

ML Handbook

Сергей Полянских

Оглавление

Предисловие

В данной книге описаны основные понятия, методы и подходы, широко используемые в современном DS и ML. Обычно, свободное владение этими понятиями необходимо для правильного понимания как основных, так и продвинутых методов ML и по умолчанию предполагается от DS специалиста.

Здесь собраны разные определения, встречавшиеся автору в научных статьях по ML и на собеседованиях. Охвачены: теория вероятностей, классическая и байесовская статистика, некоторые вопросы мат. анализа.

Освещение вопросов ни в коем случае не претендует на полноту. Основная цель книги - составить расширенный глоссарий основных понятий и подходов, встретившихся автору в процессе работы в области ML.

Обозначения

DS - дата саенс
ML - машинное обучение
RV - случайная величина

Глава 1

Математика

В этой главе описаны основные математические понятия, необходимые для правильного понимания как основных, так и продвинутых методов ML. Охвачены: теория вероятностей, классическая и байесовская статистика, некоторые вопросы мат. анализа. Освещение вопросов ни в коем случае не претендует на полноту. Основная цель - составить расширенный глоссарий основных понятий и подходов, встретившихся автору в процессе работы в области ML.

1.1 Случайная величина

Случайной величиной (RV) называется числовая функция X , определенная на некотором множестве элементарных исходов Ω (обычно подмножество \mathbb{R} или \mathbb{R}^n),

$$X : \Omega \rightarrow \mathbb{R}.$$

С прикладной точки зрения на RV обычно смотрят как на генераторы случайных чисел с заданным распределением.

Примеры:

- Рост людей, взятых из некоторой группы.
- Цвет фиксированного пикселя изображения, взятого из некоторого множества изображений.
- Некоторый признак из датасета ML задачи.

1.2 Распределение случайной величины

Если RV принимает дискретное множество значений x_1, x_2, \dots , то она полностью определяется значениями вероятностей: $p_k = \mathbb{P}(X = x_k)$.

Если множество значений RV не дискретно, то RV может быть описана своей функцией распределения (CDF, Cumulative Distribution Function): $F(x) = \mathbb{P}(X < x)$

В DS в большинстве случаев CDF дифференцируемо. Производная от CDF называется плотностью распределения случайной величины: $f(x) = F'(x)$. Таким образом, по определению

$$\mathbb{P}(a < X < b) = \int_a^b f(x)dx$$

1.3 Выборка

1.4 Закон больших чисел

1.5 Классический и байесовский подход