

ML Handbook

Сергей Полянских

Оглавление

1	Математика	5
1.1	Случайная величина	5
1.2	Распределение случайной величины	5
1.3	Выборка	6
1.4	Закон больших чисел	6
1.5	Центральная предельная теорема	6
1.6	Статистики	6
1.7	Bootstrap	7
1.8	Классический и байесовский подход	7
1.9	Метод максимального правдоподобия	7
1.10	Доверительный интервал	7
1.11	Байесовский доверительный интервал	7
1.12	Основные дискретные распределения	7
1.13	Основные непрерывные распределения	7
1.14	Матричные разложения	7
1.15	КЛ дивергенция	8
1.16	Энтропия	8
1.17	Квантили	8
1.18	Точечные оценки	8
1.19	Интервальные оценки	8
1.20	Проверка гипотез	8
1.21	Множественная проверка гипотез	8
1.22	Параметрические и непараметрические критерии, бутстреп	8
1.23	Ошибки I и II рода	8
1.24	Достигаемый уровень значимости	8
1.25	Мощность статистического критерия	8
1.26	Основные задачи статистики	8
1.27	Проверка основных гипотез	9
1.28	Корреляция Пирсона	9
1.29	Корреляция Спирмена	9
1.30	Корреляция Метьюса	9
1.31	Корреляция Крамера	9
1.32	Z-тест Фишера	9
1.33	T-тест Стьюдента	9

1.34	Критерий Пирсона χ^2	9
1.35	Точный тест Фишера	9
2	Анализ данных	10
2.1	Типы данных	10
2.2	Предобработка данных	10
2.3	Понижение размерности	10
3	Общие вопросы	11
3.1	Машинное обучение	11
3.2	Основные классы задач	11
3.3	Обнаружение аномалий	11
3.4	Контроль качества	11
3.5	Недообучение	12
3.6	Переобучение	12
3.7	Регуляризация	12
3.8	Отбор признаков	12
3.9	Параметры алгоритма	12
3.10	Подбора метапараметров	12
3.11	Основные типы алгоритмов	12
3.12	Многоклассовая классификация	12
3.13	Дисбаланс классов	12
3.14	Ансамбли алгоритмов	12
3.15	Метрики классификации	12
3.16	ROC-AUC метрика	12
3.17	Метрики регрессии	12
3.18	Метрики кластеризации	12
3.19	Разложение ошибки алгоритма	12
3.20	Кривые валидации	12
3.21	Кривые обучения	12
3.22	Метрические методы	12
3.23	Метод ближайших соседей	12
3.24	Линейные методы	12
3.25	Линейная регрессия	12
3.26	Логистическая регрессия	12
3.27	SVM	13
3.28	Ядра и спрямляющие пространства	13
3.29	Решающие деревья	13
3.30	Случайный лес	13
3.31	Градиентный бустинг	13
3.32	Байесовские методы	13
4	Нейросети	14

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 и в некоторых случаях на строгость. Основная цель книги - составить расширенный глоссарий основных понятий и подходов, встретившихся автору в процессе работы в области ML.

Обозначения

DS	- наука о данных
ML	- машинное обучение
RV	- случайная величина
CDF	- функция распределения случайной величины
PDF	- плотность распределения случайной величины
CLT	- центральная предельная теорема
EX	- среднее случайной величины X
DX	- дисперсия случайной величины X
$X \sim Y$	- случайные величины X и Y одинаково распределены

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.

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 часто 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 их 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 большинстве прикладных случаев CDF оказывается дифференцируемой функцией. Производная от CDF называется плотностью распределения случайной величины (PDF, Probability density function): $f(x) = F'(x)$. Таким образом, по определению

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

1.3 Выборка

Выборкой объема n из генеральной совокупности X называется последовательность независимых и распределенных как X случайных величин:

$$X_1, X_2, \dots, X_n, \quad X_k \sim X$$

На практике под выборкой понимают конкретные реализации величин X_k , то есть последовательность чисел x_1, x_2, \dots, x_n .

1.4 Закон больших чисел

Закон больших чисел утверждает, что если X_1, X_2, \dots, X_n - выборка объема n из генеральной совокупности X , то ее среднее с ростом n стабилизируется к среднему значению X :

$$\frac{X_1 + X_2 + \dots + X_n}{n} \approx EX, \quad n \rightarrow \infty.$$

1.5 Центральная предельная теорема

Центральная предельная теорема (CLT) является в некотором смысле уточнением закона больших чисел. В упрощенном варианте она утверждает, что если X_1, X_2, \dots, X_n - выборка объема n из генеральной совокупности X , то ее распределение ее среднего при больших n очень близко к нормальному,

$$\frac{X_1 + X_2 + \dots + X_n}{n} \approx N(\mu, \sigma^2/n), \quad \mu = EX, \sigma^2 = DX, \quad n \rightarrow \infty.$$

Заметим, что если совокупность распределена нормально, $X \sim N(\mu, \sigma^2)$, то предыдущая формула обращается в точное равенство при любых n .

1.6 Статистики

Пусть X_1, X_2, \dots, X_n - выборка объема n . Статистикой называется произвольная RV, являющаяся функцией выборки:

$$T = T(X_1, X_2, \dots, X_n).$$

Часто статистикой выборки называют конкретное значение $T(x_1, x_2, \dots, x_n)$, полученное на данной реализации x_1, x_2, \dots, x_n выборки.

Примеры:

- $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$ - выборочное среднее.
- $X_{(n)} = \max(X_1, X_2, \dots, X_n)$ - максимальное значение в выборке.
- медиана, перцентили.

1.7 Bootstrap

1.8 Классический и байесовский подход

1.9 Метод максимального правдоподобия

1.10 Доверительный интервал

1.11 Байесовский доверительный интервал

1.12 Основные дискретные распределения

1.13 Основные непрерывные распределения

1.14 Матричные разложения

...может разделить главу на части...

- 1.15 КЛ дивергенция
- 1.16 Энтропия
- 1.17 Квантили
- 1.18 Точечные оценки
- 1.19 Интервальные оценки
- 1.20 Проверка гипотез
- 1.21 Множественная проверка гипотез
- 1.22 Параметрические и непараметрические критерии, бутстреп
- 1.23 Ошибки I и II рода
- 1.24 Достигаемый уровень значимости
- 1.25 Мощность статистического критерия
- 1.26 Основные задачи статистики

...из лекций новосиба курсера...

- 1.27 Проверка основных гипотез
- 1.28 Корреляция Пирсона
- 1.29 Корреляция Спирмена
- 1.30 Корреляция Метьюса
- 1.31 Корреляция Крамера
- 1.32 Z-тест Фишера
- 1.33 Т-тест Стьюдента
- 1.34 Критерий Пирсона χ^2
- 1.35 Точный тест Фишера

Глава 2

Анализ данных

Анализ и предобработка данных - первая задача, успешное решение которой зачастую определяет успех в решении любых задач ML. В этой главе описываются основные подходы....

2.1 Типы данных

2.2 Предобработка данных

2.3 Понижение размерности

Глава 3

Общие вопросы

В этой главе приводятся основные понятия ML и DS.

3.1 Машинное обучение

Машинное обучение (ML) - область искусственного интеллекта, изучающая самообучающиеся модели, то есть решающие поставленную задачу не по заранее запрограммированному алгоритму, а предварительно настраивая свое поведение согласно имеющимся данным.

Обычно методы ML содержат свободные параметры, подбор которых наилучшим (в смысле имеющихся данных) образом и составляет процесс обучения алгоритма.

3.2 Основные классы задач

3.3 Обнаружение аномалий

3.4 Контроль качества

...оценка обобщающей способности...

- 3.5 Недообучение
- 3.6 Переобучение
- 3.7 Регуляризация
- 3.8 Отбор признаков
- 3.9 Параметры алгоритма
- 3.10 Подбора метапараметров
- 3.11 Основные типы алгоритмов
- 3.12 Многоклассовая классификация
- 3.13 Дисбаланс классов
- 3.14 Ансамбли алгоритмов
- 3.15 Метрики классификации
- 3.16 ROC-AUC метрика
- 3.17 Метрики регрессии
- 3.18 Метрики кластеризации
- 3.19 Разложение ошибки алгоритма
- 3.20 Кривые валидации
- 3.21 Кривые обучения
- 3.22 Метрические методы
- 3.23 Метод ближайших соседей
- 3.24 Линейные методы
- 3.25 Линейная регрессия
- 3.26 Логистическая регрессия

3.27 SVM

3.28 Ядра и спрямляющие пространства

3.29 Решающие деревья

3.30 Случайный лес

...отличие от беггинга над решающими деревьями...

3.31 Градиентный бустинг

3.32 Байесовские методы

Глава 4

Нейросети

В данной главе приводится обзор основных понятий и методов, связанных с нейросетями.