как

$$\begin{aligned}\gamma(\mathbf{t}_n) &= \hat{\alpha}(\mathbf{t}_n) \hat{\beta}(\mathbf{t}_n) \\ \xi(\mathbf{t}_{n-1}, \mathbf{t}_n) &= \frac{1}{c_n} \hat{\alpha}(\mathbf{t}_{n-1}) p(\mathbf{x}_n | \mathbf{t}_n) p(\mathbf{t}_n | \mathbf{t}_{n-1}) \hat{\beta}(\mathbf{t}_n)\end{aligned}$$

План лекции

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

Метод максимального правдоподобия для СММ

Алгоритм «вперед-назад»

Устойчивые формулы для алгоритма «вперед-назад»

Модельный пример

Модельный пример I

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

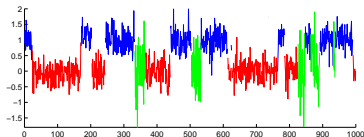
Модельный
пример

В качестве иллюстрации работы EM-алгоритма обучения СММ рассмотрим простой модельный пример. $X \in \mathbb{R}$, $K = 3$, данные сгенерированы из нормальных распределений со следующими параметрами:

$$\begin{aligned}\mu_1 &= 0, & \mu_2 &= 0, & \mu_3 &= 1 \\ \sigma_1^2 &= 0.1, & \sigma_2^2 &= 0.5, & \sigma_3^2 &= 0.1\end{aligned}$$

Априорные вероятности и матрица перехода выбраны следующим образом:

$$\pi = [0.3, 0.2, 0.5], \quad A = \begin{bmatrix} 0.98 & 0.01 & 0.01 \\ 0.01 & 0.97 & 0.02 \\ 0.01 & 0.01 & 0.98 \end{bmatrix}$$



Модельный пример II

Лекция 4.
Скрытые
марковские
модели. Часть 2.

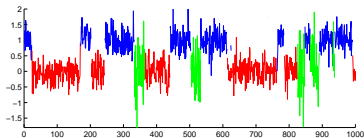
Ветров

Метод
максимального
правдоподобия
для СММ

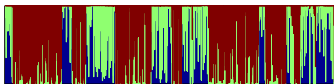
Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

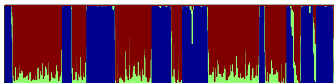
Модельный
пример



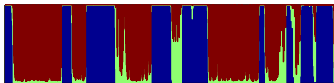
Ит. 1:



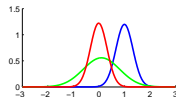
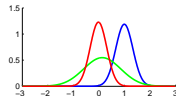
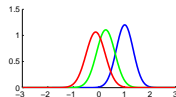
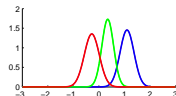
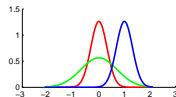
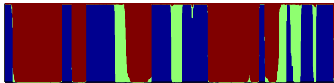
Ит. 5:



Ит. 20:



Ит. 54:



Модельный пример III

Лекция 4.
Скрытые
марковские
модели. Часть 2.

Ветров

Метод
максимального
правдоподобия
для СММ

Алгоритм
«вперед-назад»

Устойчивые
формулы для
алгоритма
«вперед-назад»

Модельный
пример

После 54-ой итерации EM-алгоритма значения параметров были следующие:

$$\boldsymbol{\pi} = [10^{-190}, 10^{-125}, 1], \quad A = \begin{bmatrix} 0.984 & 0.004 & 0.012 \\ 0.013 & 0.949 & 0.038 \\ 0.011 & 0.011 & 0.978 \end{bmatrix}$$

$$\mu_1 = -0.01, \quad \sigma_1^2 = 0.11$$

$$\mu_2 = 0.1, \quad \sigma_2^2 = 0.51$$

$$\mu_3 = 1, \quad \sigma_3^2 = 0.11$$

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

Лекция 6. Снижение размерности в данных. Метод главных компонент.

Курс «Математические основы теории
прогнозирования»

План лекции

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

Метод главных компонент (РСА)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для РСА

Выбор числа главных компонент с помощью байесовского п

Другие модели

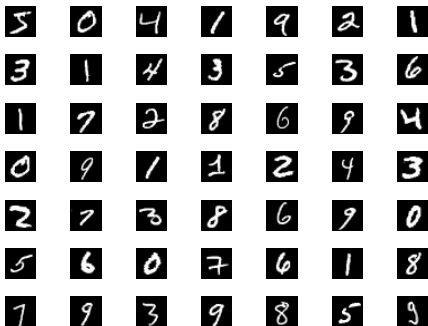
Практический пример: распознавание рукописных цифр

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели



Задача: выбор адекватного признакового описания изображения $\mathbf{x} = (x_1, \dots, x_D)$ для дальнейшего использования алгоритмов распознавания.

Выбор признакового пространства

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

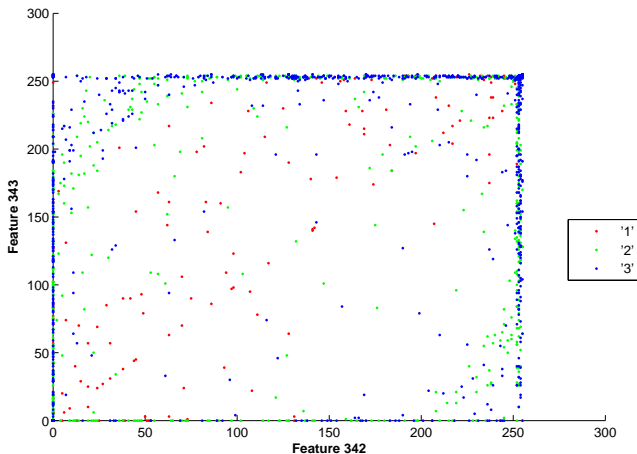
Прямой способ: вытянуть матрицу интенсивностей в вектор.

Такой способ описания данных не является адекватным в силу ряда причин:

- Слишком большое число признаков (для картинки 28×28 получается 784 признака). Как следствие, большое время обработки, проблемы переобучения и т.д.
- Близкие в признаковом пространстве объекты не соответствуют одному классу. Гипотеза компактности является одним из основных предположений большинства методов распознавания.

Иллюстрация

Близкие в признаковом пространстве объекты не соответствуют одному классу.



Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

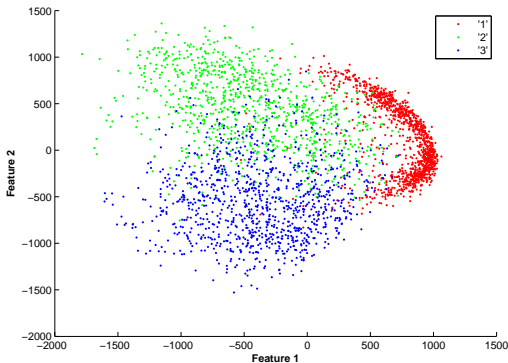
Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Преобразование признаков пространства

После преобразования признаков (с помощью метода главных компонент) количество признаков уменьшается (с 784 до нескольких десятков), причем объекты одного класса образуют компактные области в признаковом пространстве.



Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Сокращение размерности в данных

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

Пусть имеется некоторая выборка $X = \{\mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$.

Цель — представить выборку в пространстве меньшей размерности $d < D$, причем в новом пространстве «схожие» объекты должны образовывать компактные области.

Причины сокращения размерности:

- уменьшение вычислительных затрат при обработке данных
- борьба с переобучением
- сжатие данных для более эффективного хранения информации
- визуализация данных
- извлечение признаков
- интерпретация данных
- ...

План лекции

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для PCA

Выбор числа главных компонент с помощью байесовского п

Другие модели

Идея метода

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Идея метода главных компонент (разложение Карунена-Лоева, Principal Component Analysis, PCA) — проекция данных на гиперплоскость с наименьшей ошибкой проектирования. Эквивалентная формулировка: поиск проекции на гиперплоскость с сохранением большей части дисперсии в данных.

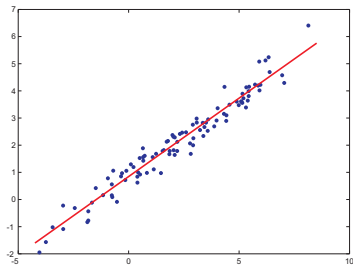


Иллюстрация решения PCA

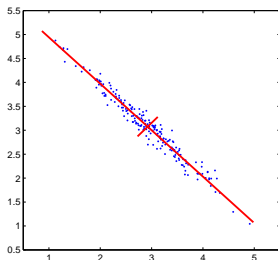
Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

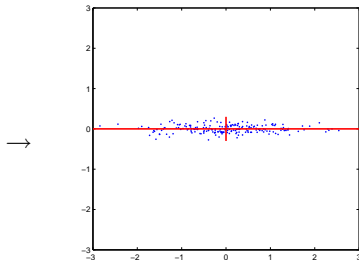
Вероятностный
метод главных
компонент

Другие модели

Выборка в
исходном пространстве



Выборка в пространстве
декоррелированных признаков



- Переносим начало координат в центр выборки
- Поворачиваем оси координат так, чтобы признаки стали декоррелированными
- Отбрасываем направления с низкой дисперсией

Линейное преобразование

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

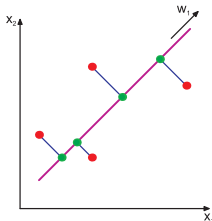
Другие модели

Пусть D — размерность исходного пространства, d — размерность искомого пространства. Будем искать преобразование данных в семействе линейных функций:

$$\mathbf{x} = \boldsymbol{\mu} + t_1 \mathbf{w}_1 + \dots + t_d \mathbf{w}_d = W \mathbf{t} + \boldsymbol{\mu}$$

Здесь $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{t} \in \mathbb{R}^d$ — новые координаты объекта, $W = (\mathbf{w}_1 | \dots | \mathbf{w}_d) \in \mathbb{R}^{D \times d}$, $\boldsymbol{\mu} \in \mathbb{R}^D$.

Графическая интерпретация — проекция данных на гиперплоскость, $\mathbf{w}_1, \dots, \mathbf{w}_d$ — базис гиперплоскости.



Критерий поиска гиперплоскости

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Критерий выбора гиперплоскости — минимальная ошибка проектирования в смысле суммы квадратов отклонений исходных точек и их проекций.

Пусть $\mathbf{w}_1, \dots, \mathbf{w}_D$ — ортонормированный базис в пространстве \mathbb{R}^D , первые d компонент которого — базис искомой гиперплоскости.

Тогда

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{w}_i) \mathbf{w}_i \text{ — исходные точки}$$

$$\hat{\mathbf{x}}_n = \sum_{i=1}^d t_{ni} \mathbf{w}_i + \sum_{i=d+1}^D \mu_i \mathbf{w}_i \text{ — приближение}$$

Критерий

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2 \rightarrow \min_{\mathbf{w}_1, \dots, \mathbf{w}_d, t_1, \dots, t_N, \mu}$$

Решение задачи

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2 \rightarrow \min_{\mathbf{w}_1, \dots, \mathbf{w}_d, t_1, \dots, t_N, \mu}$$

Можно показать, что решением задачи будет следующее:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \text{ — выборочное среднее}$$

$$S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \text{ — выборочная матрица ковариации}$$

$\mathbf{w}_1, \dots, \mathbf{w}_d$ — ортонормированный базис из собственных векторов матрицы S , отвечающих d наибольшим собственным значениям $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$.

$$t_{ni} = \mathbf{x}_n^T \mathbf{w}_i$$

$$\mu_i = \bar{\mathbf{x}}^T \mathbf{w}_i$$

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Преобразование критерия J

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

$$\begin{aligned} J &= \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \left(\sum_{i=1}^D (\mathbf{x}_n^T \mathbf{w}_i) \mathbf{w}_i - \sum_{i=1}^d t_{ni} \mathbf{w}_i - \right. \\ &\quad \left. - \sum_{i=d+1}^D \mu_i \mathbf{w}_i \right)^T \left(\sum_{i=1}^D (\mathbf{x}_n^T \mathbf{w}_i) \mathbf{w}_i - \sum_{i=1}^d t_{ni} \mathbf{w}_i - \sum_{i=d+1}^D \mu_i \mathbf{w}_i \right) = \\ &\quad \frac{1}{N} \sum_{n=1}^N \left(\sum_{i=1}^D (\mathbf{x}_n^T \mathbf{w}_i)^2 + \sum_{i=1}^d t_{ni}^2 + \sum_{i=d+1}^D \mu_i^2 - \right. \\ &\quad \left. 2 \sum_{i=1}^d (\mathbf{x}_n^T \mathbf{w}_i) t_{ni} - 2 \sum_{i=d+1}^D \mu_i (\mathbf{x}_n^T \mathbf{w}_i) \right) \end{aligned}$$

$$\frac{\partial}{\partial t_{ni}} : -2(\mathbf{x}_n^T \mathbf{w}_i) + 2t_{ni} = 0 \Rightarrow t_{ni} = \mathbf{x}_n^T \mathbf{w}_i$$

$$\frac{\partial}{\partial \mu_i} : \sum_{n=1}^N (-2(\mathbf{x}_n^T \mathbf{w}_i) + 2\mu_i) = 0 \Rightarrow \mu_i = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n^T \mathbf{w}_i = \bar{\mathbf{x}}^T \mathbf{w}_i$$

Преобразование критерия II

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

$$t_{ni} = \mathbf{x}_n^T \mathbf{w}_i, \quad \mu_i = \bar{\mathbf{x}}^T \mathbf{w}_i$$

$$\mathbf{x}_n - \hat{\mathbf{x}}_n = \sum_{i=d+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{w}_i\} \mathbf{w}_i$$

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=d+1}^D (\mathbf{x}_n^T \mathbf{w}_i - \bar{\mathbf{x}}^T \mathbf{w}_i)^2 = \sum_{i=d+1}^D \mathbf{w}_i^T S \mathbf{w}_i \rightarrow \min_{\mathbf{w}_{d+1}, \dots, \mathbf{w}_D}$$

Здесь $S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$ — выборочная матрица ковариации данных.

Решение I

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

$$J = \sum_{i=d+1}^D \mathbf{w}_i^T S \mathbf{w}_i \rightarrow \min_{\mathbf{w}_{d+1}, \dots, \mathbf{w}_D}$$
$$\mathbf{w}_i^T \mathbf{w}_j = \delta_{ij}$$

Можно показать, что оптимум достигается, когда в качестве $\mathbf{w}_{d+1}, \dots, \mathbf{w}_D$ выбираются собственные вектора матрицы ковариации S , отвечающие наименьшим собственным значениям $\lambda_{d+1} \geq \dots \geq \lambda_D$. Тогда

$$J = \sum_{i=d+1}^D \lambda_i$$

Решение II

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

Рассмотрим оптимизацию по w_i :

$$\tilde{J} = w_i^T S w_i \rightarrow \min_{w_i}$$
$$w_i^T w_i = 1$$

Функция Лагранжа:

$$L = w_i^T S w_i + \lambda_i (1 - w_i^T w_i) \rightarrow \text{extr}$$

Приравнивая производную к нулю, получаем:

$$S w_i = \lambda_i w_i$$

Подставляя обратно в критерий, получаем:

$$\tilde{J} = w_i^T S w_i = \lambda_i$$

Альтернативная интерпретация PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Рассмотрим в качестве критерия проектирования максимизацию дисперсии спроектированных данных.

Пусть $d = 1$, т.е. необходимо найти некоторое направление \mathbf{w}_1 : $\mathbf{w}_1^T \mathbf{w}_1 = 1$. Каждый объект выборки \mathbf{x}_n проектируется в $\mathbf{w}_1^T \mathbf{x}_n$. Тогда дисперсия проекций вычисляется как

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^T \mathbf{x}_n - \mathbf{w}_1^T \bar{\mathbf{x}})^2 = \mathbf{w}_1^T S \mathbf{w}_1$$

Здесь

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_i, \quad S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

Альтернативная интерпретация II

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

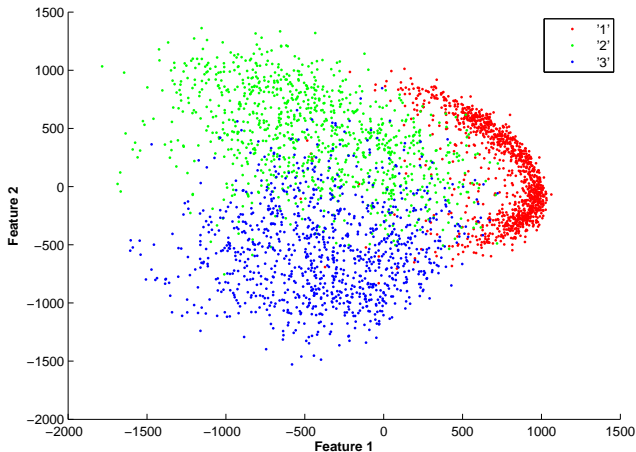
Другие модели

Теперь необходимо максимизировать дисперсию проекций $\mathbf{w}_1^T S \mathbf{w}_1$ при ограничении $\mathbf{w}_1^T \mathbf{w}_1 = 1$. Таким образом, мы приходим к точно такому же критерию, который возникал в случае минимизации невязки. Максимальная дисперсия достигается для собственного вектора выборочной матрицы ковариации, отвечающего максимальному собственному значению.

Рассматривая аналогично проектирование на новые направления, ортогональные уже найденным, получим, что наилучшее d -мерное линейное подпространство определяется собственными векторами выборочной матрицы ковариации, отвечающие d максимальным собственным значениям.

Пример с распознаванием рукописных цифр

Проекция данных на первые два признака (два собственных вектора матрицы ковариации, отвечающих наибольшему собственным значениям).



Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Интерпретация новых признаков

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

По построению каждой точке $\mathbf{x} \in \mathbb{R}^D$ соответствует некоторая картинка. Новые признаки $\mathbf{t} \in \mathbb{R}^d$ — проекции на выбранные направления (собственные вектора). Изменение одного признака в новом пространстве \mathbb{R}^d соответствует движению вдоль собственного вектора в исходном пространстве \mathbb{R}^D .

Движение вдоль первого признака:



Интерпретация: изменение ширины цифры

Движение вдоль второго признака:



Интерпретация: изменение характерного наклона цифры

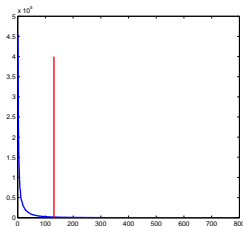
Выбор размерности подпространства d в PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели



- Рассматриваем собственные значения в порядке убывания $\lambda_1 \geq \lambda_2 \geq \dots \lambda_D$
- Выбираем d так, что $\sum_{i=d+1}^D \lambda_i \ll \sum_{i=1}^d \lambda_i$ либо, чтобы

$$\sum_{i=1}^d \lambda_i \geq \gamma \sum_{i=1}^D \lambda_i$$

Здесь γ — доля объясняемой дисперсии, типичные значения 0.95, 0.99.

Неоднозначность решения PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

- Если все собственные значения выборочной матрицы ковариации различны, т.е. $\lambda_1 > \lambda_2 > \dots > \lambda_D$, то ортонормированный базис из собственных векторов определен однозначно, т.е. гиперплоскость проекции определена однозначно. Тем не менее, в полученной гиперплоскости базис может быть выбран произвольным образом.
- Если существуют одинаковые собственные вектора, то тогда гиперплоскость проекции определена не однозначно.

В PCA существует произвол в выборе координат объектов в новом пространстве

План лекции

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для PCA

Выбор числа главных компонент с помощью байесовского подхода

Другие модели

Мотивация

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

РСА может быть сформулирован как вероятностная модель с латентными переменными, для оптимизации которой используется метод максимального правдоподобия.

Преимущества вероятностной постановки:

- Вычисление правдоподобия на тестовой выборке позволяет непосредственно сравнивать различные вероятностные модели
- EM-алгоритм для РСА позволяет быстро находить решения в ситуациях, когда требуется небольшое число лидирующих главных компонент, а также позволяет избежать вычисление выборочной матрицы ковариации в качестве промежуточного шага
- Возможность применения байесовского подхода для автоматического определения количества главных компонент (по аналогии с методом релевантных векторов)
- Можно работать со смесью РСА и использовать вариант EM-алгоритма для обучения в такой модели

Вероятностная модель PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Пусть имеется выборка $\{\mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$. Рассмотрим для каждого объекта \mathbf{x}_n латентную переменную $\mathbf{t}_n \in \mathbb{R}^d$ как координаты объекта \mathbf{x}_n в подпространстве, вообще говоря, меньшей размерности $d < D$. Определим априорное распределение в пространстве латентных переменных как

$$p(\mathbf{t}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, I)$$

Модель наблюдаемых переменных \mathbf{x} представляет собой линейное преобразование латентной переменной с добавлением гауссовского шума с единой дисперсией по всем направлениям:

$$p(\mathbf{x}|\mathbf{t}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{t} + \boldsymbol{\mu}, \sigma^2 I)$$

Здесь $\mathbf{W} \in \mathbb{R}^{D \times d}$, $\boldsymbol{\mu} \in \mathbb{R}^D$.

Иллюстрация вероятностной модели PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

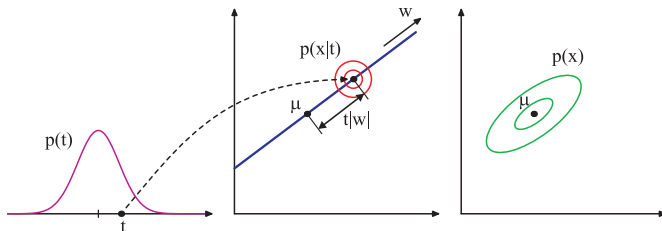
Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели



Функция правдоподобия

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Функция правдоподобия может быть вычислена как

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{t})p(\mathbf{t})d\mathbf{t}$$

Этот интеграл является сверткой двух нормальных распределений и может быть вычислен аналитически. В результате получается снова нормальное распределение

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, C), \quad C = WW^T + \sigma^2I$$

Действительно,

$$\begin{aligned}\mathbb{E}\mathbf{x} &= \mathbb{E}(W\mathbf{t} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}) = \boldsymbol{\mu} \\ \mathbb{E}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T &= \mathbb{E}(W\mathbf{t} + \boldsymbol{\varepsilon})(W\mathbf{t} + \boldsymbol{\varepsilon})^T = \\ &= W\mathbb{E}\mathbf{t}\mathbf{t}^T W^T + \mathbb{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T = WW^T + \sigma^2I\end{aligned}$$

Инвариантность относительно поворота координат в пространстве латентных переменных

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для PCA

Выбор числа главных компонент с помощью байесовского подхода

Другие модели

Пусть R — некоторая ортогональная матрица. Рассмотрим матрицу базисных векторов $\tilde{W} = WR$. В этом случае вероятностная модель PCA приводит к тому же распределению, что и в случае матрицы W . Действительно, т.к. $RR^T = I$, то

$$C = \tilde{W}\tilde{W}^T + \sigma^2 I = WRR^T W^T + \sigma^2 I = WW^T + \sigma^2 I$$

Решение

Логарифм правдоподобия для рассматриваемой модели выглядит как

$$\begin{aligned}\log p(X|\mu, W, \sigma^2) &= \sum_{n=1}^N \log p(\mathbf{x}_n|\mu, W, \sigma^2) = \\ &= -\frac{ND}{2} \log 2\pi - \frac{N}{2} \log \det C - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T C^{-1} (\mathbf{x}_n - \mu)\end{aligned}$$

Дифференцируя правдоподобие и приравнявая производную к нулю, получаем:

$$\begin{aligned}\mu_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, & \sigma_{ML}^2 &= \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i \\ W_{ML} &= U_d(L_d - \sigma^2 I)^{1/2} R\end{aligned}$$

Здесь $U_d \in \mathbb{R}^{D \times d}$ — матрица, состоящая из d собственных векторов выборочной матрицы ковариации, отвечающие d наибольшим собственным значениям $\lambda_1, \dots, \lambda_d$,

$L_d = \text{diag}(\lambda_1, \dots, \lambda_d)$, R — произвольная ортогональная матрица.

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Апостериорное распределение латентной переменной

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для PCA

Выбор числа главных компонент с помощью байесовского подхода

Другие модели

Можно показать, что

$$p(\mathbf{t}|\mathbf{x}) = \mathcal{N}(\mathbf{t}|M^{-1}W^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2M^{-1})$$

Здесь $M = W^T W + \sigma^2 I$. Для решения задачи визуализации данных в пространстве латентных переменных требуется мат. ожидание латентной переменной по своему апостериорному распределению:

$$\mathbb{E}(\mathbf{t}) = M^{-1}W_{ML}(\mathbf{x} - \bar{\mathbf{x}})$$

EM-алгоритм в общем виде

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Пусть имеется вероятностная модель, в которой:

- часть переменных X — известна
- часть переменных T — не известна
- имеется некоторый набор параметров Θ

Требуется оценить набор параметров Θ с помощью метода максимального правдоподобия:

$$p(X|\Theta) = \int p(X, T|\Theta)dT \rightarrow \max_{\Theta}$$

Схема EM-алгоритма

Требуется найти максимум правдоподобия в вероятностной модели со скрытыми переменными:

$$p(X|\Theta) = \int p(X, T|\Theta)dT \rightarrow \max_{\Theta} \Leftrightarrow \log \left(\int p(X, T|\Theta)dT \right) \rightarrow \max_{\Theta}$$

- **Е-шаг.** Фиксируется значение параметров Θ_{old} .
Оценивается апостериорное распределение на скрытые переменные $p(T|X, \Theta_{old})$, и полное правдоподобие усредняется по полученному распределению:

$$\mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta) = \int \log p(X, T|\Theta) p(T|X, \Theta_{old}) dT$$

- **М-шаг.** Фиксируется апостериорное распределение $p(T|X, \Theta_{old})$, и производится поиск новых значений параметров Θ_{new} :

$$\Theta_{new} = \arg \max_{\Theta} \mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta)$$

- Шаги Е и М повторяются до сходимости.

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Нижняя оценка для логарифма правдоподобия

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (РСА)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для РСА

Выбор числа главных компонент с помощью байесовского подхода

Другие модели

Можно показать, что

$$\log p(X|\Theta) \geq \mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta) + Const \text{ и } \Leftrightarrow \Theta = \Theta_{old}$$

Здесь *Const* - некоторая константа, не зависящая от Θ .

Иллюстрация EM-алгоритма

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

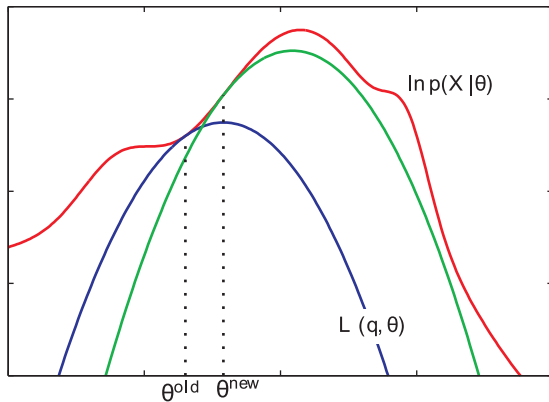
Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели



EM-алгоритм для PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

$$\mathbb{E}\mathbf{t}_n = M^{-1}W^T(\mathbf{x} - \bar{\mathbf{x}})$$

$$\mathbb{E}\mathbf{t}_n\mathbf{t}_n^T = \sigma^2M^{-1} + \mathbb{E}\mathbf{t}_n\mathbb{E}\mathbf{t}_n^T$$

$$W_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})\mathbb{E}\mathbf{t}_n^T \right] \left[\sum_{n=1}^N \mathbb{E}\mathbf{t}_n\mathbf{t}_n^T \right]^{-1}$$

$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N [\|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2\mathbb{E}\mathbf{t}_n^T W_{new}^T (\mathbf{x}_n - \bar{\mathbf{x}}) + \text{tr}(\mathbb{E}\mathbf{t}_n\mathbf{t}_n^T W_{new}^T W_{new})]$$

Мотивация EM-алгоритма для PCA

Несмотря на существование решения для W и σ^2 в явном виде, использование EM-алгоритма может быть предпочтительным в ряде случаев:

- EM-алгоритм избегает вычисления выборочной матрицы ковариации (сложность $O(ND^2)$) и поиска ее собственных значений (сложность $O(D^3)$). Самые сложные операции в EM-алгоритме требуют $O(NDd)$ и $O(d^3)$, что может дать существенный выигрыш в скорости для данных больших размерностей.
- EM-алгоритм может быть применен для модели факторного анализа, для которой не существует решения в явном виде. Эта модель полностью повторяет вероятностную модель для PCA, однако различные факторы могут иметь разную дисперсию:

$$p(\mathbf{x}|\mathbf{t}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{t} + \boldsymbol{\mu}, \Psi)$$

Здесь Ψ — диагональная матрица.

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для PCA

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Проблема выбора количества главных компонент d

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (РСА)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для РСА

Выбор числа главных компонент с помощью байесовского подхода

Другие модели

- Параметр d является типичным структурным параметром и должен задаваться до старта EM-алгоритма
- Одним из возможных способов выбора d является визуальное оценивание распределения собственных значений выборочной матрицы ковариации с выделением области, в которой собственные значения не отличаются существенно от нуля. К сожалению, на практике такой порог часто провести не удается.
- Разумной альтернативой является байесовский подход

Априорное распределение на базисные вектора

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Выберем в качестве априорного распределения на веса нормальное распределение с отдельными коэффициентами регуляризации для каждого базисного вектора:

$$p(W|\alpha) = \prod_{i=1}^D \left(\frac{\alpha_i}{2\pi}\right)^{D/2} \exp\left(-\frac{1}{2}\alpha_i \mathbf{w}_i^T \mathbf{w}_i\right)$$

Здесь \mathbf{w}_i — i -ая колонка матрицы W . Для подбора коэффициентов α можно использовать максимизацию обоснованности:

$$p(X|\alpha, \mu, \sigma^2) = \int p(X|W, \mu, \sigma^2)p(W|\alpha)dW$$

Данный интеграл не берется аналитически, поэтому возможными путями являются приближение Лапласа и вариационный подход.

Процедура обучения

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Вероятностная
формулировка
метода

EM-алгоритм
для РСА

Выбор числа
главных
компонент с
помощью
байесовского
подхода

Другие модели

Использование приближения Лапласа приводит к следующей формуле пересчета коэффициентов регуляризации:

$$\alpha_i^{new} = \frac{D}{\mathbf{w}_i^T \mathbf{w}_i}$$

Таким образом, процедура обучения является итерационной. На каждой итерации для текущих значений α с помощью EM-алгоритма оцениваются W, μ, σ^2 , а затем коэффициенты α пересчитываются. В EM-алгоритме формула для пересчета W выглядит следующим образом:

$$W_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E} \mathbf{t}_n^T \right] \left[\sum_{n=1}^N \mathbb{E} \mathbf{t}_n \mathbf{t}_n^T + \sigma^2 A \right]^{-1}$$

Здесь $A = \text{diag}(\alpha_1, \dots, \alpha_D)$.

План лекции

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Вероятностная формулировка метода

EM-алгоритм для PCA

Выбор числа главных компонент с помощью байесовского п

Другие модели

Недостатки PCA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

- Метод главных компонент способен находить только линейные подпространства исходного пространства, которые «объясняют» данные с высокой точностью. На практике поверхность, вдоль которой располагаются данные, может быть существенно нелинейной
- PCA является инвариантным относительно поворота координат в пространстве латентных переменных. Это означает, что восстановление значений латентных переменных является неоднозначным. В ряде случаев такая ситуация может быть неадекватной, например, в задаче разделения независимых источников, представленных линейной смесью с неизвестными коэффициентами.

Метод независимых компонент, ICA

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Пусть наблюдаемые данные $\mathbf{x} = (x_1, \dots, x_d)$ являются линейной комбинацией независимых источников $\mathbf{t} = (t_1, \dots, t_d)$:

$$\mathbf{x} = W\mathbf{t},$$

причем размерности \mathbf{x} и \mathbf{t} совпадают, а матрица W неизвестна.

Если преобразование W является невырожденным, то $\mathbf{t} = W^T \mathbf{x}$. Рассмотрим задачу поиска одного источника t_1 . Будем искать t_1 как линейную комбинацию наблюдаемых данных $t_1 = \mathbf{a}^T \mathbf{x}$. Задача – определить вектор \mathbf{a} .

Если \mathbf{a} совпадает с колонкой матрицы W , то t_1 будет одним из искомым источников.

$$t_1 = \mathbf{a}^T \mathbf{x} = \mathbf{a}^T W \mathbf{t} = \{\mathbf{z} \triangleq W^T \mathbf{a}\} = \mathbf{z}^T \mathbf{t} = z_1 t_1 + \dots + z_d t_d$$

По центральной предельной теореме сумма случайных величин приближается к нормальному распределению. Следовательно, чем больше слагаемых в сумме $z_1 t_1 + \dots + z_d t_d$, тем более она похожа на гауссиану. Отсюда

\mathbf{a} должен быть таким, чтобы $\mathbf{a}^T \mathbf{x}$ было как можно меньше похоже на гауссиану.

Критерии негауссовости

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

- Экссесс.

$$\text{kurt} = \mathbb{E}t^4 - 3(\mathbb{E}t^2)^2$$

В предположении, что случайная величина t нормализована, $\text{kurt} = \mathbb{E}t^4$. Для гауссианы значение эксцесса равно 0, для остальных распределений может быть как положительным, так и отрицательным. Поэтому в качестве меры негауссовости выбирают модуль или квадрат эксцесса.

- Негэнтропия.

$$H(t) = - \int \log p(t)p(t)dt$$

У гауссовского распределения энтропия $H(t)$ является максимальной среди всех распределений с одинаковой матрицей ковариации. Поэтому в качестве меры негауссовости используется негэнтропия

$$H(t_{\text{gauss}}) - H(t)$$

Здесь t_{gauss} — гауссиана с той же матрицей ковариации, что и t . Негэнтропия является неотрицательной, и ее нужно максимизировать для поиска независимых источников.

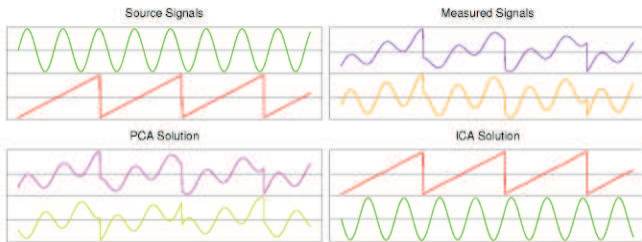
ICA vs. PCA. Пример.

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели



Анализ независимых факторов, Independent Factor Analysis, IFA

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (РСА)

Вероятностный метод главных компонент

Другие модели

В методе РСА и факторном анализе в качестве априорного распределения латентных переменных предполагается выбирать нормальное распределение с центром в нуле и некоторой дисперсией, возможно, различной для разных направлений. В анализе независимых факторов в качестве априорного распределения предлагается выбрать независимое:

$$p(\mathbf{t}) = \prod_{i=1}^d p(t_i), \quad p(\mathbf{x}|\mathbf{t}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{t} + \boldsymbol{\mu}, \Psi)$$

С одной стороны, данная модель является существенным обобщением РСА и позволяет находить, вообще говоря, нелинейные зависимости, с другой стороны, модель является достаточно простой, чтобы можно было предложить варианты EM-алгоритма для ее оптимизации.

Одним из возможных вариантов реализации данного метода является моделирование распределений отдельных факторов с помощью смеси гауссиан с разным числом компонент для разных факторов.

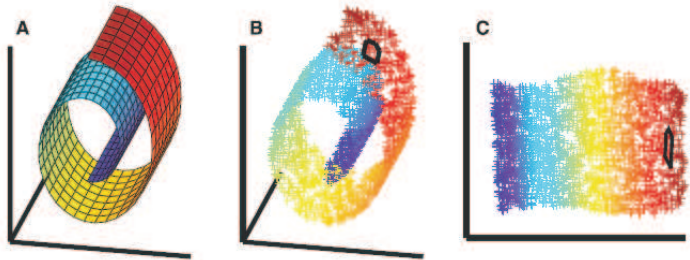
Локальное линейное погружение, Local Linear Embedding, LLE

Лекция 6.
Снижение размерности в данных. Метод главных компонент.

Метод главных компонент (PCA)

Вероятностный метод главных компонент

Другие модели



Если данные располагаются вдоль некоторой поверхности, то разумно искать такое преобразование признакового пространства, при котором близкие объекты на этой поверхности были бы близки в новом пространстве. При этом близкие в евклидовом смысле объекты в результате такого преобразования могут оказаться очень далекими.

PCA не позволяет находить преобразования с подобными свойствами!

Схема LLE

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели

Основная идея — искать преобразование с сохранением отношения соседства между объектами

- Восстановление структуры соседства на обучающей выборке:

$$\varepsilon(W) = \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{j \in \text{Neighbours}(\mathbf{x}_n)} W_{nj} \mathbf{x}_j \right\|^2 \rightarrow \min_W$$
$$\sum_j W_{nj} = 0$$

- Поиск координат в новом пространстве со схожей структурой соседства:

$$\Phi(T) = \sum_{n=1}^N \left\| \mathbf{t}_n - \sum_{j \in \text{Neighbours}(\mathbf{x}_n)} W_{nj} \mathbf{t}_j \right\|^2 \rightarrow \min_T$$
$$\sum_{n=1}^N \mathbf{t}_n = \mathbf{0}, \quad \frac{1}{N} \sum_{n=1}^N \mathbf{t}_n \mathbf{t}_n^T = I$$

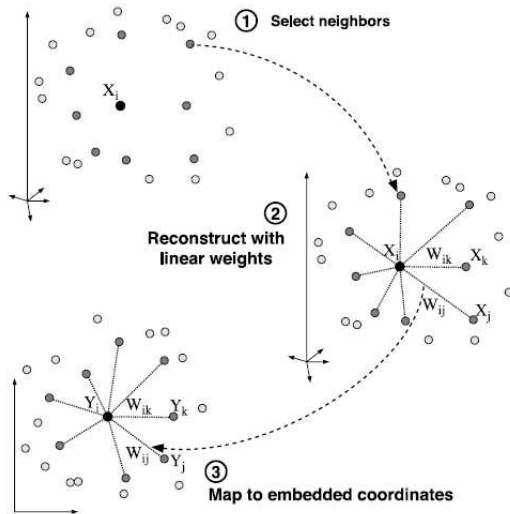
Иллюстрация LLE

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

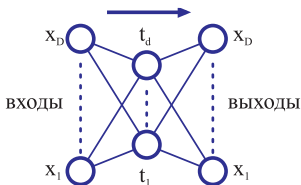
Вероятностный
метод главных
компонент

Другие модели



Автоассоциативные нейронные сети

Рассмотрим многослойный персептрон следующего вида:



Здесь $d < D$. Задача нейронной сети — объяснить как можно точнее входные данные с помощью выходных. Таким образом, функционалом качества обучения является следующий:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(\mathbf{x}_n, \mathbf{w}) - \mathbf{x}_n\|^2$$

Если в сети только один скрытый слой, то даже с использованием нелинейной сигмоидной функции активации результат обучения сети совпадает с результатом PCA

Лекция 6.
Снижение
размерности
в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

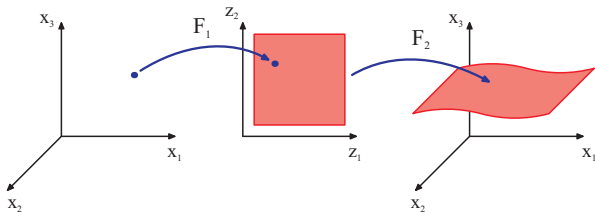
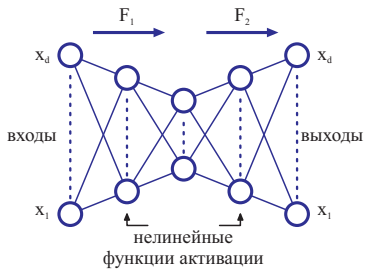
Другие модели

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели



Generative Topographic Mapping, GTM

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (РСА)

Вероятностный
метод главных
компонент

Другие модели

Метод GTM используется, в основном, для визуализации данных и является вероятностным обобщением самоорганизующихся сетей Кохонена.

В модели GTM предполагается, что данные образуются в результате нелинейного преобразования латентных переменных с добавлением нормального шума:

$$p(\mathbf{x}|\mathbf{t}) = \mathcal{N}(\mathbf{x}|\mathbf{y}(\mathbf{t}, W), \sigma^2 I)$$

Пространство латентных переменных является двухмерным, а распределение латентных переменных представляет собой сумму дельта-функций с центрами в некоторой регулярной сетке:

$$p(\mathbf{t}) = \frac{1}{l} \sum_{j=1}^l \delta(\mathbf{t} - \mathbf{t}_j)$$

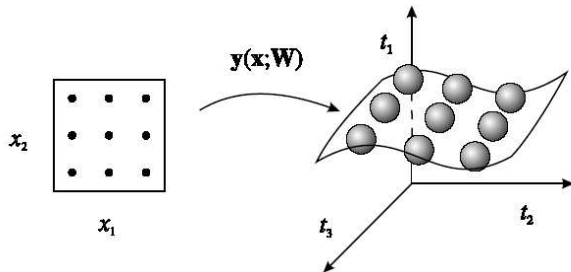
Иллюстрация GTM

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

Вероятностный
метод главных
компонент

Другие модели



Функция правдоподобия выглядит следующим образом:

$$\log p(X|W, \sigma^2) = \sum_{n=1}^N \log \left[\frac{1}{l} \sum_{j=1}^l p(\mathbf{x}_n | \mathbf{t}_j, W, \sigma^2) \right]$$

Иллюстрация GTM

Лекция 6.
Снижение
размерности в
данных. Метод
главных
компонент.

Метод главных
компонент (PCA)

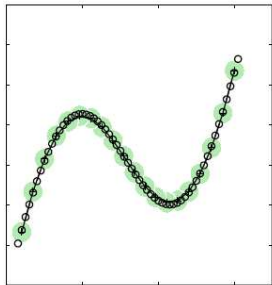
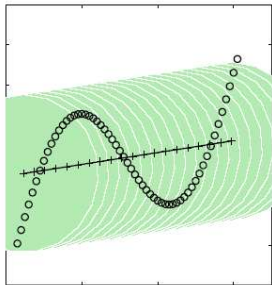
Вероятностный
метод главных
компонент

Другие модели

В качестве функции регрессии предлагается взять обобщенную линейную:

$$\mathbf{y}(\mathbf{t}, W) = W\phi(\mathbf{t})$$

Здесь $\{\phi_j(\mathbf{t})\}_{j=1}^d$ — набор фиксированных базисных функций, в качестве которых могут выступать, например, гауссианы с центрами в регулярной сетке. Тогда можно предложить эффективный EM-алгоритм для оптимизации такой модели.



Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные
функции

Альфа-
расширение

Разрезы графов

Ю. И. Журавлев¹, Д. П. Ветров¹

¹МГУ, ВМиК, каф. ММП

Курс «Математические основы теории
прогнозирования»

План

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

Потоки в сетях

Разрезы графов

Ветров

Ликбез

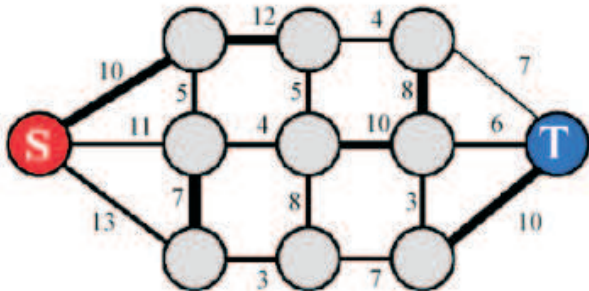
Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Рассмотрим неориентированный граф с двумя выделенными вершинами (стоком t и истоком s)
- Пусть с каждым ребром $(u, v) \in E$ ассоциировано некоторое неотрицательное число $c(u, v) \geq 0$ — пропускная способность



Потоки в сетях

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Назовем потоком неотрицательную функцию $f(u, v)$, определенную на ребрах графа, такую что

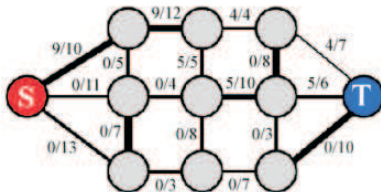
$$f(u, v) \leq c(u, v)$$

$$\sum_{v:(u,v) \in E} f(u, v) - \sum_{v:(v,u) \in E} f(v, u) = 0, \quad \forall u \neq \{s, t\}$$

- Первое условие ограничивает поток через ребро его пропускной способностью, а второе гарантирует отсутствие источников и стоков вне выделенной пары вершин
- Задача поиска максимального потока состоит в максимизации величины

$$M(f) = \sum_{v:(s,v) \in E} f(s, v) \rightarrow \max_f$$

по всем допустимым потокам



Разрезы графов

Разрезы графов

Ветров

Ликбез

Марковские сети

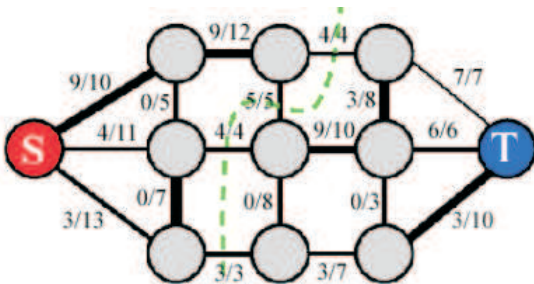
Разрезы графов

Субмодулярные функции

Альфа-расширение

- $(s - t)$ -разрезом графа называется разбиение вершин графа на два непересекающихся множества S и T , такие что $s \in S$, $t \in T$
- Величиной разреза называется сумма пропускных способностей всех ребер, один конец которых находится в множестве S , а другой — в множестве T

$$c(S, T) = \sum_{(u,v) \in E, u \in S, v \in T} c(u, v)$$



Теорема Форда-Фалкерсона

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Известная теорема Форда-Фалкерсона гласит, что максимальный поток в сети равен ее минимальному разрезу
- Существует эффективный (полиномиальной сложности) алгоритм решения задачи поиска минимального разреза в графе
- Задачи поиска максимального потока и минимального разреза являются двойственными



Марковские решетки

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

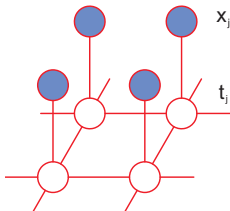
Субмодулярные функции

Альфа-расширение

- В дальнейшем будем рассматривать т.н. марковские решетки, в которых размер максимальной клики не превосходит двух
- В этом случае совместное распределение переменных марковской сети выражается формулой

$$p(Y) = \frac{1}{Z} \prod_{(y_i, y_j) \in E} \psi_{ij}(y_i, y_j)$$

- Наиболее типичным примером таких сетей являются изображения



Марковские решетки с бинарными переменными

Разрезы графов

Ветров

Ликбез

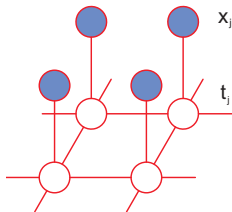
Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- В дальнейшем будем рассматривать марковские решетки такого вида



- Необходимо по наблюдаемым переменным X восстановить наиболее вероятные значения скрытых переменных T

$$T_{MP} = \arg \max_T P(T|X)$$

- Остановимся на важном частном случае, когда скрытые переменные бинарные $t \in \{0, 1\}$

Марковские решетки с бинарными переменными

Разрезы графов

Ветров

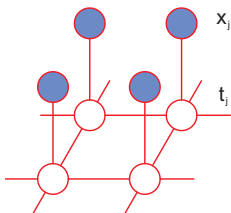
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Распределение скрытых переменных марковской сети в этом случае выглядит так

$$p(T|X) = \frac{p(X, T)}{p(X)} = \prod_{(i,j) \in E} \psi_{ij}(t_i, t_j) \prod_i \psi_i(x_i, t_i) \times \text{Const}$$

- В энергетической нотации задача максимизации этого распределения принимает вид минимизации энергии, что более удобно с вычислительной точки зрения

$$E(T|X) = \sum_{(i,j) \in E} E_{ij}(t_i, t_j) + \sum_i E_i(x_i, t_i) = - \sum_{(i,j) \in E} \log \psi_{ij}(t_i, t_j) - \sum_i \log \psi_i(x_i, t_i) \rightarrow \min_T$$

Пример задачи сегментации

Разрезы графов

Ветров

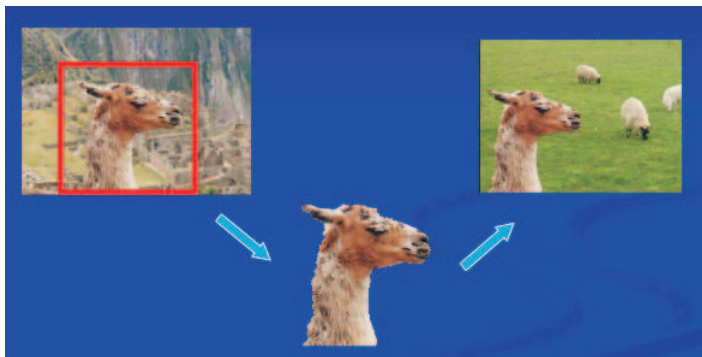
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



Значения скрытых переменных кодируют принадлежность каждого пикселя к объекту либо к фону. Использование графических моделей позволяет учесть, что соседние пиксели чаще всего относятся к одному классу

Репараметризация

Разрезы графов

Ветров

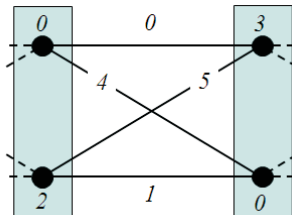
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Рассмотрим две соседние вершины t_i и t_j , каждая из которых может принимать два значения
- Тогда функция $E_i(x_i, t_i)$ задается двумя значениями (при известном x_i), а функция $E_{ij}(t_i, t_j)$ — четырьмя
- Для сведения к задаче о поиске минимального разреза нам понадобится выполнить т.н. репараметризацию, сделав «веса» горизонтальных ребер нулевыми

Репараметризация

Разрезы графов

Ветров

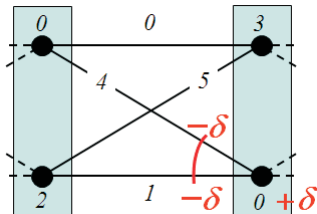
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Вычитая одинаковые значения из двух ребер, сходящихся в одну вершину и прибавляя это же значение к весу самой вершины, мы получаем эквивалентный функционал энергии
- Значения такой энергии в каждой точке совпадают со значением исходной энергии
- Применение этой процедуры позволяет путем изменения функции E_i получить эквивалентный энергетический функционал, в котором $E_{ij}(0,0) = E_{ij}(1,1) = 0$ и $E_{ij}(0,1) = E_{ij}(1,0)$ для всех $(i,j) \in E$

Сведение к разрезу графов

Разрезы графов

Ветров

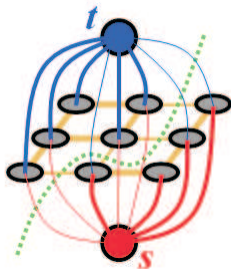
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Определим на следующем графе пропускную способность таким образом:
 $c(S, t_i) = E_i(x_i, 0)$, $c(T, t_i) = E_i(x_i, 1)$, $c(t_i, t_j) = E_{ij}(0, 1)$, $\forall (i, j) \in E$
- Поиск минимального разреза в таком графе отвечает минимизации энергии

$$E(T|X) = \sum_{(i,j) \in E} E_{ij}(t_i, t_j) + \sum_i E_i(x_i, t_i)$$

т.е. поиску наиболее вероятных значений T

Сегментация с семенами

Разрезы графов

Ветров

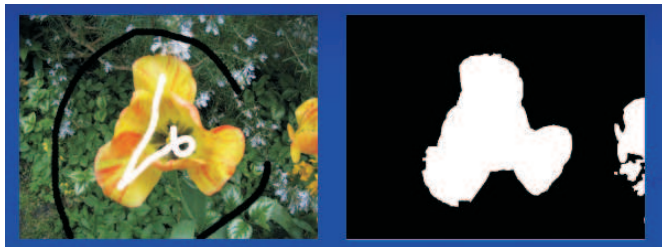
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Часто для некоторых пикселей известно заранее, к какому классу они принадлежат
- Например, пользователь может задать фрагменты изображения и фона (семена)

Сведение к разрезу графов

Разрезы графов

Ветров

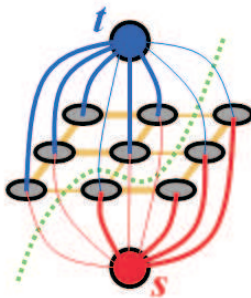
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Пусть O — семена объекта ($t_i = 1$), а B — семена фона ($t_i = 0$)
- Тогда достаточно задать $c(S, t_i) = +\infty, \forall t_i \in O$ и $c(T, t_i) = +\infty, \forall t_i \in B$
- Этим мы запретим соответствующие разрезы, сделав невозможным отнесение семян объекта к фону и наоборот

Способы введения унарного слагаемого

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

По смыслу унарное (относительно скрытых переменных) слагаемое $E_i(x_i, t_i)$ задает насколько данный пиксель соответствует тому или иному классу. Оно может отражать следующую информацию:

- Цветовая модель — показывает насколько появление тех или иных цветов более вероятно в данном классе
- Позиционная модель — показывает априорные предположения о положении данного класса на изображении
- Текстурная модель — показывает насколько текстура окрестности пикселя вероятна для данного пикселя

Способы введения парного слагаемого

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

Парное слагаемое $E_{ij}(t_i, t_j)$ отражает степень взаимозависимостей классов соседних пикселей. Наиболее распространенными примерами являются

- Модель Поттса: $E_{ij}(t_i, t_j) = 1 - \delta(t_i, t_j)$ — штраф за несовпадение классов соседних пикселей
- Штраф за несовпадение классов с учетом контраста

$$E_{ij}(t_i, t_j) = \exp\left(-\frac{(x_i - x_j)^2}{2\sigma^2}\right) (1 - \delta(t_i, t_j))$$

Чем сильнее различаются цвета (интенсивности) пикселей x_i , тем меньше штраф за несовпадение классов

- Заметим, что во втором случае парное слагаемое зависит от наблюдаемых переменных x_i и x_j - такая конструкция называется условным случайным полем

Субмодулярность

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Назовем энергию субмодулярной, если для всех ее парных слагаемых верно

$$E_{ij}(0, 0) + E_{ij}(1, 1) \leq E_{ij}(0, 1) + E_{ij}(1, 0)$$

- Условие субмодулярности является в некотором смысле аналогом выпуклости для функций бинарного переменного
- Унарное слагаемое при этом может быть произвольным
- Легко показать, что с помощью разрезов графов можно оптимизировать именно субмодулярную энергию

Доказательство необходимости

Разрезы графов

Ветров

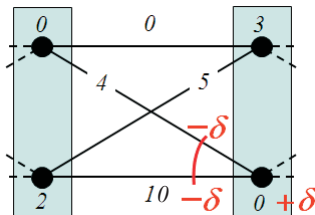
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Пусть имеется парное слагаемое, не являющееся субмодулярной функцией

$$E_{ij}(0, 0) + E_{ij}(1, 1) > E_{ij}(0, 1) + E_{ij}(1, 0)$$

- Тогда и только тогда в результате репараметризации нулевые веса горизонтальных связей приведут к возникновению отрицательных диагональных связей

Доказательство необходимости

Разрезы графов

Ветров

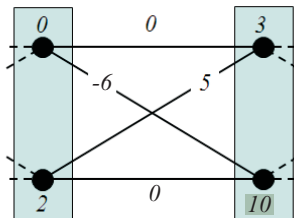
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Но это означает, что пропускная способность некоторых ребер в графе, разрез которого мы будем минимизировать, станет отрицательной!
- В этой ситуации классические полиномиальные алгоритмы поиска минимального разреза в графе неприменимы
- Задача оптимизации энергии стала NP-трудной
- Существуют некоторые обобщения полиномиального алгоритма на случаи, когда энергия может быть сведена в субмодулярной путем замены части переменных на свои отрицания: $t_i \rightarrow (1 - t_i)$

Случай небинарных переменных

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- До сих пор рассматривался случай, когда скрытые переменные бинарные
- Теперь рассмотрим ситуацию, когда скрытые переменные t_i могут принимать одно из K значений
- Физически это соответствует делению изображения на K областей
- Такие задачи возникают при построении карт диспаратетов, коллажах, семантической сегментации и пр.

Семантическая сегментация

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



Пример задачи семантической сегментации

Анатомическая разметка

Разрезы графов

Ветров

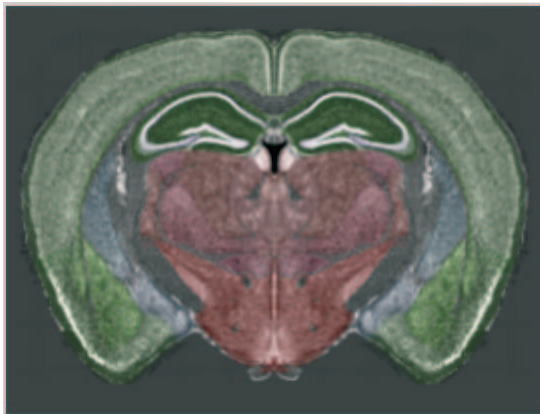
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



Пример задачи автоматического выделения анатомических зон головного мозга мыши

Итерационная схема

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

$$E(T|X) = \sum_{(i,j) \in E} E_{ij}(t_i, t_j) + \sum_i E_i(x_i, t_i) \rightarrow \min_T, \quad t_i \in \{0, 1, \dots, K-1\}$$

- Задача оптимизации энергии по K -значным скрытым переменным ($K > 2$) является NP-трудной
- Тем не менее, в ряде случаев можно построить итерационную процедуру, сходящуюся к близкому к глобальному оптимуму ответу
- Наибольшее распространение получил алгоритм т.н. α -расширения

Итерационная схема

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Начинаем с произвольного начального приближения
- В цикле для каждой метки $\alpha \in \{0, \dots, K - 1\}$ заменяем часть других меток на данную, так чтобы минимизировать энергию (выполняем α -расширение)
- Если хотя бы для одной метки энергию удалось уменьшить, то переходим в предыдущему шагу, иначе выход

Для сходимости такого алгоритма необходимо, чтобы каждое парное слагаемое было метрикой в пространстве $\{0, \dots, K - 1\}$, т.е. удовлетворяло следующим условиям $\forall \alpha, \beta, \gamma \in \{0, \dots, K - 1\}$

- $E_{ij}(\alpha, \beta) = E_{ij}(\beta, \alpha)$ (симметричность)
- $E_{ij}(\alpha, \beta) = 0 \Leftrightarrow \alpha = \beta$ (аксиома тождества)
- $E_{ij}(\alpha, \gamma) \leq E_{ij}(\alpha, \beta) + E_{ij}(\beta, \gamma)$ (неравенство треугольника)

α -расширение

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

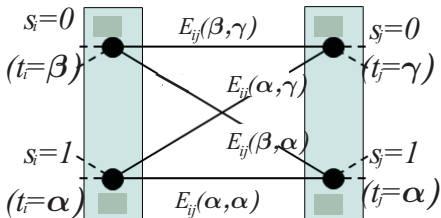
- В ходе альфа-расширения часть других меток принимает значение α , стремясь минимизировать энергию
- Введем вспомогательную марковскую решетку с бинарными переменными s_j , в которой значение 0 будет соответствовать тому, что исходные скрытые переменные t_j не изменились, а значение 1 будет означать, что исходные переменные приняли значение α

$$\forall j : t_j^{old} = \alpha \Rightarrow s_j \equiv 0$$

$$\forall j : t_j^{old} \neq \alpha \Rightarrow \begin{cases} s_j = 0 \Rightarrow t_j^{new} = t_j^{old} \\ s_j = 1 \Rightarrow t_j^{new} = \alpha \end{cases}$$

- Теперь относительно новых переменных можно построить минимальный разрез графа, минимизирующий энергию по всевозможным α -расширениям

α -расширение



- Рассмотрим некоторую пару скрытых переменных (t_i, t_j) , соединенную ребром
- Предположим, что старые значения переменных равнялись β и γ соответственно
- Для корректной репараметризации необходимо выполнение неравенства треугольника

$$E_{ij}(\beta, \gamma) \leq E_{ij}(\alpha, \gamma) + E_{ij}(\beta, \alpha)$$

- Теперь относительно новых переменных можно построить минимальный разрез графа, минимизирующий энергию по всевозможным α -расширениям

α -расширение

Разрезы графов

Ветров

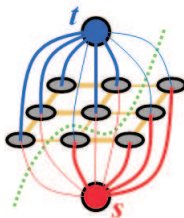
Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение



- Запретим изменять метки класса α

$$c(S, t_i) = E_i(x_i, \alpha), \quad c(T, t_i) = +\infty$$

$$c(t_i, t_j) = E_{ij}(\alpha, \alpha) = 0, \quad \forall (i, j) \in E : t_i = t_j = \alpha$$

- Для остальных вершин ($t_i^{old} \neq \alpha$) определим пропускную способность следующим образом:

$$c(S, t_i) = E_i(x_i, \alpha), \quad c(T, t_i) = E_i(x_i, t_i^{old})$$

$$c(t_i, t_j) = E(\alpha, t_i^{old}), \quad c(t_j, t_i) = E(\alpha, t_j^{old}), \quad \forall (i, j) \in E$$

- Вершинам, попавшим в тот же подграф, что и исток S , будет присвоена метка α , энергия при этом уменьшится

Точность получающегося решения

Разрезы графов

Ветров

Ликбез

Марковские сети

Разрезы графов

Субмодулярные функции

Альфа-расширение

- Алгоритм α -расширения является итерационным, полиномиальным, поэтому не гарантирует достижение глобального оптимума (NP-трудная задача)
- Можно показать, что значение энергии, получившейся в результате альфа-расширения, лежит в интервале

$$E(T^*) \leq E(T) \leq 2kE(T^*),$$

где T^* — оптимальное (наиболее вероятное) значение скрытых переменных, а

$$k = \frac{\max E_{ij}(\beta, \gamma)}{\min_{\beta \neq \gamma} E_{ij}(\beta, \gamma)}$$

степень контрастности парной энергии

Московский государственный университет
имени М. В. Ломоносова
Факультет вычислительной математики и кибернетики

Математические основы
теории прогнозирования
(курс лекций)

лектор — академик РАН Ю. И. Журавлев

Оглавление

1		3
1.1	Стандартная задача распознавания	3
1.2	Алгоритм “Кора” (Вайнцвайг, Бонгарт)	4
1.3	Тестовый алгоритм (Ю. И. Журавлев)	6
2		9
2.1	Логические алгоритмы распознавания	9
3		16
3.1	Алгоритмы вычисления оценок	16
3.2	Эффективные формулы вычисления оценок	20
4		24
4.1	Вычисление характеристик, определяющих алгоритм вычисления оценок . .	24
4.2	Алгебры над алгоритмами	26
5		28
5.1	Построение алгоритмов распознавания, корректных для заданной контрольной выборки	28

Курс лекций, прочитанный для 3 потока IV курса факультета ВМиК, набран в системе L^AT_EX студентами:

1 лекция — П. Клеменков

2 лекция — А. Гудков

3 лекция — К. Симонян

4 и 5 лекции — А. Фокин

Лекция 1

1.1 Стандартная задача распознавания

Пусть дано множество M , являющееся суммой подмножеств K_1, \dots, K_l , называемых обычно классами.

$$M = \bigcup_{j=1}^l K_j$$

Различают случаи а) $K_u \cap K_v = \emptyset$, б) $K_u \cap K_v$, вообще говоря, не пусто. В случае а) говорят о задаче с пересекающимися, в случае б) — непересекающимися классами (множества $K_j, j = 1 \dots l$ принято называть классами).

В дальнейшем рассматриваются только M специального вида: элементы M являются наборами длины n : $\tilde{a} \in M, \tilde{a} = (a_1, a_2, \dots, a_i, \dots, a_n)$. При этом для каждого номера $i, i = 1 \dots n$, определено множество допустимых значений M_i , являющееся метрическим пространством с метрикой ρ_i , т.е. выполнены аксиомы: $\rho_i(c, d) \geq 0, \rho_i(c, c) = 0, \rho_i(c, d) = \rho_i(d, c), \rho_i(c, e) + \rho_i(e, d) \geq \rho_i(c, d), c, d, e \in M_i$. В некоторых случаях выполнения последней аксиомы (аксиомы треугольника) не требуют. Тогда говорят, что в M_i введена полуметрика.

В качестве исходной информации задаются некоторые сведения о множестве M и классах K_1, \dots, K_l .

В дальнейшем в качестве исходной информации рассматривается так называемая стандартная обучающая информация I : выделяется конечное множество S_1, \dots, S_m элементов из M : $S_i = (a_{i1}, a_{i2}, \dots, a_{it}, \dots, a_{in}), i = 1 \dots m, a_{it} \in M_t$, для которых известно, в какие из K_1, \dots, K_l они входят. Последнее оформляется заданием информационного вектора $\tilde{\alpha}(S_i) = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{ij}, \dots, \alpha_{in}), (\alpha_{ij} = 1) \rightarrow S_i \in K_j, (\alpha_{ij} = 0) \rightarrow S_i \notin K_j, i = 1 \dots m, j = 1 \dots l$.

Для удобства данные об элементах S_i и их информационных векторах представляют в виде таблиц:

	1	2	...	t	...	n
S_1	a_{11}	a_{12}	...	a_{1t}	...	a_{1n}
S_2	a_{21}	a_{22}	...	a_{2t}	...	a_{2n}
...
S_i	a_{i1}	a_{i2}	...	a_{it}	...	a_{in}
...
S_m	a_{m1}	a_{m2}	...	a_{mt}	...	a_{mn}
$\underbrace{\hspace{10em}}_{T_1}$						

K_1	K_2	\dots	K_j	\dots	K_l	
α_{11}	α_{12}	\dots	α_{1j}	\dots	α_{1l}	$\tilde{\alpha}(S_1)$
α_{21}	α_{22}	\dots	α_{2j}	\dots	α_{2l}	$\tilde{\alpha}(S_2)$
\dots	\dots	\dots	\dots	\dots	\dots	\dots
α_{i1}	α_{i2}	\dots	α_{ij}	\dots	α_{il}	$\tilde{\alpha}(S_i)$
\dots	\dots	\dots	\dots	\dots	\dots	\dots
α_{m1}	α_{m2}	\dots	α_{mj}	\dots	α_{ml}	$\tilde{\alpha}(S_m)$

T_2

Совокупность таблиц T_1, T_2 называется стандартной обучающей информацией I , таблица T_2 называется информационной матрицей.

Стандартная задача распознавания: пусть задан элемент $S \in M, S \notin \{S_1, S_2, \dots, S_i, \dots, S_m\}$. Найти алгоритм A , который, используя только I и представление S строит информационный вектор $\tilde{\alpha}(S) = (\tilde{\alpha}_1(S), \dots, \tilde{\alpha}_j(S), \dots, \tilde{\alpha}_l(S))$.

$$A(I, S) = \tilde{\alpha}(S).$$

В задачах распознавания часто обучающая информация I оказывается недостаточной для построения “правильного” вектора $\tilde{\alpha}(S)$. Поэтому допускаются и широко используются эвристические алгоритмы, допускающие ошибки и отказы при вычислении координат информационных векторов.

Такие алгоритмы $\tilde{A}(I, S)$ строят квази-информационные векторы $\tilde{\beta}(S) = (\beta_1(S), \dots, \beta_j(S), \dots, \beta_l(S))$. При этом возможно, что: $\beta_j(S) \neq \tilde{\alpha}_j(S)$ (ошибка в распознавании), $\beta_j(S) = \Delta$ — так кодируется отказ от вычисления j -й координаты информационного вектора.

В литературе описано значительное число таких эвристических алгоритмов, допускающих небольшое (допустимое при практическом применении) число ошибок и отказов при решении достаточно узких классов реальных прикладных задач. Мы опишем два таких алгоритма, получивших большое распространение при прогнозировании в геологии, медицине, технике и т.п.

В дальнейшем координаты $1, 2, \dots, n$, задающие n -мерные объекты в M , будем называть признаками.

1.2 Алгоритм “Кора” (Вайнцвайг, Бонгарт)

Применяется для M , элементами которых являются бинарные признаки: $M_i = \{0, 1\}$, $i = 1 \dots n$, в основном для задач с двумя непересекающимися классами: $M = K_1 \cup K_2$, $K_1 \cap K_2 = \emptyset$.

В таблице $\|a_{ij}\|_{m \times n}$, задающей объекты с известной классовой принадлежностью, пусть S_1, \dots, S_q принадлежат K_1 , S_{q+1}, \dots, S_m принадлежат K_2 . Просматриваем все тройки признаков r, u, v (число таких троек, очевидно, равно $\binom{n}{3}$) и анализируем часть таблицы T_1 ,

составленной только из столбцов r, u, v :

a_{1r}	a_{1u}	a_{1v}
a_{2r}	a_{2u}	a_{2v}
\dots	\dots	\dots
a_{ir}	a_{iu}	a_{iv}
\dots	\dots	\dots
a_{qr}	a_{qu}	a_{qv}
a_{q+1r}	a_{q+1u}	a_{q+1v}
a_{q+2r}	a_{q+2u}	a_{q+2v}
\dots	\dots	\dots
a_{jr}	a_{ju}	a_{jv}
\dots	\dots	\dots
a_{mr}	a_{mu}	a_{mv}

Среди первых q строк выделяем и фиксируем все тройки, не совпадающие ни с одной из троек в строках $q + 1, \dots, m$. Формируем множество таких троек $\{(a_{ir}, a_{iu}, a_{iv})\}$. Аналогично выделяем все тройки (a_{jr}, a_{ju}, a_{jv}) , не совпадающие ни с одной из первых q троек. Множества $\{(a_{ir}, a_{iu}, a_{iv})\}, \{(a_{jr}, a_{ju}, a_{jv})\}$ назовем, соответственно, характеристиками классов K_1, K_2 . Такие характеристики формируем для всех троек (r, u, v) .

Пусть задан для распознавания объект $S = (b_1 \dots b_r \dots b_u \dots b_v \dots b_n)$. Сравниваем все характеристики всех троек для K_1 с соответствующими тройками в распознаваемом объекте S . Число совпадений $(a_{ir}, a_{iu}, a_{iv}) = (b_r, b_u, b_v)$ обозначаем $\Gamma(S, K_1)$ — число голосов, поданных для S за класс K_1 . Аналогично формируем величину $\Gamma(S, K_2)$: число совпадений $(a_{jr}, a_{ju}, a_{jv}) = (b_r, b_u, b_v)$. Вводим пороговый параметр ν .

Если $\Gamma(S, K_1) - \nu > \Gamma(S, K_2)$, относим S классу K_1 , при $\Gamma(S, K_2) - \nu > \Gamma(S, K_1)$ — в класс K_2 . В остальных случаях алгоритм отказывается от классификации. На практике часто полагают $\nu = 0$.

Пример 1

Дана таблица T_1

	1	2	3	4	5
S_1	1	0	1	0	0
S_2	0	1	0	1	0
S_3	0	0	1	0	1
S_4	1	0	0	1	0
S_5	1	0	0	0	1
S_6	0	1	0	0	1

Имеем $\binom{5}{3} = 10$ троек признаков. Перечислим характеристики для K_1 и K_2 .
Характеристики для K_1 :

- | | |
|-----------------------------------|-----------------------------------|
| 1.(1, 2, 3) : (101), (001) | 2.(1, 2, 4) : (011), (000) |
| 3.(1, 2, 5) : (010), (001) | 4.(1, 3, 4) : (110), (001), (010) |
| 5.(1, 3, 5) : (100), (000), (011) | 6.(1, 4, 5) : (100), (010) |
| 7.(2, 3, 4) : (010), (101) | 8.(2, 3, 5) : (010), (100), (011) |
| 9.(2, 4, 5) : (000), (110) | 10.(3, 4, 5) : (100), (101) |

Характеристики для K_2

1.(1, 2, 3) : (100)	2.(1, 2, 4) : (101), (010)
3.(1, 2, 5) : (101), (011)	4.(1, 3, 4) : (101), (100), (000)
5.(1, 3, 5) : (100), (101), (001)	6.(1, 4, 5) : (110), (101)
7.(2, 3, 4) : (001), (000), (100)	8.(2, 3, 5) : (000), (001), (101)
9.(2, 4, 5) : (010), (101)	10.(3, 4, 5) : (011)

Очевидно, что с увеличением n (числа признаков) число троек в характеристиках растет весьма быстро. Поэтому при реальных решениях обязательно использование компьютеров.

Заметим, что для объекта $S = (00000) : \Gamma(S, K_1) = \Gamma(S, K_2) = 3$. Поэтому алгоритм не классифицирует этот объект.

При распознавании объекта $S = (10101)$ имеют место совпадения с элементами характеристики K_1 : (1, 2, 3), (101); (1, 3, 4), (110); (2, 3, 4), (010); (2, 3, 5), (011); (3, 4, 5), (101). Следовательно, $\Gamma(S, K_1) = 5$. Легко проверить, что $\Gamma(S, K_2) = 2$. Алгоритм “Кора” заносит S в класс K_1 .

1.3 Тестовый алгоритм (Ю. И. Журавлев)

Пусть задана бинарная таблица $\|a_{ij}\|_{m \times n}$, строки которой S_1, \dots, S_m разделены на два класса, причем S_1, \dots, S_q — строки первого класса K_1 , S_{q+1}, \dots, S_m — строки второго класса K_2 .

Набор столбцов с номерами k_1, \dots, k_l образует тест, если после удаления из таблицы всех столбцов за исключением вышеозначенных, ни одна из строк из K_1 не совпадет ни с одной из строк класса K_2 . Тест называется тупиковым, если при удалении из него хотя бы одного столбца, хотя бы одна из строк K_1 совпадет хотя бы с одной из строк K_2 .

Предположим, что построены все тупиковые тесты T_1, \dots, T_r бинарной таблицы, $T_i = \{n_{i1}, \dots, n_{ip(i)}\}$, $i = 1 \dots r$, n_{uv} — номера столбцов, входящих в тест.

Распознаваемый объект $S = (a_1 \dots a_n)$ последовательно совмещается с тупиковыми тестами. При работе с тестом T_i набор $a_{n_{i1}}, \dots, a_{n_{ip(i)}}$ сравнивается по столбцам теста со всеми строками таблицы $\|a_{ij}\|_{m \times n}$. При этом возможно совпадение со строкой не более чем в одном из классов K_1, K_2 (это следует из определения теста). Число совпадений суммируется отдельно для классов K_1, K_2 . Полученные суммы $\Gamma(S, K_1), \Gamma(S, K_2)$ используются для классификации объекта S так же, как в алгоритме “Кора”.

Понятия “тест” и “тупиковый тест” нетрудно распространить для таблиц, составленных из элементов произвольной природы. Необходимо только, чтобы элементы каждого столбца содержались в метрическом пространстве. Обозначим метрику этого пространства через ρ_t . Два элемента a_{it}, a_{jt} назовем различимыми, если $\rho_t(a_{it}, a_{jt}) > \epsilon_t$; в противном случае a_{it}, a_{jt} неразличимы. В определении теста достаточно заменить слова “равны” и “не равны” на “неразличимы” и “различимы”. Значение ϵ_t задается из “содержательных” соображений или определяется при решении другой задачи.

Построение совокупности тупиковых тестов связано с решением системы базовых уравнений.

Пусть дана система

$$f_i(x_1 \dots x_n) = 1, i = 1 \dots k. \quad (1.1)$$

Система (1.1) эквивалентна одному уравнению

$$\prod_{i=1}^k f_i(x_1 \dots x_n) = 1. \quad (1.2)$$

Представим f_i в виде дизъюнктивной нормальной формы (д.н.ф.)

$$f_i = \mathfrak{D}_i = \bigvee_i \mathcal{K}_{it(i)},$$

где $\mathcal{K}_{it(i)}$ — элементарные конъюнкции, т.е. произведения вида $x_{j_1}^{\sigma_1} \cdot \dots \cdot x_{j_p}^{\sigma_p}$, $x^\sigma = x$ при $\sigma = 1, \bar{x}$ при $\sigma = 0$.

Выполним в (1.2) операции логического умножения и получим финальную д.н.ф.

$$\bigvee_{i=1}^r Q_i = 1. \quad Q_i = x_{i_1}^{\delta_1} \cdot \dots \cdot x_{i_p}^{\delta_p}.$$

Последовательно решаем уравнения $Q_i = 1$; $x_{i_1} = \delta_1, \dots, x_{i_p} = \delta_p$, остальные $x_j = 0, 1$. Совокупность всех решений и есть совокупность всех решений системы (1.1).

Выведем систему уравнений для построения всех тупиковых тестов таблицы $\|a_{ij}\|_{m \times n}$, в которой строки S_1, \dots, S_q принадлежат K_1 , а строки S_{q+1}, \dots, S_m — классу K_2 .

Сопоставим столбцам $1, 2, \dots, n$ булевы переменные x_1, \dots, x_n . Напишем систему из $q \cdot (m - q)$ булевых уравнений. Паре $S_i = (a_{i1}, \dots, a_{in}) \in K_1$, $S_j = (a_{j1}, \dots, a_{jn}) \in K_2$ сопоставим уравнение

$$f_{ij} = \bigvee_{a_{it} \neq a_{jt}} x_t, \quad i = 1, \dots, q; \quad j = q + 1, \dots, m$$

$$\prod_{\substack{i=1, \dots, q \\ j=q+1, \dots, m}} f_{ij} = 1.$$

Выполняем умножение и приходим к д.н.ф. $\bigvee x_{u_1} \cdot \dots \cdot x_{u_p}$, в которой проводим все упрощения $\mathcal{K} \vee \mathcal{K} \cdot \bar{\mathcal{K}} = \mathcal{K}$.

Финальное уравнение

$$\bigvee x_{l_1} \cdot \dots \cdot x_{l_v}$$

определяет все тупиковые тесты (l_1, \dots, l_v) .

Обоснование алгоритма см. в лекции 4.

Процесс умножения при переходе к финальному уравнению и реализация функций в классе д.н.ф. весьма трудоемки. В тестовом алгоритме последнее отсутствует, т.к. уравнения сразу задаются в виде д.н.ф. Процесс умножения можно существенно упростить, используя специфику булевой алгебры. Укажем несколько упрощающих приемов.

- а) из двух уравнений $f = 1$, $f \cdot \tilde{f} = 1$. Второе можно удалить, т.к. оно является следствием первого.
- б) пусть даны уравнения $f_0 \vee f_i = 1$, $i = 1 \dots k$; тогда $\prod_{i=1}^k (f_0 \vee f_i) = f_0 \vee f_1 \cdot \dots \cdot f_k$. Действительно: $f_0 \cdot f_0 = f_0$, $f_0 \vee f_0 \cdot f_i = f_0$. (правило поглощения)
- в) $\mathcal{K} \cdot (\mathcal{K} \cdot \mathcal{K}') = \mathcal{K} \cdot \mathcal{K}'$, $Q \vee Q = Q$.

Существует большое число других упрощающих правил. Эффективность трех приведенных выше продемонстрируем на примере, рассмотренном в алгоритме “Кора”.

	1	2	3	4	5	
S_1	1	0	1	0	0	$S_1, S_2, S_3 \in K_1; \quad S_4, S_5, S_6 \in K_2$
S_2	0	1	0	1	0	
S_3	0	0	1	0	1	
S_4	1	0	0	1	0	
S_5	1	0	0	0	1	
S_6	0	1	0	0	1	

Имеем 9 уравнений, получаемых при сравнении строк $S_i, i = 1, 2, 3$ со строками $S_j, j = 4, 5, 6$.

$$\begin{aligned} (S_1, S_4) : x_3 \vee x_4 = 1; & \quad (S_1, S_5) : x_3 \vee x_5 = 1; & \quad (S_1, S_6) : x_1 \vee x_2 \vee x_3 \vee x_5 = 1 \\ (S_2, S_4) : x_1 \vee x_2 = 1; & \quad (S_2, S_6) : x_4 \vee x_5 = 1; & \quad (S_2, S_5) : x_1 \vee x_2 \vee x_4 \vee x_5 = 1 \\ (S_3, S_5) : x_1 \vee x_3 = 1; & \quad (S_3, S_6) : x_2 \vee x_3 = 1; & \quad (S_3, S_4) : x_1 \vee x_3 \vee x_4 \vee x_5 = 1 \end{aligned}$$

По правилу а) удаляются уравнения $(S_1, S_6), (S_2, S_5), (S_3, S_4)$. Среди оставшихся выделим уравнения $x_3 \vee x_4 = 1, x_3 \vee x_5 = 1, x_1 \vee x_3 = 1, x_2 \vee x_3 = 1$. По правилу б) произведение левых частей даст $x_3 \vee x_1 \cdot x_2 \cdot x_4 \cdot x_5 = 1$. Перемножаем оставшиеся два уравнения

$$(x_1 \vee x_2) \cdot (x_4 \vee x_5) = x_1x_4 \vee x_1x_5 \vee x_2x_4 \vee x_2x_5$$

Перемножая левые части двух полученных уравнений и, используя в), имеем:

$$x_1x_3x_4 \vee x_1x_3x_5 \vee x_2x_3x_4 \vee x_2x_3x_5 \vee x_1x_2x_4x_5 = 1.$$

К последней д.н.ф. правило поглощения неприменимо, поэтому наборы

$$(1, 3, 4), (1, 3, 5), (2, 3, 4), (2, 3, 5), (1, 2, 4, 5) \quad \text{образуют все тупиковые тесты.}$$

Классифицируем, как в алгоритме “Кора” набор $S = (10101)$. По тесту $(1, 3, 4)$ имеем совпадение с $S_1 \in K_1$; по $(1, 3, 5)$ нет совпадений ни с одной строкой; по $(2, 3, 4)$ — совпадение с $S_1, S_3 \in K_1$; по $(2, 3, 5)$ — с $S_3 \in K_1$; по $(1, 2, 4, 5)$ — с $S_5 \in K_2$. Следовательно

$$\Gamma(S, K_1) = 3, \quad \Gamma(S, K_2) = 1.$$

Вывод: при $\nu < 2 : S \in K_1$, при $\nu \geq 2$ алгоритм откажется от распознавания.

Лекция 2

2.1 Логические алгоритмы распознавания

Для избежания громоздких выкладок и привлечения теории функций k -значной логики ограничимся задачей распознавания с двумя непересекающимися классами K_1, K_2 , причем признаки будут принимать только значения 0,1.

В дальнейшем объекты исходной информации I задаются бинарными наборами $S_1, \dots, S_r, S_{r+1}, \dots, S_m$, где $S_i = (\alpha_{i1} \dots \alpha_{ik} \dots \alpha_{in})$, $i = 1 \dots n$. Объекты S_i , $i = 1 \dots r$ принадлежат K_1 , объекты S_i , $i = r + 1, \dots, m$ — классу K_2 .

Напомним некоторые сведения из теории булевых функций (функций алгебры логики).

Каждая $f(x_1, \dots, x_n)$, вообще говоря, неоднозначно представима дизъюнктивной нормальной формой д.н.ф. $\bigvee_i \mathcal{K}_i$, где \mathcal{K}_i — элементарные конъюнкции. Если $\mathcal{K}_i = x_1^{\sigma_1}, \dots, x_k^{\sigma_k}$, то k — ранг конъюнкции,

$$x^\sigma = \begin{cases} x, & \text{если } \sigma = 1; \\ \bar{x}, & \text{если } \sigma = 0. \end{cases}$$

Если \mathcal{N}_f — множество единиц f , $\mathcal{N}_{\mathcal{K}}$ — интервал, \mathcal{K} — множество единиц конъюнкции \mathcal{K} , то $f = \bigvee_{i=1}^t \mathcal{K}_i \Leftrightarrow \mathcal{N}_{\mathcal{K}_i} = \bigcup_{i=1}^t \mathcal{N}_{\mathcal{K}_i}$. Интервал называется максимальным, а соответствующая ему элементарная конъюнкция — простой импликантой, если не существует $\mathcal{N}_{\mathcal{K}'} : \mathcal{K} \subset \subset \mathcal{N}_{\mathcal{K}'} \subset \mathcal{N}_f$.

Пусть $\mathcal{N}_{\mathcal{K}_1}, \dots, \mathcal{N}_{\mathcal{K}_c}$ — совокупность всех максимальных интервалов функции f . Д.н.ф. $D_c(f)$ называют сокращенной д.н.ф. функции f . Каждая д.н.ф. минимальной сложности получается удалением из $D_c(f)$ некоторых э.к. ($D_c(f) = \bigvee_{i=1}^l \mathcal{K}_i$).

Напомним, что сложностью д.н.ф. называется сумма рангов входящих в нее э.к.

Построение минимальных д.н.ф. подразделяется на следующие этапы:

- I. строится произвольная д.н.ф. D_f , реализующая f
- II. к D_f применяются преобразования $x_i \mathcal{K}_u \vee \bar{x}_i \mathcal{K}_v \rightarrow x_i \mathcal{K}_u \vee \bar{x}_i \mathcal{K}_v \vee \mathcal{K}_u \mathcal{K}_v$; $\mathcal{K} \vee \mathcal{K} \mathcal{K}' \rightarrow \mathcal{K}$ до тех пор, пока это возможно. Построенная д.н.ф. называется сокращенной.
- III. в д.н.ф. $D_c(f) = \bigvee_{i=1}^l \mathcal{K}_i$ выбирается произвольный интервал $\mathcal{N}_{\mathcal{K}_j}$, $1 \leq j \leq l$, такой, что $\mathcal{N}_{\mathcal{K}_j} \subset \bigcup_{i \neq j} \mathcal{N}_{\mathcal{K}_i}$. Э.к. \mathcal{K}_j удаляются из $D_c(f)$, оставшиеся интервалы образуют покрытие \mathcal{N}_f ; поэтому $D_c(f) \setminus \mathcal{K}_j$ реализует f . Процесс повторяется до тех пор, пока в покрытии не остаются только интервалы, не покрываемые суммой остальных.

ных $\mathcal{N}_{\mathcal{K}_{U_1}}, \dots, \mathcal{N}_{\mathcal{K}_{U_p}}$. Соответствующая д.н.ф называется тупиковой для $f : D_1(f) = \bigvee_{i=1}^p K_{U_i}$.

Процесс удаления конъюнкций (их интервалов из покрытия) не однозначен, поэтому число тупиковых д.н.ф. может быть велико. Очевидно, что среди тупиковых содержатся и все минимальные д.н.ф.

Для дальнейшего нам понадобятся несколько утверждений и алгоритмов:

- I. Как относительно нетрудоемко построить д.н.ф. приемлемой сложности, реализующую f
- II. Найти аналитический критерий, позволяющий легко проверить:

$$\mathcal{N}_{\mathcal{K}} \subseteq \bigcup_{i=1}^q \mathcal{N}_{\mathcal{K}_i},$$

что необходимо для построения тупиковых д.н.ф.

- III. Как влияют на соотношения $\mathcal{N}_{\mathcal{K}} \subseteq \bigcup \mathcal{N}_{\mathcal{K}_i}$, $\mathcal{N}_{\mathcal{K}} \not\subseteq \bigcup \mathcal{N}_{\mathcal{K}_i}$ преобразования $x_i \rightarrow x_j^{\sigma_{ij}}$, $\binom{i}{j}$ — подстановка (преобразование взаимно однозначно), $\sigma_{ij} \in \{0, 1\}$

Проверка соотношения $\mathcal{N}_{\mathcal{K}} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$. Не ограничивая общности считаем, что в \mathcal{K} и \mathcal{K}_i , $i = 1 \dots l$ нет переменных x_t в различных степенях. Т.е. если $x_t^\sigma \in \mathcal{K}$, то в \mathcal{K}_i , $i = 1 \dots l$ нет сомножителей $x_t^{\bar{\sigma}}$. Действительно, если бы $x_t^{\bar{\sigma}}$ находился в \mathcal{K}_i , то $\mathcal{K}\mathcal{K}_i \equiv 0$, ($\mathcal{N}_{\mathcal{K}} \cap \mathcal{N}_{\mathcal{K}_i} = \emptyset$) и $\mathcal{N}_{\mathcal{K}_i}$ не влиял бы на выполнимость проверяемого соотношения.

Представим \mathcal{K}_i в виде $\mathcal{K}_i^1 \cdot \mathcal{K}_i^2$; в \mathcal{K}_i^1 входят все сомножители, общие с \mathcal{K} , в \mathcal{K}_i^2 — оставшиеся. Если оставшихся нет, полагаем $\mathcal{K}_i^2 = 1$.

Теорема 1 (Критерий поглощения) $\mathcal{N}_{\mathcal{K}} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$ тогда и только тогда, когда $\bigvee_{i=1}^l \mathcal{K}_i^2 \equiv 1$.

Доказательство. Достаточность. Пусть $\tilde{\alpha} \in \mathcal{N}_{\mathcal{K}} : \mathcal{K}(\tilde{\alpha}) = 1$. Но тогда очевидно $\mathcal{K}_i^1 = 1$, $i = 1 \dots l$. Выделим в $\tilde{\alpha}$ поднабор из координат, соответствующих переменным из \mathcal{K}_i^2 , $i = 1 \dots l$. Так как $\bigvee_{i=1}^l \mathcal{K}_i^2 = 1$, то найдется \mathcal{K}_u^2 , $1 \leq u \leq l$, равное 1 на этом наборе. Но

в этом случае $\mathcal{K}_u^1 \cdot \mathcal{K}_u^2(\tilde{\alpha}) = 1$, $\mathcal{K}_u(\tilde{\alpha}) = 1$. Следовательно $\tilde{\alpha} \in \mathcal{N}_{\mathcal{K}_u} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$. Достаточность доказана.

Необходимость. Пусть $\bigvee_{i=1}^l \mathcal{K}_i^2 \not\equiv 1$. Тогда найдется поднабор $\tilde{\beta}$ координат переменных, не входящих в \mathcal{K} , такой, что $\mathcal{K}_i^2(\tilde{\beta}) = 0$, $i = 1 \dots l$.

Пусть $\mathcal{K} = X_{t_1}^{\sigma_1}, \dots, X_{t_v}^{\sigma_v}$. Сформулируем набор $\tilde{\gamma}$, положив в нем координаты t_1, \dots, t_v , равные, соответственно $\sigma_1, \dots, \sigma_v$, добавим значения координат из набора $\tilde{\beta}$; остальные координаты зададим произвольно. Тогда $\mathcal{K}(\tilde{\gamma}) = 1$, $\mathcal{K}_i^1(\tilde{\gamma}) \cdot \mathcal{K}_i^2(\tilde{\gamma}) = 0$, $i = 1 \dots l$, так $\mathcal{K}_i^2(\tilde{\gamma}) = 0$, $i = 1 \dots l$ ($\mathcal{K}_i^2(\tilde{\beta}) = 0$, а $\tilde{\beta}$ — часть набора $\tilde{\gamma}$). Имеем: $\tilde{\gamma} \in \mathcal{N}_{\mathcal{K}}$, $\tilde{\gamma} \notin \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$. Необходимость доказана.

Пусть π — преобразование $x_i \rightarrow y_j^{\sigma_{ij}}$, $\binom{i}{j}$ — подстановка, $\pi(\mathcal{K})$ — результат преобразования \mathcal{K} с помощью π .

Теорема 2 $\mathcal{N}_{\mathcal{K}} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i} \leftrightarrow \mathcal{N}_{\pi(\mathcal{K})} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\pi(\mathcal{K}_i)}$

Доказательство. Пусть $\mathcal{N}_{\mathcal{K}} \subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$. Тогда $\bigvee_{i=1}^l \mathcal{K}_i^2 \equiv 1$ (по т. 1). Но любая подстановка в функцию $f(z_1, \dots, z_m)$, $z_i = \varphi_i(y_{i_1}, \dots, y_{i_{k_i}})$ приводит к $f(\varphi_1, \dots, \varphi_m) \equiv 1$, если $f(z_1, \dots, z_m) \equiv 1$. Проверка последнего тривиальна. Если $\mathcal{N}_{\mathcal{K}} \not\subseteq \bigcup_{i=1}^l \mathcal{N}_{\mathcal{K}_i}$, то $\bigvee_{i=1}^l \mathcal{K}_i^2 \not\equiv 1$ (т. 1), и существует набор $\tilde{\beta}$: $\mathcal{K}_i^2(\tilde{\beta}) = 0$, $i = 1 \dots l$. Это значит, что в каждой \mathcal{K}_i^2 есть сомножитель x_r^σ , а r -я координата в $\tilde{\beta}$ равна $\bar{\sigma}$. Если при $\pi : x_r \rightarrow y_t$, то в новом наборе t -я координата равна $\bar{\sigma}$, а соответствующий сомножитель: y_t^σ . Очевидно, $\pi(\mathcal{K}_i^2(\tilde{\beta})) = 0$. Случай $x_r \rightarrow \bar{y}_t$ разбирается аналогично.

Сказанное выше применимо ко всем \mathcal{K}_i^2 , имеющим сомножитель x_r^σ . Остальные \mathcal{K}_u^2 либо не имеют сомножителя от переменной x_r , либо имеют сомножитель x_r^σ . Следовательно, в каждой из этих \mathcal{K}_u^2 найдется сомножитель от x_q , $q \neq r$, для которого проходят предыдущие выкладки. Таким образом $\pi(\bigvee_{i=1}^l \mathcal{K}_i^2) \not\equiv 1$ и $\mathcal{N}_{\pi(\mathcal{K})} \not\subseteq \bigcup_{i=1}^l \mathcal{N}_{\pi(\mathcal{K}_i)}$.

Переходим к реализации I. Докажем сначала:

$$\begin{aligned} & (x_1 \vee \dots \vee x_n) \cdot (\bar{x}_1 \vee \dots \vee \bar{x}_n) = \\ & = x_1 \cdot x_2 \vee x_2 \cdot x_3 \vee \dots \vee x_i \cdot \bar{x}_{i+1} \vee x_{i+1} \cdot \bar{x}_{i+2} \vee \dots \vee x_{n-1} \cdot \bar{x}_n \vee x_n \cdot \bar{x}_1 \end{aligned}$$

Легко видеть, что левая часть реализует функцию, равную на наборах $(0 \dots 0 \dots 0)$, $(1 \dots 1 \dots 1)$. Действительно, на любом другом наборе либо найдется пара координат i , $i+1$, таких, что $\alpha_i = 1$, $\alpha_{i+1} = 0$ и тогда $x_i \cdot \bar{x}_{i+1} = 1$, либо $\alpha_n = 1$, $\alpha_1 = 0$. И тогда $x_n \cdot \bar{x}_1 = 1$ (рассматривается набор $\tilde{\alpha} = (\alpha_1 \dots \alpha_n)$).

Заметим, что умножение двух конъюнкций $(x_1^{\sigma_1} \vee \dots \vee x_n^{\sigma_n}) \cdot (x_1^{\delta_1} \vee \dots \vee x_n^{\delta_n})$ — произведение реализует функцию, равную 0 только на наборах $(\bar{\sigma}_1 \dots \bar{\sigma}_n)$, $(\bar{\delta}_1, \dots, \bar{\delta}_n)$ — с помощью преобразования π можно привести к умножению двух конъюнкций

$$\begin{aligned} & (y_1 \vee \dots \vee y_k \vee \bar{y}_{k+1} \vee \dots \vee \bar{y}_l \vee y_{l+1} \vee \dots \vee y_n) \cdot \\ & \cdot (y_1 \vee \dots \vee y_k \vee \bar{y}_{k+1} \vee \dots \vee \bar{y}_l \vee \bar{y}_{l+1} \vee \dots \vee \bar{y}_n) = \\ & = y_1 \vee \dots \vee y_k \vee \bar{y}_{k+1} \vee \dots \vee \bar{y}_l \vee y_{l+1} \cdot \bar{y}_{l+2} \vee \dots \vee y_{n+1} \cdot \bar{y}_n \vee y_n \cdot \bar{y}_{l+1} \end{aligned}$$

Таким образом, д.н.ф., реализующую функцию с двумя нулями можно построить, используя только n конъюнкций. Оказывается, что, обобщив приведенные выше построения, можно легко получить д.н.ф. относительно невысокой сложности, если число нулей функции невелико.

Пусть таблица нулей булевой функции имеет вид

$$\begin{aligned} (\alpha_{11} \dots \alpha_{1i} \dots \alpha_{1n}) &= \tilde{\alpha}_1 \\ (\alpha_{21} \dots \alpha_{2i} \dots \alpha_{2n}) &= \tilde{\alpha}_2 \\ &\dots \\ (\alpha_{k1} \dots \alpha_{ki} \dots \alpha_{kn}) &= \tilde{\alpha}_k \end{aligned}$$

Формулы, реализующая функцию, нуля которой суть наборы $\tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_k$ имеет вид

$$\prod_{i=1}^k (x_1^{\alpha_{i1}} \vee x_2^{\alpha_{i2}} \vee \dots \vee x_n^{\alpha_{in}}) \quad (2.1)$$

Приведем произведение (2.1) к д.н.ф. Сразу исключим из таблицы нулевой и единичный столбцы, так соответствующие им переменные в (2.1) можно сразу вынести за скобки и написать $x_{u_1} \vee \dots \vee x_{u_p} \vee \bar{x}_{t_1} \vee \dots \vee \bar{x}_{t_v}$, соответственно, для нулевых и единичных столбцов.

Выполним преобразование $x_i \rightarrow x_j^{\sigma_{ij}}$, $\sigma_{ij} \in \{0, 1\}$, $\binom{i}{j}$ — подстановка, таким образом, чтобы строка $\tilde{\alpha}_1$ перешла в строку $(0 \ 0 \ \dots \ 0 \ \dots \ 0 \ 0) = \tilde{0}$, одинаковые столбцы таблицы $\|\alpha_{ij}\|_{k \times n}$ получили последовательные номера, объединившись в блоки, число которых не может превосходить $2^{k-1} - 1$ (столбцы образует $k - 1$ строка и нет нулевого столбца).

Пример 2 Исходная таблица T_1 :

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
0	1	0	1	0	1	0	1	0	1
1	0	1	0	1	0	1	0	1	0
1	1	1	1	1	1	0	0	0	0

Проведем сначала преобразование $x_1 \rightarrow x_1$, $x_2 \rightarrow \bar{x}_2$, $x_3 \rightarrow x_3$, $x_4 \rightarrow \bar{x}_4$, $x_5 \rightarrow x_5$, $x_6 \rightarrow \bar{x}_6$, $x_7 \rightarrow x_7$, $x_8 \rightarrow \bar{x}_8$, $x_9 \rightarrow x_9$, $x_{10} \rightarrow \bar{x}_{10}$, а затем номера переменных преобразуем подстановкой

$$\binom{1 \ 3 \ 5 \ 8 \ 10 \ 2 \ 4 \ 6 \ 7 \ 9}{1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10}$$

Получим таблицу:

y_1	y_2	y_3	y_4	y_5	y_6	y_7	y_8	y_9	y_{10}
0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	0	0	0	0	0

Если рассматривать блок отдельно, то имеем функцию, принимающую значения 0 на нулевом и единичном наборах. Напишем формулу:

$$\varphi = y_1 \bar{y}_2 \vee y_2 \bar{y}_3 \vee y_3 \bar{y}_4 \vee y_4 \bar{y}_5 \vee y_5 \bar{y}_1 \vee y_6 \bar{y}_7 \vee y_7 \bar{y}_8 \vee y_8 \bar{y}_9 \vee y_9 \bar{y}_{10} \vee y_{10} \bar{y}_6$$

Функция φ равно 0 на всех трех строках таблицы, но также имеются и "лишние" нули, например:

$$(0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1)$$

$\varphi = 0$ на любом наборе, имеющем одинаковые первые 5 и вторые 5 координат. Чтобы отбросить лишние нули, построим функцию ψ : она зависит от первых переменных каждого из блоков (в нашем случае от y_1 и y_6). Образует таблицу нулей функции ψ из столбцов значений этих переменных. На остальных наборах $\psi = 1$ (в нашем случае это набор $(0 \ 1)$). Реализуем ψ с помощью д.н.ф.: $\psi = \bar{y}_1 \cdot y_6$.

y_1	y_6	ψ
0	1	0
1	1	0
1	0	0
0	1	1

Формула $\varphi \vee \psi$ реализует функцию, нули которой заданы таблицей T_1 : φ равна 0 на любых наборах, которые в каждом из блоков имеют одинаковые значения координат, ψ "отбрасывает" лишние нули, оставляя те, которые реально задаются таблицей.

Теперь достаточно выполнить обратное преобразование, и мы получим формулу для функции, нули которой заданы таблицей T_1 . Заметим, что для функции от 10 переменных мы получили д.н.ф. из $10 + 1$ конъюнкции.

В "наихудшем случае" в преобразованной таблице получится $2^{3-1} - 1 = 3$ блока, и функция ψ примет вид:

y_1	y_{i_1}	y_{i_2}	ψ
0	0	0	1
0	1	1	0
1	0	1	0

На остальных наборах $\psi = 1$. Легко написать д.н.ф. для ψ :

$$\psi = \bar{y}_1 \bar{y}_{i_1} y_{i_2} \vee y_1 \bar{y}_{i_2} \vee \bar{y}_{i_2} y_{i_1} \vee y_1 \bar{y}_{i_2}$$

Таким образом для функций с тремя нулями д.н.ф. для ψ не может состоять более, чем из 4 элементарных конъюнкций.

В общем случае, для функций с k нулевыми наборами $\tilde{\alpha}_1, \dots, \tilde{\alpha}_k$ процесс построения д.н.ф. полностью повторяет действия примера.

Находим преобразование, которое приводит таблицу нулей к виду: 1) первая строка становится нулевой, 2) таблица состоит из блоков одинаковых столбцов (различные блоки состоят из разных столбцов). В каждом блоке присутствуют только строки $\tilde{0} = (0 \dots 0)$ и $\tilde{1} = (1 \dots 1)$. Пусть таблица разбита на блоки $B_1, B_2, \dots, B_i, \dots, B_q$, $q \leq 2^{k-1} - 1$. Пусть блок B_i состоит из столбцов значений переменных

$$x_{t_i}, x_{t_i+1}, \dots, x_{t_i+v_i}$$

Запишем формулу $\mathcal{X}(B_i) = x_{t_i} \cdot \bar{x}_{t_i+1} \vee x_{t_i+1} \cdot \bar{x}_{t_i+2} \vee \dots \vee x_{t_i+v_i-1} \cdot \bar{x}_{t_i+v_i} \vee x_{t_i+v_i} \cdot \bar{x}_{t_i}$. Д.н.ф. $D = \bigvee_{i=1}^q \mathcal{X}(B_i)$ реализует функцию, равную 0 на (и только на) наборах, таких, что значения координат в каждом блоке одинаковы.

Для исключения "лишних" нулей образуем (как в примере) функцию, зависящую от переменных, взятых по одному из каждого блока (например, от первых переменных блока). Соответствующие столбцы образуют множество нулей функции ψ (столбцы значений выбранных переменных). Остальные наборы образуют множество единиц функции ψ , что исключает "лишние" нули и позволяет написать для функции с нулями $\tilde{\alpha}_1, \dots, \tilde{\alpha}_k$ формулу

$$\bigvee_{i=1}^q \mathcal{X}(B_i) \vee \psi$$

Написав д.н.ф. для ψ (ψ зависит не более, чем от $2^{k-1} - 1$ переменной), получим искомую д.н.ф. Применяя к ней преобразования

$$x_i \mathcal{K} \vee \bar{x}_i \mathcal{K}' \rightarrow x_i \mathcal{K} \vee \bar{x}_i \mathcal{K}' \vee \mathcal{K} \mathcal{K}', \quad \mathcal{K} \vee \mathcal{K} \mathcal{K}' \rightarrow \mathcal{K}$$

построим сокращенную д.н.ф. с нулями $\tilde{\alpha}^1, \dots, \tilde{\alpha}^k$.

Замечание. Умножение "тестовых" уравнений приводит к д.н.ф., в которой нет переменных с отрицаниями. Поэтому сокращенная д.н.ф. (совокупность всех простых импликант) получается применением только преобразования

$$\mathcal{K} \vee \mathcal{K} \mathcal{K}' \Rightarrow \mathcal{K}$$

Но $\mathcal{K} \cdot \mathcal{K}'$ не может соответствовать тупику тест, так как тест, соответствующих \mathcal{K} , получается из него удалением некоторых столбцов. Заметим также, что полученная д.н.ф. является единственной тупикувой (для доказательства применить критерий поглощения и заметить, что д.н.ф., не содержащая отрицаний переменных, не может быть равной 1 на всех наборах значений переменных).

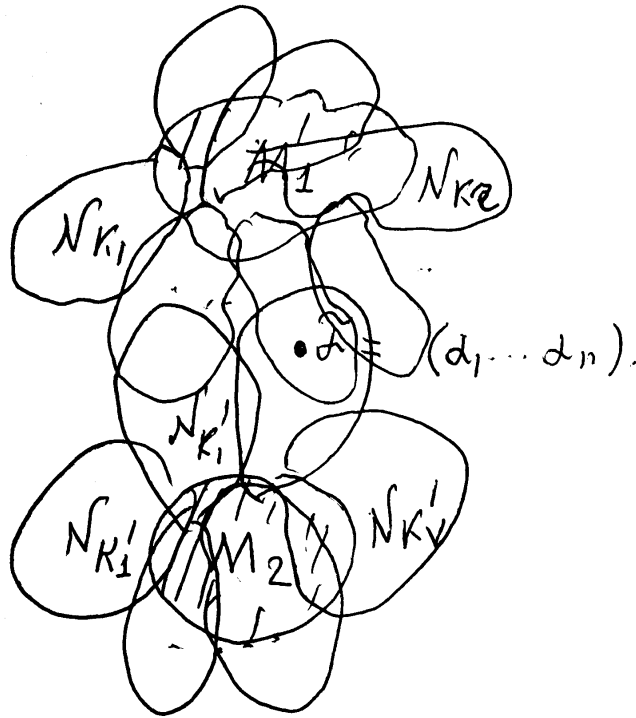
Применение к построению алгоритмов распознавания. Пусть в исходной бинарной таблице наборы $\tilde{\alpha}^i = (\alpha_{i1}, \dots, \alpha_{in})$, $i = 1 \dots k$ принадлежат K_1 , наборы $\tilde{\beta}^j = (\beta_{j1}, \dots, \beta_{jn})$, $j = 1 \dots l$ принадлежат K_2 , $K_1 \cap K_2 = \emptyset$. Построим семейство логических алгоритмов распознавания. Введем две не всюду определенные булевские функции:

$$F_1(x_1, \dots, x_n) = \begin{cases} 1 & \text{на наборах } \tilde{\alpha}^i, i = 1 \dots k \\ 0 & \text{на наборах } \tilde{\beta}^j, j = 1 \dots l \end{cases}$$

$$F_2(x_1, \dots, x_n) = \begin{cases} 1 & \text{на наборах } \tilde{\beta}^j, j = 1 \dots l \\ 0 & \text{на наборах } \tilde{\alpha}^i, i = 1 \dots k \end{cases}$$

С помощью описанной выше процедуры по нулям F_1 и F_2 построим д.н.ф. для функций, равных 0 только на этих наборах, и с помощью преобразований Блейка построим сокращенные д.н.ф. этих функций (на практике число добавляемых конъюнкций ограничивается). Из сокращенной д.н.ф. удалим все конъюнкции, интервалы которых имеют пустое пересечение с множеством единиц соответствующей функции.

Пусть после этих процедур для F_1 остались интервалы $\mathcal{N}_{K_1}, \dots, \mathcal{N}_{K_r}$, имеющие непустое пересечение с $M_1 = \{\tilde{\alpha}^1 \dots \tilde{\alpha}^k\}$, и для F_2 — $\mathcal{N}'_{K'_1}, \dots, \mathcal{N}'_{K'_q}$ с непустым пересечением с $M_2 = \{\tilde{\beta}^1, \dots, \tilde{\beta}^l\}$. Схематически это можно изобразить так:



Здесь $\tilde{\alpha} = (\alpha_1, \dots, \alpha_n)$ — распознаваемый объект. Построим совокупность тушиковых д.н.ф. D_{11}, \dots, D_{1p} для F_1 и D_{21}, \dots, D_{2q} — для F_2 . При этом можно пользоваться критерием поглощения. Но можно проверить условия

$$M_1 \cap \mathcal{N}_{K_i} \subseteq \bigcup_{j \neq i} \mathcal{N}_{K_j} \cap M_1$$

$$M_2 \cap \mathcal{N}'_{K'_i} \subseteq \bigcup_{j \neq i} \mathcal{N}'_{K'_j} \cap M_2$$

и удалить интервалы, пересечение которых с M_1 (M_2) покрывается остающимися интервалами. Для такой проверки существуют более сложный критерий поглощения (см. Ю. И.

Журавлев, "Об отделимости подмножеств вершин единичного n -мерного куба [1]). Пусть построены тупиковые д.н.ф. T_{11}, \dots, T_{1K_1} для F_1 и T_{21}, \dots, T_{2K_2} для F_2 .

Пусть число выполненных равенств $T_{1i}(\tilde{\alpha}) = 1, i = 1 \dots K_1$ равно Q_1 , а выполненных равенств $T_{2i}(\tilde{\alpha}) = 1, i = 1, \dots, K_2$ равно Q_2 .

Простейшее решающее правило:

$$Q_1 > Q_2 \rightarrow \tilde{\alpha} \in K_1,$$

$$Q_1 < Q_2 \rightarrow \tilde{\alpha} \in K_2$$

при $Q_1 = Q_2$ алгоритм отказывается от распознавания.

Возможно усложнение введением порога c

$$Q_1 - c > Q_2 \rightarrow \tilde{\alpha} \in K_1,$$

$$Q_2 - c > Q_1 \rightarrow \tilde{\alpha} \in K_2$$

в остальных случаях алгоритм отказывается от распознавания.

На базе описанного выше алгоритма может быть построено параметрическое семейство — тупиковым д.н.ф. приписаны веса w , и суммируется не число вхождений $\tilde{\alpha}$, а сумма весов тупиковых д.н.ф., интервалы которых содержат точку $\tilde{\alpha}$.

Возможна оптимизация весов по текущему контролю или независимому контрольному материалу (см. подбор значений параметров алгоритмов вычисления оценок).

Лекция 3

3.1 Алгоритмы вычисления оценок

В эвристическом алгоритме “Тест” было проведено сравнение объектов из таблицы обучения и распознаваемого объекта по подмножествам признаков, образующих тупиковые тесты. Очевидно, такие сравнения могут выполняться и по другим подмножествам, например, выбираемым экспертами. При формировании величин $\Gamma(S, K_1), \Gamma(S, K_2)$ добавлялась 1 при выявлении любого совпадения. Однако в реальности признаки и объекты из таблицы обучения не равноценны, и при формировании $\Gamma(S, K_i)$ следует учитывать — с какими объектами по каким признакам произошло совпадение. Это можно делать, задавая различные веса признакам и объектам из таблицы обучения — эталонным объектам. При самом определении близости, если признаки не бинарные, можно также использовать различные возможности, определенные, например, параметрами $\epsilon_1, \dots, \epsilon_n$. Наконец, “поощрения” при формировании $\Gamma(S, K_i)$ возможны также за отсутствие близости к объектам другого класса, а штраф — за наличие близости к объектам класса, для которого не формируется оценка и отсутствие близости к объектам, для которых оценка формируется. Указанные обстоятельства могут реализовываться многими различными способами. Формальное описание таких способов приводит к формированию класса алгоритмов, получившего название “алгоритмы вычисления оценок”, или “алгоритмы голосования”.

Как и ранее, будем рассматривать исходную (или обучающую) информацию, заданную в виде двух таблиц:

$\|a_{ij}\|_{m \times n}$ — совокупность m объектов, заданных наборами n признаков,

$\|\alpha_{ij}\|_{m \times l}$ — информационная матрица (таблица), где строка $(\alpha_{i1} \dots \alpha_{ij} \dots \alpha_{il}) = \tilde{\alpha}(S_i)$ указывает — каким из классов K_1, \dots, K_l принадлежит или не принадлежит объект S_i , определенный строкой $(a_{i1} \dots a_{in})$ значений признаков $1, 2, \dots, n$ из таблицы обучения. Как и ранее, полагаем, что области определения признаков — это метрические пространства M_t с метриками $\rho_t, t = 1, \dots, n$.

Алгоритмы вычисления оценок определяются:

I. Заданием системы опорных множеств признаков.

Это могут быть любые множества, элементами которых являются непустые подмножества множества признаков $1, 2, \dots, n$. Такими могут быть:

- а) совокупность всех непустых подмножеств множества $\{1, 2, \dots, n\}$,
- б) совокупность всех подмножеств из k элементов и т. д.

В случаях а) и б) имеем, соответственно, $2^n - 1$ и $\binom{n}{k}, k = 1, 2, \dots, n - 1$ подмножеств признаков, по которым происходит сравнение эталонных и распознаваемого объекта.

Вообще говоря, может быть задана любая совокупность $\{\Omega\}_A$ опорных множеств, задающих распознающий алгоритм A . Для удобства, в некоторых случаях, вместо

опорного множества $\Omega = \{u_1, \dots, u_k\}$ будем рассматривать характеристический вектор $\tilde{\omega} = \{\sigma_1, \dots, \sigma_n\}$, где $\sigma_{u_1} = \dots = \sigma_{u_k} = 1$, и остальные координаты равны 0. Очевидно: $\Omega \leftrightarrow \tilde{\omega}, \{\Omega\}_A \leftrightarrow \{\tilde{\omega}\}_A$.

Введем понятия $\tilde{\omega}$ -части объекта S и таблицы $\|a_{ij}\|_{m \times n}$:

$\tilde{\omega}$ -часть строки $S = (a_1, \dots, a_n)$ — обозначение $\tilde{\omega}S$ — это набор

$$\tilde{\omega}S = (a_{u_1}, \dots, a_{u_k}), \quad \tilde{\omega}(\|a_{ij}\|_{m \times n}) = \begin{pmatrix} \tilde{\omega}S_1 \\ \vdots \\ \tilde{\omega}S_m \end{pmatrix} = \begin{pmatrix} a_{1u_1} & \dots & a_{1u_k} \\ a_{2u_1} & \dots & a_{2u_k} \\ \dots & \dots & \dots \\ a_{mu_1} & \dots & a_{mu_k} \end{pmatrix}.$$

Будем также использовать обозначение $\tilde{\omega}T_1, T_1$ — таблица обучения.

II. Заданием функции близости $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i)$ между $\tilde{\omega}$ -частями распознаваемого объекта S_i и эталонного объекта S . В дальнейшем рассматриваются только функции близости, принимающие значения 0, 1. Тогда корректно введение функции $\bar{\mathcal{N}}(\tilde{\omega}S, \tilde{\omega}S_i)$, указывающей на отсутствие близости между $\tilde{\omega}S$ и $\tilde{\omega}S_i$.

Обычно рассматриваются три вида функций близости:

- 1) введем неотрицательные параметры $\epsilon_1, \dots, \epsilon_n$. Пусть $\tilde{\omega}S = (a_{u_1}, \dots, a_{u_k})$, $\tilde{\omega}S_i = (a_{iu_1}, \dots, a_{iu_k})$. Тогда

$$\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = \begin{cases} 1, & \text{если } \begin{cases} \rho(a_{u_1}, a_{iu_1}) \leq \epsilon_{u_1} \\ \dots \\ \rho(a_{u_k}, a_{iu_k}) \leq \epsilon_{u_k} \end{cases} \\ 0, & \text{если хотя бы одно из этих неравенств не выполнено.} \end{cases}$$

- 2) кроме параметров $\epsilon_1, \dots, \epsilon_n$ введем неотрицательный целочисленный параметр ν : $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$, если среди приведенных выше неравенств не более ν не выполнены, и равна 0 в остальных случаях.

Пусть $|\Omega|$ — число элементов в опорном множестве $\Omega \in \{\Omega\}_A$; пусть также

$$\min_{\Omega \in \{\Omega\}_A} |\Omega| = q.$$

Легко видеть, что следует рассматривать только значения ν , удовлетворяющие неравенству

$$0 \leq \nu \leq \left\lfloor \frac{q}{2} \right\rfloor - 1.$$

- 3) вместо параметра ν можно рассмотреть параметр ν^* и определить функцию близости следующим образом: $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$ тогда и только тогда, когда из k неравенств

$$\begin{aligned} \rho(a_{u_1}, a_{iu_1}) &\leq \epsilon_{u_1} \\ \dots & \\ \rho(a_{u_k}, a_{iu_k}) &\leq \epsilon_{u_k} \end{aligned}$$

не выполнены r неравенств, и $\frac{r}{k} < \nu^*$.

В практических системах распознавания, в основном, применяется функция близости, зависящая только от параметров $\epsilon_1, \dots, \epsilon_n$.

- III. Признакам $1, 2, \dots, n$, описывающим объекты, присваиваются веса w_1, w_2, \dots, w_n . Как правило, $w_i \geq 0$, $i = 1 \dots n$. Но последнее ограничение не является обязательным. Объектам из исходной таблицы $T_1 : S_1, \dots, S_m$ приписываются веса $w(S_1) = w^1, \dots, w(S_m) = w^m$. Здесь $w^i \geq 0$, $i = 1 \dots m$. Множеству $\tilde{\omega}S_i = (a_{iu_1}, \dots, a_{iu_k})$ приписывается вес (число голосов) $\Gamma(\tilde{\omega}S_i) = w^i \cdot (w_{u_1} + \dots + w_{u_k})$.
- IV. При сравнении $\tilde{\omega}S$ и $\tilde{\omega}S_i$ возможны следующие случаи, сведенные в таблицу и оцененные параметрами x_{ij} , $i = 0, 1$, $j = 0, 1$.

$S_i \setminus K_j$	\mathcal{N}	$\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$	$\overline{\mathcal{N}}(\tilde{\omega}S, \tilde{\omega}S_i) = 0$
$S_i \in K_j$		x_{11}	x_{10}
$S_i \notin K_j$		x_{01}	x_{00}

При формировании оценки принадлежности S классу K_j , $1 \leq j \leq l$ (величины $\Gamma_j(\tilde{\omega}S, \tilde{\omega}S_i)$) учитываются оценка $\Gamma(\tilde{\omega}S, \tilde{\omega}S_i)$ и то, какой из четырех указанных случаев имеет место: оценка $\Gamma(\tilde{\omega}S, \tilde{\omega}S_i)$ умножается на соответствующий параметр. Так, если $S_i \in K_j$, $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$, то

$$\Gamma_j(\tilde{\omega}S, \tilde{\omega}S_i) = x_{11} \cdot \Gamma(\tilde{\omega}S, \tilde{\omega}S_i).$$

Аналогично формируются оценки в остальных трех случаях. Заметим, что естественно полагать

$$x_{11} \geq 0, \quad x_{00} \geq 0, \quad x_{01} \leq 0, \quad x_{10} \leq 0, \quad x_{00} < x_{11}.$$

Действительно, близость к объекту из K_j или отсутствие близости с объектом, не принадлежащим K_j , благоприятны для оценки вхождения S в K_j , причем второй случай не более благоприятен, чем первый. Два других случая: $S_i \notin K_j$, $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$ и $S_i \in K_j$, $\overline{\mathcal{N}}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$ не должны повышать оценку вхождения S в K_j .

Заметим, что приведенные здесь объяснения не являются строгими доказательствами, но лишь эвристическими правдоподобными рассуждениями, оправдывающими (в некоторой степени) даваемое нами определение класса алгоритмов вычисления оценок.

- V. Оценка $\Gamma_j(S)$ вхождения объекта S в класс K_j задается следующей формулой

$$\Gamma_j(S) = \frac{1}{Q} \sum_{i=1}^m \sum_{\tilde{\omega} \in \{\tilde{\omega}\}_A} \Gamma_j(\tilde{\omega}S, \tilde{\omega}S_i), \quad (3.1)$$

Q — нормирующий множитель. Его величина не влияет на дальнейшие преобразования $\Gamma_j(S)$. Поэтому достаточно рассмотреть случай $Q = 1$.

- VI. Решающее правило определяется числовыми параметрами c_1, c_2 , $0 < c_1 < c_2$. Алгоритм A формирует для S информационный (квази-информационный) вектор $(\beta_1(S), \dots, \beta_j(S), \dots, \beta_l(S))$ следующим образом:
 $\beta_j(S) = 1 \rightarrow S \in K_j$, если $\Gamma_j(S) > c_2$;
 $\beta_j(S) = 0 \rightarrow S \notin K_j$, если $\Gamma_j(S) < c_1$;
в остальных случаях $\beta_j(S) = \Delta$, что означает: алгоритм A отказался распознавать вхождение S в класс K_j .

Отметим, что алгоритм вычисления оценок (АВО) подразделяется на две части. После выполнения этапа V формируется числовой вектор оценок $(\Gamma_1(S), \dots, \Gamma_j(S), \dots, \Gamma_l(S)) = \vec{\Gamma}_l(S)$. Эту часть алгоритма принято называть распознающим оператором B :

$$B(I, S) = (\Gamma_1(S), \dots, \Gamma_j(S), \dots, \Gamma_l(S)) = \vec{\Gamma}_l(S).$$

Вторая часть алгоритма (этап VI) переводит $\vec{\Gamma}_l(S)$ в квази-информационный вектор. Это — решающее правило C :

$$C(\vec{\Gamma}_l(S)) = C(\Gamma_1(S), \dots, \Gamma_j(S), \dots, \Gamma_l(S)) = (C(\Gamma_1(S)), \dots, C(\Gamma_j(S)), \dots, C(\Gamma_l(S))) = (\beta_1, \dots, \beta_j, \dots, \beta_l).$$

Резюме. Алгоритм A определяется заданием системы $\{\Omega\}_A$ опорных множеств, параметров $\epsilon_1, \dots, \epsilon_n$ (возможно и ν или ν^*), определяющих функцию близости, весов признаков w_1, \dots, w_n , весов эталонных объектов w^1, \dots, w^m , параметров x_{11}, x_{00} , “поощряющих” благоприятные ситуации, x_{01}, x_{10} , “штрафующих” за неблагоприятные ситуации, параметров c_1, c_2 порогового решающего правила, с помощью которых принимается окончательное решение о вхождении, невхождении распознаваемого объекта S в классы K_1, \dots, K_l или, для некоторых классов (может быть, для всех), — об отказе от распознавания.

Пример 3

	1	2	3	4	5	6		K_1	K_2
S_1	0	0	3	2	0.7	0.8		1	0
S_2	0	1	1	4	0.6	0.7		0	1
S_3	1	1	4	1	0.5	0.6		1	1
	I								
S	1	1	2	3	0.6	0.8			

Алгоритм определяется следующими множествами и параметрами:

$$\begin{aligned} \{\Omega\}_A &= \{(1, 2), (3, 4), (5, 6)\}; \\ \epsilon_1 = \epsilon_2 &= 0, \quad \epsilon_3 = \epsilon_4 = 1, \quad \epsilon_5 = \epsilon_6 = 0.1; \\ w^1 = w^2 &= 2, \quad w^3 = 1; \\ w_1 = w_2 &= 1, \quad w_3 = w_4 = 2, \quad w_5 = w_6 = 3; \\ x_{11} = 3, \quad x_{00} &= 1, \quad x_{01} = x_{10} = -1. \end{aligned}$$

Таблица значений функции близости

	Ω_1	Ω_2	Ω_3
S_1	0	1	1
S_2	1	0	1
S_3	0	0	0

Таблица использования параметров

$x_{ij}, i = 0, 1; j = 0, 1$

x_{10}	x_{11}	x_{11}
x_{01}	x_{00}	x_{01}
x_{01}	x_{01}	x_{01}

Таблица T_Ω оценок $\Gamma_1(\tilde{\omega}S, \tilde{\omega}S_i)$

	Ω_1	Ω_2	Ω_3
S_1	4	8	12
S_2	4	8	12
S_3	2	4	6

Таблица $T_{x_{ij}}$ значений параметров x_{ij}

-1	3	3
-1	1	-1
-1	-1	-1

Умножим поэлементно матрицу T_Ω на матрицу $T_{x_{ij}}$ и сложим все элементы матрицы $T_\Omega \cdot T_{x_{ij}}$. Получим число 36. При $c_2 < 36$ алгоритм зачислит объект S в класс K_1 . Аналогично вычисляется оценка $\Gamma_2(S)$. Вычисление предоставляется выполнить читателю.

3.2 Эффективные формулы вычисления оценок

Прямое использование формулы (3.1) для вычисления оценок $\Gamma_j(S)$, $j = 1, \dots, l$ затруднительно при большом числе множеств. Так, если $\{\Omega\}_A$ состоит из всех непустых подмножеств множества $\{1, 2, \dots, n\}$, то потребовалось бы вычислить $m \cdot (2^n - 1)$ слагаемых. Поэтому рассмотрим пути сокращения вычислений. Рассмотрим вычисления при фиксированной строке S_i :

$$\sum_{\tilde{\omega} \in \{\tilde{\omega}\}_A} \Gamma_j(\tilde{\omega}S, \tilde{\omega}S_i).$$

Разберем два различных случая: $1^\circ S_i \in K_j$ и $2^\circ S_i \notin K_j$. В случае 1° имеем два подслучая $1^\circ 1$ и $1^\circ 0$: $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$ и $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 0$.

$$1^\circ 1: \quad x_{11} \cdot w^j \cdot \sum_{\substack{\Omega: \mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i)=1 \\ \Omega \in \{\Omega\}_A}} P(\Omega), \quad P(\Omega) = w_{i_1} + \dots + w_{i_k}, \text{ если } \Omega = \{i_1, \dots, i_k\}.$$

Обозначим через $R(\mathcal{N} = 1, t)$ число опорных множеств, содержащих t и участвующих в суммировании. Очевидно, величина w_t встретится при суммировании $R(\mathcal{N} = 1, t)$ раз, $t = 1, \dots, n$. Поэтому формулу для $1^\circ 1$ можно переписать:

$$x_{11} \cdot w^j \cdot \sum_{t=1}^n w_t \cdot R(\mathcal{N} = 1, t).$$

В подслучае $1^\circ 0$ вместо множителя x_{11} появится x_{10} , а величина $R(\mathcal{N} = 1, t)$ заменится на $R(\mathcal{N} = 0, t)$, где $R(\mathcal{N} = 0, t)$ — число опорных подмножеств Ω , содержащих t и таких, что $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 0$. Окончательно при $S_i \in K_j$ получаем:

$$w^j(x_{11} \cdot \sum_{t=1}^n w_t \cdot R(\mathcal{N} = 1, t) + x_{10} \cdot \sum_{t=1}^n w_t \cdot R(\mathcal{N} = 0, t)) = \Gamma_j^1(S_i \in K_j). \quad (3.2)$$

При 2° , действуя в точности как в 1° , выводим:

$$w^j(x_{00} \cdot \sum_{t=1}^n w_t \cdot R(\mathcal{N} = 0, t) + x_{10} \cdot \sum_{t=1}^n w_t \cdot R(\mathcal{N} = 1, t)) = \Gamma_j^0(S_i \notin K_j). \quad (3.3)$$

Окончательно в формуле

$$\sum_{\Omega \in \{\Omega\}_A} \Gamma_j(\tilde{\omega}S, \tilde{\omega}S_i)$$

суммирование по опорным множествам заменяется на два суммирования n слагаемых и вычисление значений $R(\mathcal{N} = 1, t)$, $R(\mathcal{N} = 0, t)$.

Нами доказана

Теорема 3

$$\Gamma_j(S) = \frac{1}{Q} \left(\sum_{S_i \in K_j} \Gamma_j^1(S_i \in K_j) + \sum_{S_i \notin K_j} \Gamma_j^0(S_i \notin K_j) \right),$$

где величины Γ_j^1 , Γ_j^0 определяются по формулам (3.2), (3.3).

1. Пусть система $\{\Omega\}_A$ состоит из всех k -элементных подмножеств множества $\{1, 2, \dots, n\}$. Функция близости определяется только параметрами $\epsilon_1, \epsilon_2, \dots, \epsilon_n$. Пусть также $S_i = (a_{i1}, \dots, a_{ik}, \dots, a_{in})$, $S = (a_1, \dots, a_k, \dots, a_n)$. Выпишем неравенства

$$\begin{aligned} \rho_1(a_1, a_{i1}) &\leq \epsilon_1 \\ \dots \\ \rho_k(a_k, a_{ik}) &\leq \epsilon_k \\ \dots \\ \rho_n(a_n, a_{in}) &\leq \epsilon_n \end{aligned}$$

Совокупность номеров признаков, для которых выполнены или, соответственно, не выполнены неравенства, обозначим через M^+ , M^- , а мощности соответствующих множеств через $|M^+|$, $|M^-|$.

Легко видеть, что $R(\mathcal{N} = 1, t)$ равно $\binom{|M^+|-1}{k-1}$ для $t \in M^+$, и равно 0 для $t \in M^-$. Так как число подмножеств, содержащих t в $\{\Omega\}_A$, равно $\binom{n-1}{k-1}$, то

$$R(\mathcal{N} = 0, t) = \begin{cases} \binom{n-1}{k-1} - \binom{|M^+|-1}{k-1}, & \text{для } t \in M^+ \\ \binom{n-1}{k-1}, & \text{для } t \in M^-. \end{cases}$$

Из последнего видно, что $\sum_{t=1}^n w_t \cdot R(\mathcal{N} = 1, t)$, $\sum_{t=1}^n w_t \cdot R(\mathcal{N} = 0, t)$ в рассматриваемом случае заменится более “простыми” суммами, так как величины $R(\mathcal{N} = 1, t)$, $R(\mathcal{N} = 0, t)$ для каждой строки S_i принимают только два различных значения. Так,

$$\sum_{t=1}^n w_t \cdot R(\mathcal{N} = 1, t) = \left(\sum_{t \in M^+} P_t \right) \cdot \binom{|M^+|-1}{k-1}.$$

Аналогично упрощаются и другие введенные ранее формулы.

2. Рассмотрим ту же, что и в 1, систему опорных множеств и функцию близости с параметрами $\epsilon_1, \dots, \epsilon_n$, ν , $\nu < \lfloor \frac{k}{2} \rfloor - 1$. Аналогично предыдущему пункту, при сравнении строк S и S_i образуем подмножества признаков-координат M^+ и M^- .

Пусть $t \in M^+$. Тогда функция близости равна 1, если она содержит $0, 1, 2, \dots, \min(\nu, |M^-|)$ признаков из M^- и, соответственно, $k-1, k-2, \dots, k - \min(\nu, |M^-|) - 1$ признаков из M^+ (для упрощения выкладок мы рассматриваем только случай $k \geq \min(\nu, |M^-|) + 1$). Тогда число опорных подмножеств, содержащих t и таких, что $\mathcal{N}(\tilde{\omega}S, \tilde{\omega}S_i) = 1$, очевидно, равно

$$\begin{aligned} &\binom{|M^+|-1}{k-1} \cdot \binom{|M^-|}{0} + \binom{|M^+|-1}{k-2} \cdot \binom{|M^-|}{1} + \dots + \binom{|M^+|-r}{k-1-r} \cdot \binom{|M^-|}{r} + \dots + \\ &+ \binom{|M^+|-1}{k-1-\nu} \cdot \binom{|M^-|}{\nu}. \end{aligned}$$

Все t из M^+ имеют одинаковый только что выписанный коэффициент при P_t . Остальные случаи вычисляются так же просто.

3. В качестве опорных рассматриваются все непустые подмножества множества $\{1, 2, \dots, n\}$. Тогда

$$R(\mathcal{N} = 1, t) = \begin{cases} 2^{|M^+|-1}, & \text{для } t \in M^+ \\ 0, & \text{для } t \in M^-. \end{cases}$$

$$R(\mathcal{N} = 0, t) = \begin{cases} 2^{|M^+|-1} \cdot (2^{|M^-|} - 1), & \text{для } t \in M^+ \\ 2^{n-1}, & \text{для } t \in M^-. \end{cases}$$

Подставляя полученные значения в (3.2), получаем

$$w^i \left(x_{11} \cdot \left(\sum_{t \in M^+} w_t \right) \cdot 2^{|M^+|-1} + x_{10} \left(\left(\sum_{t \in M^+} w_t \right) \cdot (2^{|M^+|-1}) \cdot (2^{|M^-|} - 1) + \left(\sum_{t \in M^-} w_t \right) \cdot 2^{n-1} \right) \right)$$

Аналогичный вид принимает после соответствующих подстановок формула (3.3).

4. Совокупность характеристических векторов опорных множеств образует интервал в E^n , т.е. удовлетворяет условию: $\mathcal{K} = 1$, \mathcal{K} — элементарная конъюнкция. Не ограничивая общности, можно полагать

$$\mathcal{K} = \overline{x_1} \cdot \dots \cdot \overline{x_r} \cdot x_{r+1} \cdot \dots \cdot x_{r+k}.$$

Тогда в систему опорных множеств не войдут признаки $1, \dots, r$; в каждое из опорных множеств войдут признаки $r+1, \dots, r+k$, и к ним последовательно присоединятся все подмножества (включая пустое, если $\mathcal{K} \neq 0$) или все непустые подмножества множества $r+k+1, \dots, n$, если $\mathcal{K} = 0$. Будем рассматривать $\mathcal{K} \neq 0$ и рассмотрим функцию близости с параметрами $\epsilon_1, \dots, \epsilon_n$.

Если не выполнено включение $\{r+1, \dots, r+k\} \subseteq M^+$, то $R(\mathcal{N} = 1, t) = 0$, $R(\mathcal{N} = 0, t) = 2^{n-(r+k+1)}$. Последнее следует из того, что $|\{\Omega\}_A| = 2^{n-(r+k)}$, и хотя бы один из признаков любого опорного подмножества принадлежит M^- .

Пусть имеет место: $\{r+1, \dots, r+k\} \subseteq M^+$. Тогда для $t \in \{r+1, \dots, r+k\}$:

$$R(\mathcal{N} = 1, t) = 2^{|M^+|-k},$$

$$R(\mathcal{N} = 0, t) = 2^{|M^-|} - 1.$$

Для $t \in M^+ \setminus \{r+1, \dots, r+k\}$:

$$R(\mathcal{N} = 1, t) = 2^{|M^+|-k-1},$$

$$R(\mathcal{N} = 0, t) = 2^{|M^+|-k-1} \cdot (2^{n-(k+r)} - 1).$$

Наконец, для $t \in M^-$: $R(\mathcal{N} = 1, t) = 0$, $R(\mathcal{N} = 0, t) = 2^{|M^-|-1}$.

После соответствующих подстановок, получаем эффективные формулы вычисления $\Gamma_j(S)$, $j = 1, \dots, l$.

Последний случай (4) можно использовать при решении прикладных задач. Рассмотрим булевскую функцию, равную 1 на элементах $\{\tilde{\omega}\}_A$. Если реализовать ее дизъюнктивной нормальной формой $\bigvee_{i=1}^r \mathcal{K}_i$, где \mathcal{K}_i — элементарные конъюнкции, и

$\mathcal{N}_{\mathcal{K}_u} \cap \mathcal{N}_{\mathcal{K}_v} = \emptyset$, то, написав по 4 формулы для каждой \mathcal{K}_i и сложив их, получим формулу для вычисления $\Gamma_j(S)$, $j = 1, \dots, l$. Достаточно и реализации в классе д.н.ф., в которых интервалы некоторых пар конъюнкций пересекаются (интервалы конъюнкций из разных пар не пересекаются). Пусть оценка по системе опорных множеств, задаваемых конъюнкцией \mathcal{K} , есть $\Gamma_j(\mathcal{K}, S)$. Пусть, также, $\{\tilde{\omega}\}_A = \mathcal{K}_1 \cup \mathcal{K}_2$, $\mathcal{K}_1 \cdot \mathcal{K}_2 \neq 0$. Тогда, очевидно, $\Gamma_j(S) = \Gamma_j(S, \mathcal{K}_1) + \Gamma_j(S, \mathcal{K}_2) - \Gamma_j(S, \mathcal{K}_1 \cdot \mathcal{K}_2)$.

Лекция 4

4.1 Вычисление характеристик, определяющих алгоритм вычисления оценок

Как было показано в лекции 3, алгоритм распознавания определяется заданием системы опорных множеств $\{\Omega\}_A$ и числовых параметров $\epsilon_1, \dots, \epsilon_n, w_1, \dots, w_n, w^1, \dots, w^m, x_{11}, x_{00}, x_{01}, x_{10}, c_1, c_2$.

В своей работе алгоритм использует исходную (обучающую) информацию, состоящую из таблицы обучения $\|a_{ij}\|_{m \times n}$ и ее информационной матрицы $\|\alpha_{ij}\|_{m \times l}$. Рассматривается задача с l , вообще говоря, пересекающимися классами K_1, \dots, K_l .

Параметры подбираются таким образом, чтобы обеспечить максимальную точность распознавания на определенном заранее множестве объектов.

I. Текущий контроль. Из исходной матрицы последовательно изымаются строки $S_i, i = 1 \dots m$, вместе с информационным вектором $\tilde{\alpha}(S_i)$, и для строки S_i по оставшейся исходной информации строится квази-информационный вектор $\tilde{\beta}(S_i) = (\beta_{i1} \dots \beta_{il}), i = 1 \dots m$

В матрице $\|\alpha_{ij} - \beta_{ij}\|_{m \times l}$ определяется число единиц. Операция $\alpha_{ij} - \beta_{ij}$ определяется следующим образом:

$\alpha_{ij} \backslash \beta_{ij}$	0	1	Δ
0	0	1	1
1	1	0	1

Алгоритм A подбирается таким образом, чтобы число единиц (т.е. сумма ошибок и отказов) была бы минимальной.

II. Независимый контроль. формируется контрольное множество $S^1, \dots, S^q, S_i = (b_{i1} \dots b_{in}), i = 1 \dots q$ и таблицы информационных векторов $\|\beta_{ij}\|_{q \times l}$. С использованием алгоритма A и всей исходной информации формируется совокупность $\tilde{\gamma} = (\gamma_{i1} \dots \gamma_{il})$ квази-информационных векторов и минимизируется число единиц в матрице $\|\beta_{ij} - \gamma_{ij}\|_{q \times l}$.

- 1) Система опорных множеств определяется (задается) экспертами и последовательно рассматриваются алгоритмы в набором всех k -элементарных подмножеств: $k = 2, 3, \dots, r$. Как правило, достаточно ограничиться $r \leq 3\sqrt{\pi}$.
- 2) Параметры x_{ij} определяются перебором вариантов. Как правило, рассматриваются целочисленные значения $x_{11} = 1, 2, 3, 4, 5, 6, 7; x_{00} \leq [\frac{1}{2}x_{11}], x_{01}, x_{10}$ принимают отрицательные или нулевые значения. В большинстве действующих систем полагают: $x_{11} = 1, x_{00} = x_{01} = x_{10} = 0$.

- 3) Пусть определения все характеристика за исключением $\epsilon_1, \dots, \epsilon_n$ (о них позднее).

Рассмотрим задачу с независимым контролем и напишем систему неравенств

$$\begin{aligned} \Gamma_j(S^i) &> c_2, \forall \beta_{ij} = 1 \\ \Gamma_j(S^i) &< c_1, \forall \beta_{ij} = 0 \\ i &= 1, 2, \dots, q, j = 1, \dots, l \end{aligned} \quad (4.1)$$

В левой части системы неравенств (4.1) находятся билинейные формы $\Sigma w_u \cdot w^v$. Полагаем $w_1 = \dots = w_n = 1$. Получаем систему линейных относительно $w^1 \dots w^q$ неравенств. Находим максимальную совместную подсистему и ее решение w_1^1, \dots, w_1^q . Подставляем эти значения в левую часть (4.1) и получаем линейную относительно w_1, \dots, w_n систему. Находим совместную максимальную подсистему и ее решение w_{11}, \dots, w_{n1} . Подставляем эти значения левую часть (4.1) и т.д. Процесс заканчивается либо когда удастся получить наборы параметров, удовлетворяющие всем неравенствам (4.1), либо когда после очередной итерации число неравенств в совместной подсистеме уменьшится.

- 4) Для определения $\epsilon_1, \dots, \epsilon_n$ существует большое число эвристических методов. Приведем один из них (может быть, не лучший). Оставим в таблице обучения и контроля только k -е столбцы:

$$\begin{pmatrix} a_{1k} \\ a_{2k} \\ \vdots \\ a_{mk} \end{pmatrix}, \begin{pmatrix} b_{1k} \\ b_{2k} \\ \vdots \\ b_{lk} \end{pmatrix}$$

Для каждой пары $(a_{ik}b_{jk})$ вычислим $\|\tilde{\alpha}(S_i) + \tilde{\beta}(S^k)_{mod2}\|$, то есть число различных координат в этом векторе.

Если $l - \|\tilde{\alpha}(S_i) + \tilde{\beta}(S^k)\| > \|\tilde{\alpha}(S_i) + \tilde{\beta}(S^k)\|$, то формируем неравенство

$$|a_{ik} - b_{jk}| < \epsilon_k$$

При изменении знака:

$$|a_{ik} - b_{jk}| > \epsilon_k.$$

При $l - \|\tilde{\alpha}(S_i) + \tilde{\beta}(S^k)\| = \|\tilde{\alpha}(S_i) + \tilde{\beta}(S^k)\|$ выписываем одно из неравенств

$$|a_{ik} - b_{jk}| \leq \epsilon_k, \quad |a_{ik} - b_{jk}| \geq \epsilon_k.$$

Значение ϵ_k находится из условия: данное значение удовлетворяет наибольшему числу построенных неравенств.

4.2 Алгебры над алгоритмами

Ранее мы видели, что алгоритм вычисления оценок делится на две части: распознающий оператор B и решающее правило C : $A = B \cdot C$. $A(I, S) = (\Gamma_1(S), \dots, \Gamma_l(S)) = \vec{\Gamma}_l(S)$, $C(\vec{\Gamma}_l(S)) = (C(\Gamma_1(S)), \dots, C(\Gamma_l(S)))$

$$C(\Gamma_j(S)) = \begin{cases} 1, & \Gamma_j(S) > c_2 \\ 0, & \Gamma_j(S) < c_1 \\ \Delta, & c_1 \leq \Gamma_j(S) \leq c_2 \end{cases} \quad c_1 < c_2$$

Оказывается, что подобное представление имеет место для большого класса алгоритмов.

Пусть A — алгоритм, работающий с исходной информацией $I \in \{I\}$, $A \in \{A\}$, и каждый из A по любой $I \in \{I\}$ должен получить ответ на фиксированные вопросы Q_1, \dots, Q_l , причем число возможных ответов равно трем: 1, 0, Δ . Тогда

Теорема 4 *Каждый A может быть представлен в виде $A = B \cdot C$, $B(I) = (a_1 \dots a_j \dots a_l)$ — числовой вектор \vec{a} , $C(\vec{a}) = (C(a_1), \dots, C(a_l))$, причем*

$$C(a_i) = \begin{cases} 1, & a_i > c_2 \\ 0, & a_i < c_1 \\ \Delta, & c_1 \leq a_i \leq c_2 \end{cases}$$

где c_1 и c_2 — константы, фиксированные для всех алгоритмов. Заметим, что $A(I) = (C(a_1), \dots, C(a_l)) = (\delta_1 \dots \delta_l) = \vec{\delta}$

Доказательство. Введем вспомогательный оператор $C^{-1}(\vec{\delta}) = (C^{-1}(\beta_1) \dots C^{-1}(\beta_l)) = (a_1, \dots, a_l) = \vec{a}$. Тогда $B = A \cdot C^{-1}$, $A = (A \cdot C^{-1}) \cdot C$. Теорема доказана.

Класс алгоритмов $\{A\}$ порождает класс операторов $\{B\}$, которые можно складывать, умножать, умножать на число.

Действительно, если $B_1(I) = (a_{11} \dots a_{1l})$, $B_2(I) = (a_{21} \dots a_{2l})$, то $(B_1 + B_2)(I) = (a_{11} + a_{21}, \dots, a_{1l} + a_{2l})$, $B_1 \cdot B_2(I) = (a_{11} \cdot a_{21}, \dots, a_{1l} \cdot a_{2l})$, $(d \cdot B_1)(I) = (d \cdot a_{11}, \dots, d \cdot a_{1l})$.

Нетрудно видеть, что используя операторы из $\{B\}$, $A = B \cdot C$, можно построить полиномы

$$\tilde{B} = \sum c_{i_1 \dots i_k} \cdot B_{i_1}^{r_1} \cdot \dots \cdot B_{i_k}^{r_k},$$

где роль переменных играют операторы; $c_{i_1 \dots i_k}$ — константы.

Совокупность $\{\tilde{B}\}$ называют алгебраическим замыканием семейства $\{B\}$, а $\{\tilde{B}\} \cdot C$ — алгебраическим замыканием класса алгоритмов $\{A\}$ — обозначения $\mathfrak{U}\{B\}$, $\mathfrak{U}\{A\}$.

Оказывается (Ю. И. Журавлев), что в $\mathfrak{U}\{A\}$ при выполнении простых легко проверяемых условий можно построить алгоритм, не делающий ошибок на контрольной совокупности.

Алгоритм имеет вид:

$$(d \cdot \sum (c_i B_i)^{k_i}) \cdot C.$$

B_i представимы линейными формами от распознающих операторов вычисления оценок.

Условия:

Пусть $I = \{\|a_{ij}\|_{m \times n} \| \alpha_{ij} \|_{m \times l}\}$.

Контрольный материал: $\|b_{uv}\|_{q \times n}$, $\|\beta_{ij}\|_{q \times l}$.

1. в матрице $\|\alpha_{ij}\|_{m \times l}$ нет одинаковых столбцов.
2. для каждой пары S^u, S^v контрольных объектов, $S^u = (b_{u1} \dots b_{un})$, $S^v = (b_{v1} \dots b_{vn})$ найдется $S_r \in I$, $S_r = (a_{r1} \dots a_{rl})$ и признак k , $1 \leq k \leq n$, $r = r(u, v)$, $k = k(u, v)$ такие, что

$$\rho_k(a_{rk}, b_{uk}) \neq \rho_k(a_{rk}, b_{vk})$$

Лекция 5

5.1 Построение алгоритмов распознавания, корректных для заданной контрольной выборки

Рассматривается задача распознавания (или прогноза) со стандартной обучающей информацией:

$$\begin{aligned} I_0 &= \{S_1, \dots, S_m, \tilde{\alpha}(S_1), \dots, \tilde{\alpha}(S_m)\}, \\ S_i &= (a_{i1}, \dots, a_{in}), \quad i = 1, 2, \dots, m, \\ a_{ij} &\in M_j, \quad j = 1, 2, \dots, n. \end{aligned}$$

Здесь S_1, \dots, S_m — описания объектов, составляющих обучающий материал, $\tilde{\alpha}(S_i)$, $i = 1, 2, \dots, m$, — информационные векторы объектов S_i по свойствам $P_j \equiv S_i \in K_j$. Другими словами, если $\alpha(S_i) = (\alpha_{i1}, \dots, \alpha_{ij}, \alpha_{il})$, $j = 1, 2, \dots, l$, то

$$\alpha_{ij} = \begin{cases} 1, & S_i \in K_j, \\ 0, & S_i \notin K_j. \end{cases}$$

Задача распознавания Z определяется начальной информацией I_0 и конечной выборкой $\tilde{S}^q = (S^1, \dots, S^q)$, $S^i = (b_{i1}, \dots, b_{in})$, $i = 1, 2, \dots, q$, т.е. $Z = \{I_0, \tilde{S}^q\}$.

Требуется для каждого объекта S^i из \tilde{S}^q вычислить его информационный вектор $\tilde{\beta}(S^i)$ или, что то же самое, значение свойства $P_j(S^i) \equiv S^i \in K_j$, $j = 1, 2, \dots, l$.

Далее будем считать, что информационные векторы для \tilde{S}^q известны, и, основываясь на этом, строить алгоритм, который правильно вычисляет эти свойства.

- I. Пусть дано множество $\{A\}$, вообще говоря, некорректных алгоритмов для решения задач распознавания, представленных в виде $A = B \cdot C$, где B — распознающий оператор, C — решающее правило. Напомним, что $B(I_0, S^q) = \|\Gamma_{ij}\|_{q \times l}$. Здесь Γ_{ij} — действительные числа, $C(\|\Gamma_{ij}\|_{q \times l}) = \|\beta_{ij}\|_{q \times l}$; $\beta_{ij} \in \{0, 1, \Delta\}$, $\beta_{ij} = \Delta$ означает, что алгоритм A отказался от вычисления свойства $P_j(S^i)$; $\beta_{ij} = \beta$, $\beta \in \{0, 1\}$, означает, что алгоритм A вычислил свойство $P_j(S^i)$ равным β . При этом мы допускаем возможность ошибки.

Известно, что каждый алгоритм A может быть представлен в виде $B \cdot C$ и что по исходному семейству $\{A\}$ с помощью операций сложения, умножения и умножения на скаляр можно построить алгебраическое расширение

$$\mathfrak{U}\{A\} = \mathfrak{U}\{B\} \cdot \{C\}$$

класса алгоритмов $\{A\}$ и алгебраические расширения конечных степеней

$$\mathfrak{U}^k\{A\} = \mathfrak{U}^k\{B\} \cdot \{C\}.$$

При выполнении некоторых условий для задачи Z и исходного семейства $\{A\}$ в расширении $\mathfrak{U}^k\{A\}$ можно построить алгоритм A^* , правильно вычисляющий все значения $P_j(S^i)$, $i = 1, 2, \dots, q$, $j = 1, 2, \dots, l$. Если в качестве исходного семейства $\{A\}$ рассмотреть класс алгоритмов вычисления оценок, то искомым алгоритм A^* представим в виде

$$A^* = \left[(c_1 + c_2) \sum_{i=1}^q \sum_{j=1}^l \beta_{ij} (B_{ij})^k \right] C(c_1, c_2),$$

где c_1, c_2 — параметры решающего правила C ,

$$\begin{aligned} \beta_{ij} P_j(S^i), \quad i = 1, 2, \dots, q, \quad j = 1, 2, \dots, l, \\ B_{ij} = B_j + B_i^j, \\ B_j = B_{j1} + \dots + B_{j,j-1} + B_{j,j+1} + \dots + B_{jl}, \\ B_i^j = B_{i1}^j + \dots + B_{i,i-1}^j + B_{i,i+1}^j + \dots + B_{iq}^j; \end{aligned}$$

здесь каждый оператор B_{ij} является оператором вычисления оценок, каждый оператор B_{iv}^j является либо оператором вычисления оценок, либо разностью двух операторов вычисления оценок. Для величины k также имеется формула, приводить которую здесь нет необходимости.

Каждый оператор вычисления оценок кодируется значениями $2n + 3m + 3$ параметров, где n — число признаков, описывающих объекты, m — число объектов в I_0 . Нетрудно также видеть, что приведенная выше формула для алгоритма A^* включает в себя по крайней мере $l(l-1) + lq(q-1)$ операторов вычисления оценок. Следовательно, для полной записи кода алгоритма A^* требуется по крайней мере $(2n + m + 3)l[(l-1) + q(q-1)]$ чисел. Указанная величина может быть несколько уменьшена с помощью специальных приемов, однако она все-таки остается большой, и это неудобно при машинной реализации алгоритма, если величины n, m, l, q велики. Поэтому при реальном синтезе корректного алгоритма A^* будет использоваться только его принципиальная запись, данная выше, а реализация операторов типа B_{ij} будет проводиться другими методами. В дальнейшем будут использоваться только пороговые решающие правила $C(c_1, c_2)$: $C(\|\Gamma_{ij}\|_{q \times l}) = \|C(\Gamma_{ij})\|_{q \times l}$,

$$C(\Gamma_{ij}) = \begin{cases} 1, & \Gamma_{ij} > c_2, \\ 0, & \Gamma_{ij} < c_1, \\ \Delta, & c_1 \leq \Gamma_{ij} \leq c_2 \end{cases} \quad 0 < c_1 < c_2.$$

II. Рассмотрим информационную матрицу $\|\beta_{ij}\|_{q \times l}$ выборки \tilde{S}^q в задаче Z . Положим

$$\begin{aligned} M = \{(i, j)\}, \quad i = 1, 2, \dots, q, \quad j = 1, 2, \dots, l, \\ M_\alpha = \{(i, j)\} : \beta_{ij} = \alpha, \quad \alpha = 0, 1. \end{aligned}$$

Очевидно, $M = M_0 \cup M_1$.

Пусть B — распознающий оператор и

$$B(Z) = \|\Gamma_{rt}(B)\|_{q \times l},$$

где Γ_{rt} — действительные числа.

Определение 1 Оператор B называется допустимым для задачи Z , если существует хотя бы одна пара (u, v) из M_1 такая, что для всех (i, j) из M_0

$$\Gamma_{uv}(B) > |\Gamma_{ij}(B)|.$$

Пара (u, v) называется в этом случае отмеченной в B . Совокупность всех пар, отмеченных в B , обозначим через $M(B)$.

Пусть

$$\begin{aligned} \Gamma_{max}^0(B) &= \max_{(i,j) \in M_0} |\Gamma_{ij}(B)|, \\ \Gamma_{min}^1(B) &= \min_{(i,j) \in M(B)} \Gamma_{ij}(B). \end{aligned}$$

Положим

$$\Gamma(B) = [\Gamma_{min}^1(B)]^{-1}. \quad (5.1)$$

По оператору B построим оператор B' :

$$B' = \Gamma(B) \cdot (B). \quad (5.2)$$

Пусть $B'(Z) = \|\Gamma'_{ij}\|_{q \times l}$. Тогда из (5.2) легко следует, что

$$\Gamma'_{ij} = \Gamma(B) \cdot \Gamma_{ij}(B). \quad (5.3)$$

Лемма 1 Если (u, v) отмечена в B , то $\Gamma'_{uv}(B) \geq 1$; если (u, v) не отмечена в B , то

$$\Gamma'_{uv}(B) \leq \frac{\Gamma_{max}^0(B)}{\Gamma_{min}^1(B)} = \Gamma(B) < 1.$$

Доказательство. Если (u, v) отмечена в B , то для $\Gamma_{uv}(B)$ выполнено неравенство

$$\Gamma_{uv}(B) \geq \min_{(i,j) \in M(B)} \Gamma_{ij}(B) = \Gamma_{min}^1(B).$$

Из этого неравенства и соотношений (5.1)-(5.3) легко следует первое утверждение леммы.

Если пара (u, v) не является отмеченной в B , то

$$|\Gamma_{uv}(B)| \leq \max_{(i,j) \in M_0} \Gamma_{ij}(B) = \Gamma_{max}^0 < \min_{(i,j) \in M(B)} \Gamma_{ij}(B) = \Gamma_{min}^1(B).$$

Из последних неравенств и соотношений (5.1)-(5.3) легко следует второе утверждение леммы.

В дальнейшем положим

$$\Gamma_{max}^0(B)/\Gamma_{min}^1(B) = Q(B).$$

Пусть $\{B\}$ — произвольная конечная система распознающих операторов.

Определение 2 Система $\{B\}$ называется базисной для Z , если

$$M_1 = \bigcup_{B \in \{B\}} M(B).$$

По базисной для Z системе $\{B\}$ построим алгоритм A^* , корректный для задачи Z . Введем для операторов B' , $B' = \Gamma(B) \cdot B$, $B \in \{B\}$ целые числа $k(B')$ так, чтобы выполнялось неравенство

$$[Q(B)]^{k(B')} < \frac{c_1}{(c_1 + c_2)|\{B\}|}. \quad (5.4)$$

Для этого достаточно положить

$$k(B') = \frac{\ln(c_1 + c_2) + \ln|\{B\}| - \ln c_1}{|\ln Q(B)|} + 1.$$

Очевидно, что в этом случае неравенство (5.4) выполнено. Напомним, что величины c_1, c_2 суть пороги решающего правила $C(c_1, c_2)$.

Теорема 5 Алгоритм

$$A = \left\{ (c_1 + c_2) \sum_{B \in \{B\}} (B')^{k(B')} \right\} \cdot C(c_1, c_2)$$

является корректным для Z .

Доказательство. Распознающим оператором B^* в алгоритме A^* является оператор

$$(c_1 + c_2) \sum_{B \in \{B\}} (B')^{k(B')}.$$

Пусть $B^*(Z) = \|\Gamma_{ij}^*\|_{q \times l}$. Тогда

$$\Gamma_{ij}^* = (c_1 + c_2) \sum_{B \in \{B\}} (\Gamma'_{ij}(B))^{k(B')}.$$

Случай 1. $P_j(S^i) = \beta_{ij} = 1, 1 \leq i \leq q, 1 \leq j \leq l,$

$$\Gamma_{ij}^* = (c_1 + c_2) \left[\sum_{B \in B(i,j)} (\Gamma'_{ij}(B))^{k(B')} + \sum_{B \notin B(i,j)} (\Gamma'_{ij}(B))^{k(B')} \right].$$

Здесь $B(i, j)$ — совокупность всех операторов из $\{B\}$, в которых пара (i, j) является отмеченной. Так как $\{B\}$ — базисная система для Z , то множество $B(i, j)$ непусто. Поэтому

$$\sum_{B \in B(i,j)} (\Gamma'_{ij}(B))^{k(B')} \geq 1, \quad (5.5)$$

$$\left| \sum_{B \in B(i,j)} (\Gamma'_{ij}(B))^{k(B')} \right| < |\{B\}| \frac{c_1}{(c_1 + c_2)|\{B\}|}. \quad (5.6)$$

Неравенство (5.5) следует из определения отмеченной пары, неравенство (5.6) — из (5.4). Окончательно получаем

$$\Gamma_{ij}^* > (c_1 + c_2) \left(1 - \frac{c_1}{c_1 + c_2} \right) = c_2.$$

Но тогда из определения порогового решающего правила следует $C(\Gamma_{ij}^*) = 1 = \beta_{ij} = p_j(S^j)$.

Случай 2. $p_j(S^i) = \beta_{ij} = 0, 0 \leq i \leq q, 1 \leq j \leq l.$ В этом случае пара (i, j) не является отмеченной ни в одном операторе B . Поэтому

$$\Gamma_{ij}^* = (c_1 + c_2) \sum_{B \in B(i,j)} (\Gamma'_{ij}(B))^{k(B')} < (c_1 + c_2)|\{B\}| \times \frac{c_1}{(c_1 + c_2)|\{B\}|} = c_1.$$

Из определения $C(c_1, c_2)$ следует

$$C(\Gamma_{ij}^*) = 0 = \beta_{ij} = P_j(S^i).$$

Теорема доказана.

Определение 3 *Базисная система $\{B\}$ называется неприводимой для Z , если никакая собственная часть $\{B\}$ не является базисной для Z .*

Очевидно,

$$|\{B\}| \leq |M_1| \leq q \times l.$$

Поэтому неравенство (5.4) для неприводимых систем записывается в виде

$$Q(B)^{k(B')} < \frac{c_1}{(c_1 + c_2)|M_1|}.$$

Этому неравенству удовлетворяет

$$k(B') = \frac{\ln |M_1| + \ln(c_1 + c_2) - \ln c_1}{\ln Q(B)} + 1.$$

Формулировка теоремы 1 при новых $k(B')$ сохранится без изменений.

Из теоремы 1 следует, что для построения эффективного корректного алгоритма достаточно построить систему из небольшого числа просто выполняемых распознающих операторов B , базисную для Z . Из базисной системы затем нетрудно получить неприводимую систему. Для построения базисной системы подходит любой оператор, отмечающий непустое множество пар (i, j) таких, что $\beta_{ij} = 1$. В $[X]$ для любой задачи $Z = \{I_0, \tilde{S}^q\}$, $I_0 = \{S_1, \dots, S_m, \tilde{\alpha}(S_1), \dots, \tilde{\alpha}(S_m)\}$, в которой объекты из \tilde{S}^q попарно неизоморфны относительно I_0 и информационная матрица $\|\alpha_{ij}\|_{m \times l}$ в I_0 состоит из попарно различных столбцов, указана базисная система из $|M_1|$ операторов. Каждый оператор этой базисной системы гарантирует отметку, вообще говоря, ровно одной пары (i, j) , $\beta_{ij} = 1$. Можно указать случаи, когда один оператор гарантирует отметку существенно большего числа пар. Это возможно и для простых распознающих операторов.

Литература

- [1] Труды Мат. ин-та им. В.А.Стеклова, том 51, 1958 г.
- [2] Ю.И. Журавлев. Избранные научные труды. Москва, из-во Магистр, 1996 г., стр. 378–384
- [3] Ю.И. Журавлев, И.В. Исаев. Построение алгоритмов распознавания, корректных для заданной контрольной выборки. Журнал вычислительной математики и математической физики, том 19. №3, 1979 г.