

Машинное обучение

Лекция 1

Катя Тузова

JetBrains

Правила игры

- 13 опросов по 5 баллов в начале лекции
- 13 домашних заданий по 20 баллов при сдаче в первую неделю, 10 баллов при сдаче во вторую неделю
- 1 соревнование на платформе kaggle.com - 40 баллов
- Допуск к экзамену = 150 баллов
- Экзамен = 3 вопроса
- За каждые 80 баллов сверх 150 – минус вопрос на экзамене

Что такое машинное обучение?

Что такое машинное обучение?

(Wikipedia)

Machine learning is a scientific discipline that explores the construction and study of algorithms that can learn from data. Such algorithms operate by building a model based on inputs and using that to make predictions or decisions, rather than following only explicitly programmed instructions.

Обширный подраздел искусственного интеллекта, изучающий методы построения моделей, способных обучаться, и алгоритмов для их построения и обучения.

В чем отличия ML от AI?

Разделы AI

- Работа с естественными языками
- Представление и использование знаний
- Компьютерное зрение
- Машинное обучение
- Биологическое моделирование искусственного интеллекта
- Робототехника
- Анализ графов
- ...

Что такое машинное обучение?

Arthur Samuel (1959).

Field of study that gives computers the ability to learn without being explicitly programmed.

Пример

Чем отличается задача “найти кратчайший путь в графе” от “антиспам фильтр”?

Что такое машинное обучение?

Tom Mitchell (1998)

Well-posed Learning Problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E

E.g., Learn to play checkers

T : Play checkers

P : % of games won in world tournament

E : opportunity to play against self

Где используется ML?

Где используется ML?

- Рекомендации на Amazon, Kinopoisk
- Поисковые системы Google, Яндекс
- Выделение лиц друзей на фото Facebook
- Боты в Twitter

Какая математика понадобится?

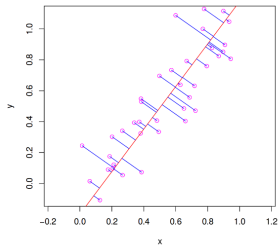
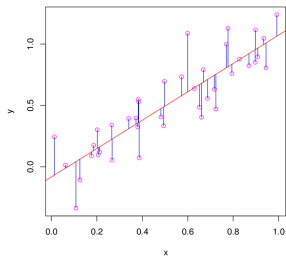
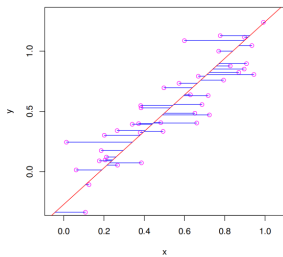
Какая математика понадобится?

- Математическая статистика
- Методы оптимизации
- Дискретная математика
- Теория вероятности
- Линейная алгебра

Вопрос

Кто из вас уже использовал методы машинного обучения?

Метод наименьших квадратов



Метод наименьших квадратов

Идея: Аппроксимировать данные линейной зависимостью

Суть: Минимизация суммы квадратов отклонений данных от прямой

Вопрос

Достаточно ли знать алгоритмы ML и математику?

Что еще надо понимать

- Когда надо применять ML
- Как сформулировать задачу в терминах ML
- Как выбрать подходящий класс алгоритмов
- Где посмотреть существующие решения
- Как настроить алгоритм
- Как оценить результаты

Тест Тьюринга

AI как наука начался с теста Тьюринга (1950).

Компьютер должен успешно выдать себя за человека в (письменном) диалоге между судьёй, человеком и компьютером

История ML

- 50-70 гг. – Базы знаний, полнотекстовый поиск, распознавание образов, нейронные сети, К ближайших соседей.
- 1973г. – “Зима” искусственного интеллекта
- 80-90гг. – Первые конференции, развитие практического применения
- 90-00гг. – Метод опорных векторов, бустинг
- Рандомизированный решающий лес (начало 2000х)
- Обучение Марковских случайных полей (2000е)
- Глубокое обучение (2006)

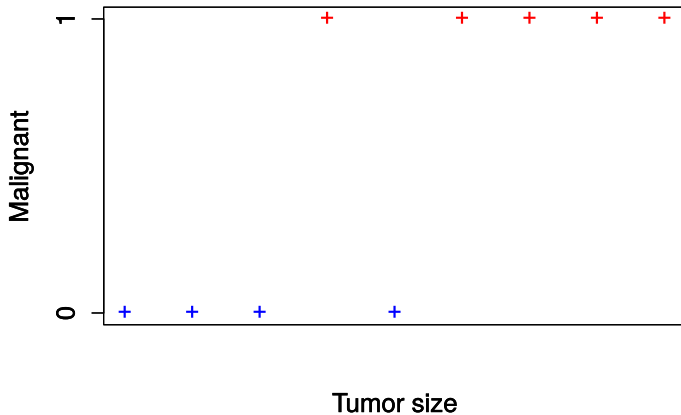
- Проблема комбинаторного взрыва
- Низкая производительность компьютеров
- Проблема представлений знаний “здорового мысла”
- Парадокс Моравца

Типы машинного обучения

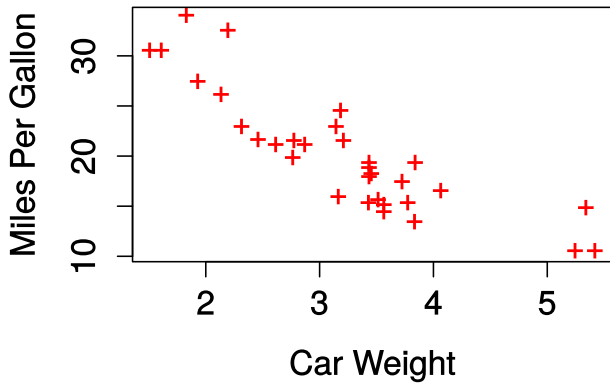
Типы машинного обучения

- С учителем
- Без учителя
- Смешанное

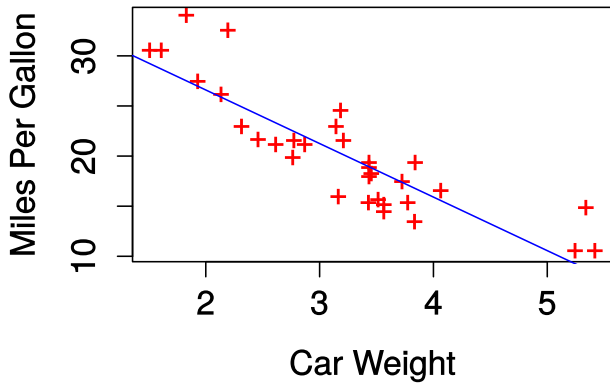
Обучение с учителем (Классификация)



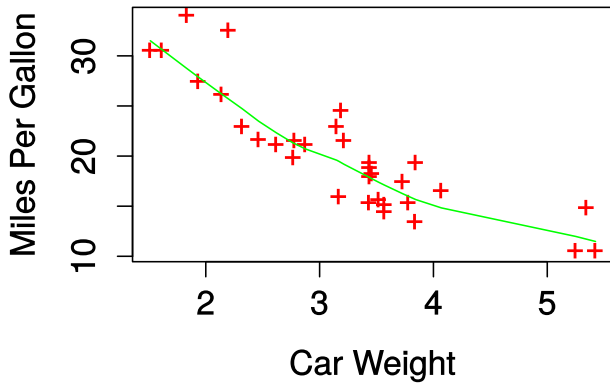
Обучение с учителем (Регрессия)



Обучение с учителем (Регрессия)



Обучение с учителем (Регрессия)



Обучение с учителем

Множество объектов X

Множество допустимых ответов Y

Прецедент - пара объект-ответ (x_i, y_i)

Обучающая выборка - совокупность пар $X_l = (x_i, y_i)_{i=1}^l$

Целевая функция $y^* : X \rightarrow Y$

Задача обучения по прецедентам:

Найти решающую функцию $a : X \rightarrow Y$

Решающая функция должна приближать целевую на всем множестве X

Пример. Метод наименьших квадратов

Линейная модель:

$$y(x, w) = \sum_{i=1}^l x_i w_i$$

Целевая функция?

Решающая функция?

Пример. Метод наименьших квадратов

Линейная модель:

$$y(x, w) = \sum_{i=1}^l x_i w_i$$

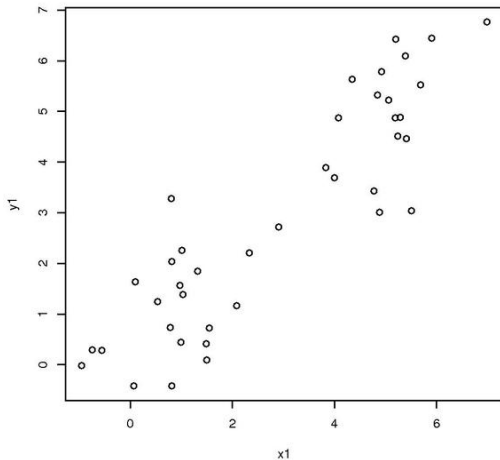
Целевая функция: $y^* = \sum_{i=1}^l (y_i - y(x, w))^2$

Решающая функция - итог оптимизации: $\arg \min_w y^*$

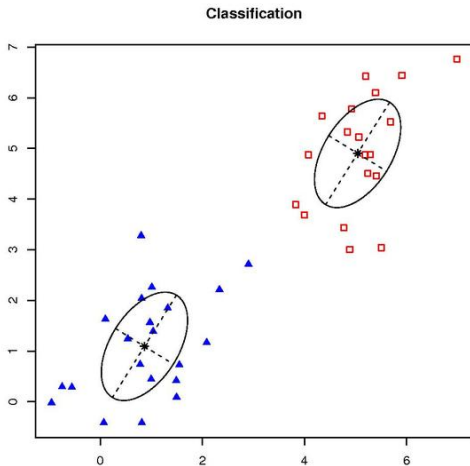
Обучение с учителем

- Классификация ($Y = 1, 2, \dots, K$ конечно) — множество X разбивается на K классов. Требуется предсказать к какому классу он принадлежит.
- Восстановление регрессии ($Y = \mathbb{R}$) — требуется найти функцию f из определенного класса, которая аппроксимирует f^*

Обучение без учителя (Пример)



Обучение без учителя (Пример)



Типы признаков

$$f : X \rightarrow D_f$$

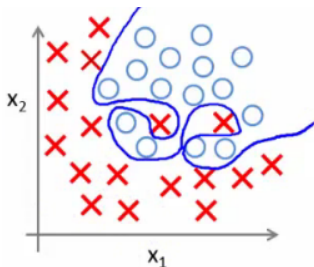
- Бинарные ($D_f = \{0, 1\}$)
- Номинальные (D_f – конечное множество)
- Порядковые (D_f – конечное упорядоченное множество)
- Количественные ($D_f = \mathbb{R}$)

Типы признаков (примеры)

- Бинарные (Пол, наличие боли в спине, в сознании ли пациент)
- Номинальные (Тип боли: колющая, режущая, ноющая)
- Порядковые (Общее состояние больного: удовлетворительное, средней тяжести, тяжелое, крайне тяжелое)
- Количественные (Температура тела, пульс, артериальное давление)

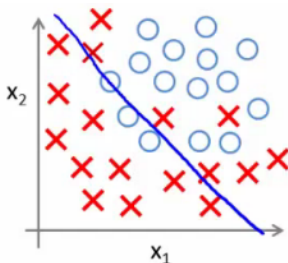
Определения

Переобучение (overfitting) это явление, при котором алгоритм слишком приспособлен для данных, на которых он обучался. Переобучение имеет место при выборе слишком сложных моделей.



Определения

Недообучение (underfitting) это явление, обратное переобучению, при котором алгоритм не полностью использует предоставленные ему для обучения данные. Недообучение имеет место при выборе недостаточно сложных моделей.



Применение ML

- Академическое
Красивые идеи, хорошая математика
- Практическое
Обеспечивает некоторое качество на множестве примеров

Что читать/смотреть

- G. James, D. Witten, T. Hastie, R. Tibshirani: "An Introduction to Statistical Learning"
- Christopher M. Bishop "Pattern Recognition and Machine Learning"
- Kevin P. Murphy "Machine Learning A Probabilistic Perspective"

- К.В. Воронцов. http://shad.yandex.ru/lectures/machine_learning.xml
- Andrew Ng <http://ml-class.org/>

На следующей лекции

- Метод ближайших соседей
- Гипотеза компактности
- Обобщенный метрический классификатор
- Проклятие размерности
- Отбор эталонов
- Автоматический выбор фич