

Основные задачи и методы теории статистических выводов

Пусть задана выборка $X_1, \dots, X_n \sim F$

Необходимо сделать выводы о распределении F

а) Параметрические и непараметрические модели

1

Статистической моделью \mathfrak{F} называется множество распределений (плотностей, регрессионных зависимостей и т.п.). Параметрическая модель \mathfrak{F} может быть параметризована конечным числом параметров

$$\mathfrak{F} = \left\{ f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}, \quad \mu \in \mathbb{R}, \quad \sigma > 0 \right\}$$

Общий вид параметрической модели:

$$\mathfrak{F} = \left\{ f(x; \theta) : \theta \in \Theta \right\}$$

Компоненты θ , не представляющие интереса для исследователя, называют мешающими параметрами

2

Непараметрическая модель — множество \mathfrak{F} , которое нельзя параметризовать конечным числом параметров

Например, \mathfrak{F} - это множество всех кумулятивных функций распределения

Пример (оценка одномерного параметра). X_1, \dots, X_n - i.i.d., имеющие распределение Бернулли с параметром p

Пример (оценка двумерного параметра). X_1, \dots, X_n - i.i.d., имеющие нормальное распределение с параметрами μ и σ

Пример (непараметрическое оценивание). X_1, \dots, X_n - i.i.d. с распределением $F \in \mathfrak{F}_{\text{ALL}}$ (класс всех функций распределения)

Пример (непараметрическая оценка плотности распределения).

3 X_1, \dots, X_n - i.i.d. с распределением F и плотностью $f = F'$. Нельзя оценить f , предполагая только, что $F \in \mathfrak{F}_{\text{ALL}}$. Необходимы дополнительные предположения о гладкости f , например, что

$$f \in \mathfrak{F} = \mathfrak{F}_{\text{DENS}} \cap \mathfrak{F}_{\text{SOB}}$$

$\mathfrak{F}_{\text{DENS}}$ - множество всех плотностей распределения

$$\mathfrak{F}_{\text{SOB}} = \left\{ f : \int (f''(x))^2 dx < \infty \right\} \quad \text{- пространство Соболева}$$

Пример (непараметрическое оценивание функционалов). X_1, \dots, X_n - i.i.d. с распределением F . Необходимо оценить

$$\mu = \mathbb{E}(X_1) = \int x dF(x)$$

предполагая только, что μ существует. Будем использовать статистический функционал

4

$$\mu = T(F) = \int x dF(x),$$

заменяв F на его оценку

Пример (регрессия, предсказание и классификация). По наблюдениям $(X_1, Y_1), \dots, (X_n, Y_n)$ необходимо восстановить функцию регрессии $r(x) = \mathbb{E}(Y|X = x)$. Класс $r \in \mathfrak{F}$ может быть как параметрическим, так и непараметрическим. Часто используют следующую запись модели:

$$Y = r(X) + \epsilon, \mathbb{E}(\epsilon) = 0$$

$$\epsilon = Y - r(X), Y = Y + r(X) - r(X) = r(X) + \epsilon$$

$$\mathbb{E}(\epsilon) = \mathbb{E}\mathbb{E}(\epsilon|X) = \mathbb{E}(\mathbb{E}(Y - r(X))|X) = \mathbb{E}(\mathbb{E}(Y|X) - r(X)) =$$

$$\mathbb{E}(r(X) - r(X)) = 0$$

Существует два подхода к статистическим выводам: частотный и байесовский

5

Основные обозначения:

$$\mathfrak{F} = \{f(x; \theta) : \theta \in \Theta\}$$

$$\mathbb{P}_\theta(X \in A) = \int_A f(x; \theta) dx$$

$$\mathbb{E}_\theta(r(X)) = \int r(x) f(x; \theta) dx$$

$$\mathbb{V}_\theta$$

б) Основные задачи: точечное оценивание, доверительные множества, тестирование гипотез

Точечное оценивание – оценка неизвестной величины (параметры в параметрической модели, распределение, плотность распределения, функций регрессии, предсказание будущих значений некоторой случайной величины) в некотором наилучшем смысле

6

Будем обозначать оценка параметра θ через $\hat{\theta}$ или $\hat{\theta}_n$. $\hat{\theta}$ - зависит от данных, то есть это случайная величина.

X_1, \dots, X_n - i.i.d. наблюдения

$$\hat{\theta}_n = g(X_1, \dots, X_n)$$

$\text{bias}(\hat{\theta}_n) = \mathbb{E}_{\theta}(\hat{\theta}_n) - \theta$ - смещение оценки

Оценка $\hat{\theta}_n$ несмещенная, если $\mathbb{E}(\hat{\theta}_n) = \theta$

Оценка $\hat{\theta}_n$ состоятельная, если $\hat{\theta}_n \xrightarrow{P} \theta$

$se = se(\hat{\theta}_n) = \sqrt{\mathbb{V}(\hat{\theta}_n)}$ - стандартное отклонение оценки или
среднеквадратическая ошибка (зависит от неизвестного распределения)

\hat{se} - оценка стандартного отклонения

7

Пример (распределение Бернулли). X_1, \dots, X_n - i.i.d., имеющие
распределение Бернулли с параметром p

$$\hat{p}_n = n^{-1} \sum_i X_i$$

$$\mathbb{E}(\hat{p}_n) = n^{-1} \sum_i \mathbb{E}(X_i) = p$$

$$se = \sqrt{\mathbb{V}(\hat{p}_n)} = \sqrt{p(1-p)/n}$$

$$\hat{se} = \sqrt{\hat{p}(1-\hat{p})/n}$$

Среднеквадратическая ошибка оценки (MSE)

$$\text{MSE} = \mathbb{E}_{\theta}(\hat{\theta}_n - \theta)^2,$$

где $\mathbb{E}_{\theta}(\cdot)$ - математическое ожидание относительно плотности выборки

$$f(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

8 Утверждение (bias-variance decomposition):

$$\text{MSE} = \text{bias}^2(\hat{\theta}_n) + \mathbb{V}_{\theta}(\hat{\theta}_n)$$

Пусть $\bar{\theta}_n = E_{\theta}(\hat{\theta}_n)$, тогда так как $\mathbb{E}_{\theta}(\hat{\theta}_n - \bar{\theta}_n) = \bar{\theta}_n - \bar{\theta}_n = 0$, то

$$\begin{aligned}\mathbb{E}_{\theta}(\hat{\theta}_n - \theta)^2 &= \mathbb{E}_{\theta}(\hat{\theta}_n - \bar{\theta}_n + \bar{\theta}_n - \theta)^2 \\ &= \mathbb{E}_{\theta}(\hat{\theta}_n - \bar{\theta}_n)^2 + 2(\bar{\theta}_n - \theta)\mathbb{E}_{\theta}(\hat{\theta}_n - \bar{\theta}_n) + \mathbb{E}_{\theta}(\bar{\theta}_n - \theta)^2 \\ &= (\bar{\theta}_n - \theta)^2 + \mathbb{E}_{\theta}(\hat{\theta}_n - \bar{\theta}_n)^2 \\ &= \text{bias}^2(\hat{\theta}_n) + \mathbb{V}(\hat{\theta}_n)\end{aligned}$$

Утверждение: Если $\text{bias} \rightarrow 0$ и $\text{se} \rightarrow 0$, то $\hat{\theta}_n$ - состоятельная оценка

Оценка $\hat{\theta}_n$ асимптотически нормальна, если

$$\frac{\hat{\theta}_n - \theta}{\text{se}} \rightsquigarrow N(0, 1)$$

9

Доверительные множества

Доверительным интервалом с доверительной вероятностью $1 - \alpha$ для

параметра θ называется интервал $C_n = (a, b)$, где

$a = a(X_1, \dots, X_n)$ и $b = b(X_1, \dots, X_n)$ - такие функции

выборки, что $\mathbb{P}_\theta(\theta \in C_n) \geq 1 - \alpha$ для всех $\theta \in \Theta$

Замечание. $C_n = (a, b)$ - случайная величина, тогда как θ - неизвестная детерминированная величина

Замечание. Если θ - векторный параметр, то C_n называется доверительным множеством

В заданном определении доверительного интервала предполагается, что

$$10 \quad \mathbb{P}_\theta(\theta \in C_n) \geq 1 - \alpha \quad \text{для всех } \theta \in \Theta$$

Поточечным асимптотическим доверительным интервалом будем называть такой доверительный интервал, что

$$\liminf_{n \rightarrow \infty} \mathbb{P}_\theta(\theta \in C_n) \geq 1 - \alpha \quad \text{для всех } \theta \in \Theta$$

Равномерным асимптотическим доверительным интервалом будем называть такой доверительный интервал, что

$$\liminf_{n \rightarrow \infty} \inf_{\theta \in \Theta} \mathbb{P}_\theta(\theta \in C_n) \geq 1 - \alpha$$

Доверительный интервал не является вероятностным утверждением о θ , поскольку θ - не случайная величина, а детерминированный неизвестный параметр. Интерпретация может быть следующей – либо повтор одного и того же эксперимента, либо построение большого количество доверительных интервалов – в первом случае $1 - \alpha$ процентов времени неизвестный параметр попадет в интервал, во втором случае доля $1 - \alpha$ построенных доверительных интервалов накроет неизвестный параметр.

Пример. X_1, \dots, X_n - i.i.d., имеющие распределение Бернулли с параметром p . $\hat{p}_n = n^{-1} \sum_i X_i$

$$C_n = (\hat{p}_n - \epsilon_n, \hat{p}_n + \epsilon_n)$$

доверительный интервал с $\epsilon_n^2 = \log(2/\alpha)/(2n)$.

12 **Неравенство Хеффдинга:** Y_1, \dots, Y_n - i.i.d., при этом $\mathbb{E}(Y_i) = 0$ и с вероятностью единица $a_i \leq Y_i \leq b_i$. Тогда для любого $t > 0$ и $\epsilon > 0$

$$\mathbb{P} \left(\sum_{i=1}^n Y_i \geq \epsilon \right) \leq e^{-t\epsilon} \prod_{i=1}^n e^{t^2 (b_i - a_i)^2 / 8}$$

$$\mathbb{P}(p \in C_n) \geq 1 - \alpha$$

Отсюда получаем, что

Утверждение. Допустим, что $\hat{\theta}_n \approx N(\theta, \hat{s}e^2)$ (оценка асимптотически нормальная). Пусть Φ - стандартная нормальная функция распределения, $z_{\alpha/2} = \Phi^{-1}(1 - (\alpha/2))$, то есть $\mathbb{P}(Z > z_{\alpha/2}) = \alpha/2$ и

$$\mathbb{P}(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha, \text{ где } Z \sim N(0, 1),$$

$$C_n = (\hat{\theta}_n - z_{\alpha/2} \hat{s}e, \hat{\theta}_n + z_{\alpha/2} \hat{s}e).$$

Тогда

$$\mathbb{P}_\theta(\theta \in C_n) \rightarrow 1 - \alpha.$$

Действительно, положим $Z_n = (\hat{\theta}_n - \theta) / \hat{s}e$. Тогда согласно

предположению $Z_n \rightsquigarrow Z$, где $Z \sim N(0, 1)$.

$$\begin{aligned}
\mathbb{P}_\theta(\theta \in C_n) &= \mathbb{P}_\theta \left(\hat{\theta}_n - z_{\alpha/2} \hat{\text{se}} < \theta < \hat{\theta}_n + z_{\alpha/2} \hat{\text{se}} \right) \\
&= \mathbb{P}_\theta \left(-z_{\alpha/2} < \frac{\hat{\theta}_n - \theta}{\hat{\text{se}}} < z_{\alpha/2} \right) \\
&\rightarrow \mathbb{P} \left(-z_{\alpha/2} < Z < z_{\alpha/2} \right) \\
&= 1 - \alpha
\end{aligned}$$

14

Доверительный интервал такого вида является поточечным асимптотическим доверительным интервалом

Пример: Пусть $X_1, \dots, X_n \sim \text{Bernoulli}(p)$, $\hat{p}_n = n^{-1} \sum_{i=1}^n X_i$
 $\mathbb{V}(\hat{p}_n) = n^{-2} \sum_{i=1}^n \mathbb{V}(X_i) = n^{-2} \sum_{i=1}^n p(1-p) = n^{-2} np(1-p)$
 $= p(1-p)/n$

$$\text{se} = \sqrt{p(1-p)/n}$$
$$\hat{\text{se}} = \sqrt{\hat{p}_n(1-\hat{p}_n)/n}$$

15

Согласно центральной предельной теореме $\hat{p}_n \approx N(p, \hat{\text{se}}^2)$

Тогда приблизительный доверительный интервал с доверительной вероятностью $1 - \alpha$ имеет вид

$$\hat{p}_n \pm z_{\alpha/2} \hat{\text{se}} = \hat{p}_n \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}}$$

Тестирование гипотез

Делается предположение о процессе, генерирующем данные и задача состоит в том, чтобы определить, содержат ли данные достаточно информации, чтобы отвергнуть это предположение. Если информации не достаточно, то говорится, что опытные данные предположению (гипотезе) не противоречат.

16

Пример. Пусть $X_1, \dots, X_n \sim \text{Bernoulli}(p)$

$$H_0 : p = 1/2 \quad (\text{основная гипотеза})$$

$$H_1 : p \neq 1/2 \quad (\text{альтернативная гипотеза})$$

Рассмотрим статистику $T = |\hat{p}_n - (1/2)|$. Если она достаточно большая (пороговое значение будет определено позже), то основная гипотеза отклоняется