

Data mining: задачи, техники и приложения

Николай Пономарев

2 мая 2024 г.

Дисклеймер: данный текст просто пересобранная презентация, саму статью можно найти в конце

1. Введение

Немного о терминологии

Две близких области: Information Retrieval (IR) и Data Mining (DM)

Определение 1. IR — наука об организации информации и алгоритмах её *быстрого* получения

Определение 2. DM — наука о методах обнаружения в данных ранее *неизвестных* и практически полезных знаний

Формально, наш курс про IR

2. Задачи

Задачи

Определение 3. Классификация (classification) — процесс разделения новых наблюдений (observation) на предопределенные классы

Определение 4. Кластеризация (clustering) — процесс разбиения данных (или наблюдений о них) на группы

Определение 5. Анализ выбросов (outlier analysis) — способ извлечения полезных знаний из выбросов

Определение 6. Ассоциативный анализ (association analysis) — процесс поиска ассоциаций среди данных, которые удовлетворяют определенным статистическим требованиям

3. Техники

Статистические подходы

Строится статистическая модель, а затем анализируется следующими способами:

Байесовские сети способ выяснения зависимостей между переменными «с помощью» формулы Байеса

Корреляция используется для установления зависимости между фактами

Регрессия установление соответствия между случайными переменными отражающими связь между зависимой переменной и независимыми переменными x

Факторный анализ используется для поиска основных источников корреляции

Машинное обучение и нейронные сети

- Развивает идеи статистических методов
- В каком-то смысле автоматизирует анализ
- Часто результаты лучше
- Модно

Ещё техники

- СУБД
- Генетические алгоритмы
- Нечёткие множества
- Визуализация

4. Приложения

Телекоммуникации

Телеком и мобильные операторы используют data mining для

- Маркетинга
 - Классификация и кластеризация \Rightarrow таргетированная реклама
- Удержания клиентов
 - Классификация и кластеризация \Rightarrow обнаружение недовольных клиентов
- Создание оптимальных тарифов
- Оптимизация использования инфраструктуры

Продажи

Торговым предприятиям нужен ДМ для исследования

- Поведения покупателей
 - Ассоциативный анализ
- Корзины
- Выбора товаров
- Расстановки продуктов на полках
 - Кластеризация
- Влияния акций

Медицина

В медицине ДМ используется для

- Обнаружения и анализа хронических заболеваний
- Поиск эффективных лекарств
- Отслеживать вероятность эпидемий

Зачем ещё это надо?

- Финансы
- Предотвращение преступлений
- Рекомендательные системы
- Реклама

Источник

Список литературы

- [1] M. K. Gupta & P. Chandra A comprehensive survey of data mining International Journal of Information Technology, 12 (2020) 1243–1257 <https://doi.org/10.1007/s41870-020-00427-7>