



Прикладная статистика и анализ данных

Съезд III



Пропуски в данных



Неслучайные пропуски

Возможные причины:

- ▶ Отказ респондента дать ответ на вопрос:

Сколько вы тратите в месяц на развлечения? <не скажу>

<не скажу> — важная информация!!!

- ▶ Признак не применим:

Посещаете ли вы рестораны? Нет

Какой ваш средний чек? <NaN>

В таких случаях:

1. Создается новый признак $x_{j'}$ = $I\{x_j \text{ пропущено}\}$;
2. Пропуски заменяются на любую константу, не встречающуюся в данных.

Случайные пропуски

Возможные причины:

- ▶ Сломалось оборудование:

Время	8:00	9:00	10:00	11:00	12:00
Температура возд.	21.4	22.1	NaN	24.2	25.5

- ▶ Было лень:



← ведро лени

Случайные пропуски

Варианты:

- ▶ Удалить строки, содержащие пропуски; **df.dropna**
- ▶ Заполнить пропуски:
 - ▶ по ближайшему соседу;
 - ▶ среднее/медиана; **df.fillna**
 - ▶ Интерполяция (временные ряды); **df.interpolate**
 - ▶ EM-алгоритм;
- ▶ Считать $X^T X$ и $X^T Y$ только по полным парам

$$\frac{1}{n} (X^T X)_{jk} = \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \approx \frac{1}{n_{jk}} \sum_{i=1}^n x_{ij} x_{ik} I\{x_{ij} \text{ и } x_{ik} \text{ не пропущены}\},$$

$$\frac{1}{n} (X^T Y)_{jk} = \frac{1}{n} \sum_{i=1}^n x_{ij} y_i \approx \frac{1}{n_j} \sum_{i=1}^n x_{ij} y_i I\{x_{ij} \text{ не пропущено}\}.$$

где n_{jk} — число полных пар (x_{ij}, x_{ik}) ; n_j — число заполненных x_{ij} .



Робастная регрессия



Метод Тейла

$$y = a + bx$$

$$\hat{b} = \text{med} \left\{ \frac{Y_j - Y_i}{x_j - x_i}, \quad 1 \leq i < j \leq n \right\}$$

$$\hat{a} = \text{med} \left\{ Y_i - \hat{b}x_i, \quad i = 1, \dots, n \right\}$$

Робастная оптимизация

$$\sum_{i=1}^n R(Y_i - x_i^T \theta) \rightarrow \min_{\theta \in \mathbb{R}^d}$$

- ▶ $R(x) = x^2$ — L_2 -норма, получаем обычный МНК;
- ▶ $R(x) = |x|$ — L_1 -норма;
- ▶ $R(x) = \frac{x^2}{2} I\{|x| < c\} + c(|x| - c/2) I\{|x| > c\}$ — Хубер;



Методы на основе ближайшего соседа



Метод ближайших соседей (kNN)

Пусть \mathcal{X} — метрическое пространство.

$x_1, \dots, x_n \in \mathcal{X}$ — обучающая выборка.

Y_1, \dots, Y_n — соответствующая целевая переменная.

Предположение:

свойства объекта меняются не сильно в его окрестности.

Тогда давайте смотреть на свойства k ближайших соседей.

Примеры.

Пусть $x \in \mathcal{X}$ — исследуемый объект.

$x_{(1)}, \dots, x_{(k)}$ — k его соседей в порядке удаления от x .

1. Классификация.

Предсказание — наиболее часто встречаемый класс.

2. Регрессия.

Предсказание — усреднение отклика по соседям.



Взвешенный метод ближайших соседей

Пусть $x \in \mathcal{X}$ — исследуемый объект.

$x_{(1)}, \dots, x_{(k)}$ — k его соседей в порядке удаления от x .

Y_1, \dots, Y_k — соответствующий отклик.

w_1, \dots, w_k — вклад k -го соседа, определяемый пользователем.

Способы определения веса:

- ▶ $w_j = 1 - j/k$ — зависящий от номера соседа
- ▶ $w_j = \|x - x_{(j)}\|^{-1}$ — зависящий от расстояния до соседа

$$\hat{y}(x) = \arg \max_y \sum_{j=1}^k w_j I\{Y_j = y\} \text{ — Классификация}$$

$$\hat{y}(x) = \frac{\sum_{j=1}^k w_j Y_j}{\sum_{j=1}^k w_j} \text{ — Регрессия}$$



Свойства

1. k — гиперпараметр модели;
2. Не редко на практике показывает хорошие результаты.
3. **Дорогое применение:**
для каждого x результат вычисляется за $O(n \ln n)$.

Способы решения проблемы:

- ▶ Хранение данных в виде дерева: K-D Tree, Ball Tree
- ▶ Приближенный поиск: Locality-sensitive hashing —
вероятностный метод понижения размерности многомерных данных.
Подбирает хеш-функций так, чтобы похожие объекты
с высокой степенью вероятности попадали в одну корзину.

<https://github.com/facebookresearch/faiss> — реализации методов



Непараметрическая регрессия



Метод локального усреднения

$$\hat{y}(x) = \sum_{i=1}^n w_i(x) Y_i \Big/ \sum_{i=1}^n w_i(x),$$

где $w_i(x) \geq 0$ убывает при удалении x от X_i .

Варианты:

1. $w_i(x) = I\{|x - X_i| \leq c\}$ — усреднение по окрестности x ;
2. $w_i(x) = (c - |x - X_i|)^k I\{|x - X_i| \leq c\}$
— взвешенное усреднение по окрестности x ;
3. Усреднение по k ближайшим соседям;
4. Ядерная оценка (см. далее).



Ядерная оценка Надарая-Ватсона

$$w_i(x) = \frac{1}{h} q\left(\frac{x - X_i}{h}\right),$$

где q — ядро (симметричная плотность),

$h > 0$ — ширина ядра.

Вероятностная интерпретация

Пусть X_1, \dots, X_n случайны и $Y_i = f(X_i, \varepsilon_i)$ — отклик. Тогда

$\frac{1}{n} \sum_{i=1}^n w_i(x) = \frac{1}{nh} \sum_{i=1}^n q\left(\frac{x - X_i}{h}\right)$ — ядерная оценка плотности X ;

$\frac{1}{n} \sum_{i=1}^n w_i(x) Y_i = \frac{1}{nh} \sum_{i=1}^n Y_i \cdot q\left(\frac{x - X_i}{h}\right)$ — "ядерное мат. ожид." EY ;

$\hat{y}(x) = \frac{\sum_{i=1}^n w_i(x) Y_i}{\sum_{i=1}^n w_i(x)}$ — "ядерное УМО" $E(Y|X)$.



Ядерная оценка: теорема о сходимости

Пусть

1. $\int_{\mathbb{R}} |q(y)| dy < \infty$;
2. $yq(y) \rightarrow 0$ при $|y| \rightarrow \infty$;
3. $EY^2 < \infty$;
4. $h_n \rightarrow 0, nh_n \rightarrow \infty$ при $n \rightarrow \infty$.

Тогда $\hat{y}(x) \xrightarrow{P} y(x)$ в точках непрерывности функции $f(x)$, плотности $p_X(x)$ и условной дисперсии $\sigma^2(x) = D(Y|X = x)$, если при этом $p(x) > 0$.

Наилучшая скорость сходимости квадратичного риска достигается при $h \sim n^{-1/5}$



Ядерная оценка: выбор ширины ядра

Функционал вида leave one out:

$$F(h) = \sum_{i=1}^n (Y_i - \hat{y}_{-i}(x_i))^2,$$

где $\hat{y}_{-i}(x)$ — ядерная оценка, построенная по выборке, из которой было исключено i -е наблюдение.

Утверждение

$$F(h) = \sum_{i=1}^n (Y_i - \hat{y}(x_i))^2 \left/ \left(1 - \frac{q(0)}{\sum_{k=1}^n q\left(\frac{x_i - x_k}{h}\right)} \right) \right.$$

Выбор h : $F(h) \rightarrow \min_h$



Ядерная оценка: доверительная лента

Предположения:

$\mu(x)$ — ожидаемый отклик;

$Y_i = \mu(x_i) + \varepsilon_i$ — наблюдаемый отклик, $E\varepsilon_i = 0$, $D\varepsilon_i = \sigma^2$;

q — гауссовское ядро.

Доверительная лента уровня доверия α :

$$\left(\hat{y}(x) - z_{3h}\delta(x), \hat{y}(x) + z_{3h}\delta(x) \right)$$

$$\delta(x) = \hat{\sigma} \sqrt{\frac{\sum_{i=1}^n w^2(x_i)}{\sum_{i=1}^n w(x_i)}}, \quad p = \frac{1 + \alpha^{1/3h}}{2},$$

$$\hat{\sigma}^2 = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (Y_{i+1} - Y_i)^2$$

Пояснение: $D(Y_{i+1} - Y_i) = D(\varepsilon_{i+1} - \varepsilon_i) = D\varepsilon_{i+1} + D\varepsilon_i = 2\sigma^2$



Ширина окна

$h = \text{const}$ не ок в случае "где-то пусто, а где-то густо".

Тогда $h(x) = \|x - X_{(k)}\|$, где $X_{(k)}$ — k -й ближайший сосед для x .

Выбросы

Оценка Надарая-Ватсона крайне неустойчива к выборсам.

См. алгоритм LOWESS.



Локальная линейная регрессионная модель

Модель $f(x) = x^T \theta(x)$

Для каждого x применяется взвешенный МНК:

$$\sum_{i=1}^n w_i(x) (Y_i - X_i^T \theta(x))^2 \rightarrow \min_{\theta(x)}$$

где $w_i(x) = q_{h(x)}(x - X_i)$.

Взвешенный МНК:

$$\hat{\theta}(x) = (X^T W(x) X)^{-1} X^T W(x) Y$$

$$W(x) = \text{diag} (w_1^2(x), \dots, w_n^2(x))$$



Проверка линейности логита в логистической регрессии



Сглаженные диаграммы рассеяния

Пусть $(x_1, Y_1), \dots, (x_n, Y_n)$ — обучающая выборка, где $Y_i \in \{0, 1\}$.

Выберем признак j и построим ядерную регрессию $y \sim x_j$:

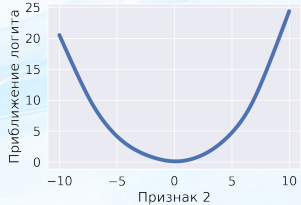
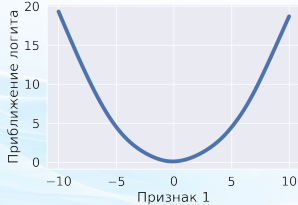
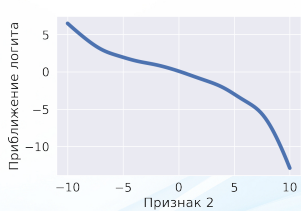
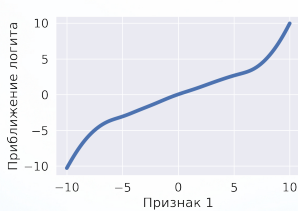
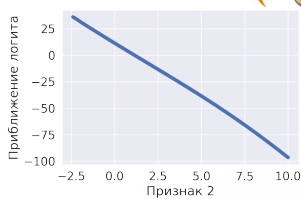
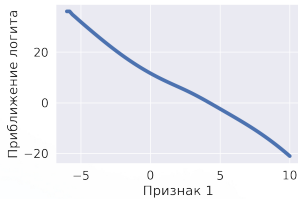
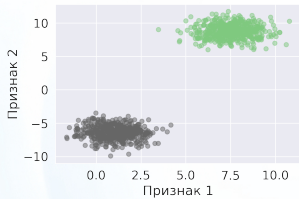
$$\hat{y}(x_j) = \sum_{i=1}^n q\left(\frac{x_j - x_{ij}}{h}\right) Y_i \bigg/ \sum_{i=1}^n q\left(\frac{x_j - x_{ij}}{h}\right),$$

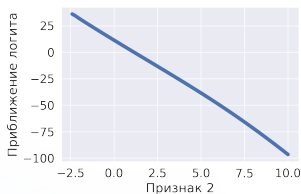
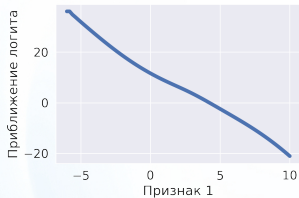
Эта регрессия — приближение вер-ти класса 1 в зависимости от x_j .

Отсюда делаем приближение логита:

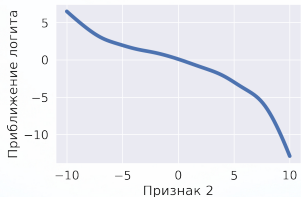
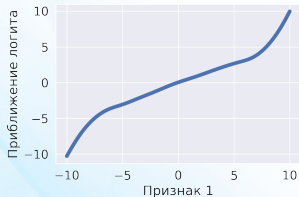
$$\text{logit}(x_j) = \log \frac{\hat{y}(x_j)}{1 - \hat{y}(x_j)}.$$

Проверка: график $\text{logit}(x_j)$ похож на прямую.

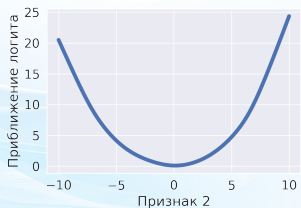
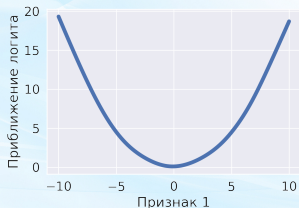




Логит линеен,
все хорошо



Классы линейно
разделимы, но зависи-
мость нелинейна



Классы не являются
линейно разделимыми



ВСЁ!