

Семинар 3. Постановка задачи машинного обучения

Надежда Чиркова
nchirkova@hse.ru, @nadiinchi (Telegram)

Что такое машинное обучение?

Машинное обучение — «обучение с помощью машины»?

Machine learning — «обучение машины»

Определение с сайта machinelearning.ru:

Машинное обучение (Machine Learning) — обширный подраздел искусственного интеллекта, изучающий методы построения алгоритмов, способных обучаться.

Зачем нужно машинное обучение?

- Заменить человека при решении задач (автоматизация);
- Поиск закономерностей в данных, которые человек не находит.

Постановка задачи машинного обучения

Задача машинного обучения:

- данные (что такое объект, какие признаки + типы признаков);
- что предсказывать;
- оценивание качества (критерий качества + способ валидации).

Матрица объекты–признаки

Числовая матрица:

	Признак 1	Признак 2	...	Признак К
Объект 1				
Объект 2				
Объект 3				
...				
Объект N				

Виды признаков

Объект — вектор в [конечномерном] пространстве признаков.

Виды признаков:

- 1 вещественные;
- 2 бинарные;
- 3 категориальные;
- 4 порядковые (упорядоченные категориальные);
- 5 подмножество супермножества;
- 6 строковые.

Что предсказываем?

Два типа обучения:

- Обучение с учителем (пытаемся понять, как зависят ответы, известные на объектах обучающей выборки, от входных данных):
 - Классификация (бинарная, multiclass, multilabel)
 - Регрессия
 - Прогнозирование временных рядов
 - Рекомендации
 - ...
- Обучение без учителя (как можем формализуем, что хотим найти в данных, и ищем).
 - Кластеризация
 - Понижение размерности
 - Визуализация
 - ...

Критерии качества

y_i — правильный ответ на i -м объекте
 $a(x_i)$ — предсказанный ответ на i -м объекте
 ℓ — число объектов в выборке

$$Q(a) = \frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a(x_i))$$

Критерии качества

y_i — правильный ответ на i -м объекте
 $a(x_i)$ — предсказанный ответ на i -м объекте
 ℓ — число объектов в выборке

$$Q(a) = \frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a(x_i))$$

Пример:

$$L(y_i, a(x_i)) = \begin{cases} 1, & y_i = a(x_i) \\ 0, & y_i \neq a(x_i) \end{cases} \quad (\text{accuracy})$$

Критерий качества

- фантазия не ограничена :)
- определяется заказчиком исходя из цели решения задачи
- должен легко вычисляться по имеющимся данным в offline

Критерий качества

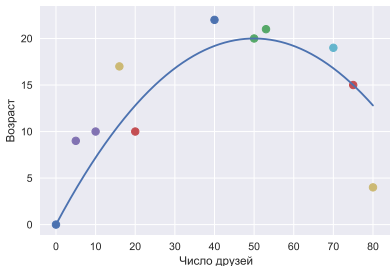
- фантазия не ограничена :)
- определяется заказчиком исходя из цели решения задачи
- должен легко вычисляться по имеющимся данным в offline

По какой выборке измеряется качество?

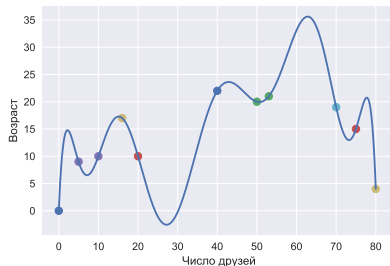
Качество на обучающей выборке

По оси абсцисс — признак, по оси ординат — целевая переменная.

Хорошая модель



Переобучение



Качество нужно оценивать по отдельной выборке!

Разделение выборки

Качество нужно оценивать по отдельной выборке!

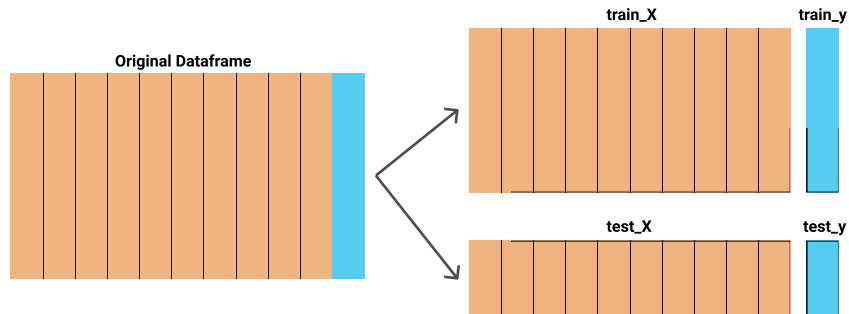
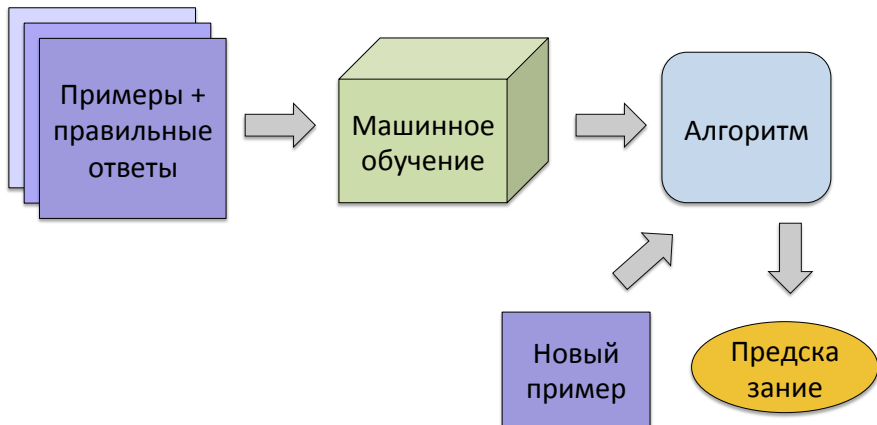


Схема построения алгоритма



Общая схема работы

- Предобработка данных и составление матрицы объекты–признаки
- Повторять, пока не устроит качество решения:
 - Построение модели
 - Оценивание качества

Простой пример

Задача кредитного скоринга: вернет ли заемщик кредит
Постановка задачи?

Простой пример

Задача кредитного скоринга: вернет ли заемщик кредит
Постановка задачи?

- данные:
 - объект?
 - признаки (с видом признака)?

Простой пример

Задача кредитного скоринга: вернет ли заемщик кредит
Постановка задачи?

- данные:
 - объект?
 - признаки (с видом признака)?
- что предсказывать (с типом задачи)?
- критерий качества?
- метод валидации решения?
- где взять данные?

Выводы

- Формальная постановка задачи — важный процесс, перевод задачи с языка прикладной области на математический язык методов решения.
- Не всегда очевидно, что является объектом, где взять данные, какой критерий качества выбрать...
- Хорошо поставленную задачу проще решать :)