

МСПС | Введение

(1)

Схема проверки курса

- 1) Решение д/з унит. в итоговой оценке. Сравн. д/з - до ^{смер.} семинара?
- 2) ~~Midtest~~ + экзамен в конце. Экзамен ~~вот~~ будет ~~состоит из 2-х частей: теорет. и прак.~~ письменный теоретик.
- 3) На Стэнке у Лусика есть много всего. Есть программа курса ~~x~~. Чтобы получить информацию, написать ему e-mail. Д/з тоже есть на Стэнке. Д/з и программа на странице курса.
- 4) Список хороших книг тоже есть на Стэнке.

Анализ, основ. задачи статистики: по данным/наблюдениям сделать выводы о распределении, которому подчиняются эти данные.

Будем изучать след. 3 направления: параметрич., непараметрич. и байесовская.

1) Параметрическая: множество рассматриваемых распределений описаны одним или неск. параметрами:

$$\mathcal{P} = \{P_\theta, \theta \in \Theta\}, \text{ где } \Theta \subset \mathbb{R}^k, k\text{-разно утн, кол-во параметров.}$$

Пример: $\mathcal{P} = \{N(a, \sigma^2)\}$, здесь $\Theta = (a, \sigma^2)$; $\Theta = \mathbb{R} \times [0, +\infty)$.

Цели хороша: достаточно оценить параметр θ , чтобы решить основную задачу статистики - найти неизвест. распределение. Она так хороша, что ~~оценку можно находить с максимальной точностью~~ часто \exists много методов нахождения оценок параметра θ , из которых можно выбрать наилучшую.

Цели плоха: По виду гистограммы/других первичных методов анализа мы не всегда можем указать нужную модель. Многие методы пар. стат. неустойчивы к выбросам.

2) Непараметрическая: мы не можем параметризовать рассматриваемое семейство распределений, т.е.

$$\mathcal{P} = \{P_\alpha, \alpha \in A\}, \text{ где } A\text{ - произвольное мн-во индексов.}$$

Пример $\mathcal{P}_1 = \{\text{все распредел., у которых нулевое матожидание}\},$

$\mathcal{P}_2 = \{\text{все симметричные распределения}\}$

важный класс

Цели хороша: позволяет рассматривать широкое семейство распределений. Кроме того, методы непараметрич. статистики, как правило, без каких-либо начальных предпос. о виде распределения являются робастными, т.е. устойчивыми к малому отклонению наблюдений, т.е. выбросам, отклонениям от норм. стат. Известно, что выбросы составляют от 0,1% до 1% точек выборки чаще всего, а есть методы типа t-критерия Стюарта, которые очень чувствительны к ошибкам и престают работать.

Цели плоха: Арсенал методов непарам. статистики гораздо беднее, чем у параметрической.

3) Байесовская.

Наиболее развивающаяся за послед. 20 лет. Здесь $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$, но параметр θ полагается случайным. Т.е. θ имеет некое априорное распределение Q , которое мы выбираем, исходя из имеющихся на текущий момент данных, о неизвестном распределении.

Пример $\mathcal{P} = \{N(a, 1), a \in \mathbb{R}\}$, где a имеет нормальное распр. $N(a, 1)$.

Цели хороша Суждение традиционной "частотной" статистики приблизительны - даётся, как правило, ответ интервал для параметра. Кроме того, мы не умеем ответить на вопрос типа $P(0,3 < \theta < 0,9) = ?$ в завис. от выборки ответ может сильно отличаться? Байес. статистика даёт точные ответы.

Цели плоха А откуда мы это априорное распр. берём? И всегда есть такой априорный распр. Кажется, что с некоторой долей основан предположение. Даётся все. Байес. статистика

Семейства распределений

1) $z \sim N(a, \sigma^2)$. $Ez = a$, $Dz = \sigma^2$; $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-a)^2}{2\sigma^2}}$.
Если $z \perp \eta$, $\eta \sim N(\mu, \sigma^2)$, то $z + \eta \sim N(a + \mu, \sigma^2 + \sigma^2)$.

Верно правило 3 σ : $P(a - 3\sigma < z < a + 3\sigma) \geq 0,99$

2) $z \sim \text{Exp}(\alpha)$. Выбавет со скоростью.

$p(x) = \alpha e^{-\alpha x} I(x > 0)$. $Ez = \frac{1}{\alpha}$; $Dz = \frac{1}{\alpha^2}$.

Если $z_1, \dots, z_n \sim \text{Exp}(\alpha)$ и независ., то $\sum_{i=1}^n z_i \sim \Gamma(\alpha, n)$.

Про множеств. корр. скажем.

3) $z \sim U[a, b]$ (или $R[a, b]$): $p(x) = \frac{1}{b-a} \cdot I(x \in [a, b])$.

$$Ez = \frac{b+a}{2}; D_z = \frac{(b-a)^2}{12}.$$

4) $z \sim \text{Cauchy}(\theta)$. Бывает со сдвигом. Ez не \exists .

• Если $z_1 \sim \text{Cauchy}(\theta_1)$; $z_2 \sim \text{Cauchy}(\theta_2)$, то $z_1 + z_2 \sim \text{Cauchy}(\theta_1 + \theta_2)$ и независ.

5) $z \sim \Gamma(\alpha, \delta)$: $p(x) = \frac{x^{\alpha-1} \delta^\alpha}{\Gamma(\alpha)} e^{-\delta x} \cdot I(x > 0)$. $Ez = \frac{1}{\delta}$; $Dz = \frac{1}{\delta^2}$.

$z_1 \perp z_2$, $z_1 \sim \Gamma(\alpha, \delta_1)$; $z_2 \sim \Gamma(\alpha, \delta_2)$, то $z_1 + z_2 \sim \Gamma(\alpha, \delta_1 + \delta_2)$.

6) $z \sim B(\alpha, \beta)$, $p(x) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} \cdot I(x \in [0, 1])$, бывает сдвигом и расщеп.

• $z \sim \text{Pareto}(\alpha)$: $p(x) = \frac{\alpha}{x^{\alpha+1}} I(x \geq 1)$.

• $z \sim \text{Weibull}(\alpha, \beta)$ $F(x) = 1 - e^{-(x)^\beta} \cdot I(x > 0)$.

Дискретное

• $z \sim \text{Pois}(\lambda)$ $Ez = \lambda$; $Dz = \lambda$. $z \perp \eta$, $z \sim \text{Pois}(\lambda_1)$; $\eta \sim \text{Pois}(\lambda_2) \Rightarrow z + \eta \sim \text{Pois}(\lambda_1 + \lambda_2)$
 $P(z=k) = \frac{\lambda^k}{k!} e^{-\lambda}$

• $z \sim \text{Geom}(p)$, $p \in (0, 1)$: $P(z=k) = p(1-p)^{k-1}$. $Ez = \frac{1}{p}$; $Dz = \frac{1-p}{p^2}$.

• $z \sim \text{Bin}(n, \theta)$, $\theta \in (0, 1)$: $P(z=k) = C_n^k \theta^k (1-\theta)^{n-k}$
 $Ez = n\theta$; $Dz = n\theta(1-\theta)$

• $z \sim \text{Bern}(\theta)$ - бинаом. с пар. $\text{Bin}(1, \theta)$.

Методы построения оценок

Гипотеза - вид оценок

Опр. Статистика - (целеричная) ф-ция от выборки.

Опр. Оценки - статистика, которая имеет значение в мн-ве Θ .
 (т.е. оценивает θ или ф-ция от него)

1) Метод моментов. $\Theta = (\theta_1, \dots, \theta_k)$

а) Пусть у нас k параметров \Rightarrow вводим k проб. ф-ций (крательно-голомно-пробных)
 $g_1(x), \dots, g_k(x)$. Пусть $m_i(\theta_1, \dots, \theta_k) = E_{\theta} g_i(X)$

б) Составим систему:

$$\begin{cases} m_1(\theta_1, \dots, \theta_k) = E_{\theta} g_1(X) = \overline{g_1(X)} \\ m_k(\theta_1, \dots, \theta_k) = E_{\theta} g_k(X) = \overline{g_k(X)} \end{cases}, \text{ где } \overline{g_i(X)} = \frac{1}{n} \sum_{j=1}^n g_i(x_j).$$

в) Находим решение этой системы $\hat{\theta}_1, \dots, \hat{\theta}_k$.

Пример $N(a, \sigma^2)$. Если $X_1, \dots, X_n \sim N(a, \sigma^2)$ - выборка

Всего функций $g_1(x) = x$; $g_2(x) = x^2$. Тогда $Eg_1(X_1) = a$; $Eg_2(X_1) = \sigma^2 + a^2 \Rightarrow$

$\Rightarrow a = \bar{x}$
 $\sigma^2 + a^2 = \overline{x^2} \Rightarrow$ решая эту сист. относительно a и σ^2 , получаем
 $\hat{a} = \bar{x}$ и $\hat{\sigma}^2 = \overline{x^2} - (\bar{x})^2$.

Оценка мет. моментов устойчива, если m^{-1} (где $m = (m_1, \dots, m_k)$) непрерывна.

(2) Метод макс. правдоподобия

Опр. Ф-ция правдоподобия выборки X_1, \dots, X_n -

$f(X_1, \dots, X_n, \theta) = \prod_{i=1}^n p_\theta(X_i)$, где $p_\theta(x)$ - обобщ. плотность X_1 .

Опр. Оценка макс. правдоподобия $\hat{\theta} \equiv \arg \max_{\theta} f(X, \theta)$.

Пример Касгаров ОМН где выборки из $\text{Exp}(\alpha)$.

Ф-ция пл-е: $f(X, \alpha) = \prod_{i=1}^n \alpha e^{-\alpha X_i} = \alpha^n \cdot e^{-\alpha \sum_{i=1}^n X_i}$

Как ищется max? Возьмём лог, а потом продифференцируем. и
 уравни к 0, тогда
 найдем экстремум.

$L(X, \alpha) = \ln f(X, \alpha) = n \ln \alpha - \alpha \sum_{i=1}^n X_i$

$\frac{\partial}{\partial \alpha} L(X, \alpha) = \frac{n}{\alpha} - \sum_{i=1}^n X_i = 0 \Rightarrow \hat{\alpha} = \frac{n}{\sum X_i}$

Теорема Если функция пл-е приравнена к 0, то из ОМН
 получается асимпт. нормальна и ас. эффективна при нек. условиях
 на распределение.

(3) Другие методы.

Опр. z_α - квантиль ф.р. $F(x)$, если $z_\alpha = \inf \{x : F(x) \geq \alpha\}$.

Опр. Вариат. ряд выборки $X_{(1)} \leq \dots \leq X_{(n)}$ - упор. выборка
 поставим по порядку (на всех ω).

Опр. Выбороч. квантиль пл-е α : $\hat{z}_\alpha = X_{([n\alpha])}$.

Теорема (о выбороч. квантилях) Пусть меткость дифференцируемая
 тогда $\sqrt{n}(\hat{z}_\alpha - z_\alpha) \xrightarrow{d} N(0, \frac{\alpha(1-\alpha)}{f^2(z_\alpha)})$.

Можно ещё рассматривать ф-лы от $F(x)$ (типа E_i, D_i)
 и порекомендовать тогда $F_n(x)$ - эмпирич. ф-ция распредел.
 (т.н. метод порекомендовки)

Пример $p(x) = \frac{1}{1+(x-\theta)^2} \cdot \frac{1}{\pi}$. Тогда из пл-е θ сущес. выбороч.
 медиана, т.е. $X_{(\frac{n}{2})}$.

Виды оценок

Будущее исправление смещенности

Опр. $\hat{\theta}$ - несмещ. для $\tau(\theta)$, если $\forall \theta \in \Theta \ E_0 \hat{\theta} = \tau(\theta)$.

Пример $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ всегда несмещ. оценивает $E_0 X_i$ (если оно есть)

Опр. Если $\forall \theta \in \Theta \ \hat{\theta}_n \xrightarrow{P_0} \theta$, то $\hat{\theta}_n$ кон. состоят. оц. θ . Борьба за экстремальность!

Пример Выбороч. медиана $\hat{\mu} = X_{(\frac{n+1}{2})} \xrightarrow{P_0} z_{1/2}$, где $F(z_{1/2}) = \frac{1}{2}$
где непрерыв. распредел. (или, точнее, $z_{1/2} = \inf\{x: F(x) \geq \frac{1}{2}\}$)

Опр. Если $\forall \theta \in \Theta \ \sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \sigma^2(\theta))$, то $\hat{\theta}_n$ кон. асимпт. нормальный. $\sigma^2(\theta)$ кон. асимпт. дисперсия, или эффективность $\hat{\theta}_n$, $q(\theta)$ (как правило, $= 0$) - асимпт. смещение. Пример: \bar{X} , у ГПТ смещен.

Доверит. интервалы

Опр. Пара статистик $(T_1(X), T_2(X))$ кон. доверит. интервалом для пар-ра θ , если $\forall \theta \in \Theta \ P_0(T_1(X) < \theta < T_2(X)) \geq 1 - \alpha$.
уровень доверия $1 - \alpha$

Опр. — // — асимпт. доверит. интервал — // — ,
если $\forall \theta \in \Theta \ P_0(T_{1n}(X) < \theta < T_{2n}(X)) \rightarrow 1 - \alpha, n \rightarrow \infty$

Пример Построим ас. доверит. интервал для пар-ра α в модели $\text{Exp}(\alpha)$.

Строится он по общей схеме - применение ГПТ, если есть вторая момент, и теорема о выбороч. квантилях.

Здесь $E X = \frac{1}{\alpha}$, $D X = \frac{1}{\alpha^2}$.

Тогда по ГПТ $\frac{\sum X_i - n \cdot \frac{1}{\alpha}}{\sqrt{n \cdot \frac{1}{\alpha^2}}} \xrightarrow{d} N(0, 1)$, т.е. асимпт. доверит. интервал такой:

$$\left(z_{\gamma} < \frac{\sum X_i - n \cdot \frac{1}{\alpha}}{\sqrt{n \cdot \frac{1}{\alpha^2}}} < z_{1-\gamma} \right) = \left(z_{\gamma} < \sqrt{n} \cdot \alpha \cdot \bar{X} - \sqrt{n} < z_{1-\gamma} \right) = \left(\frac{z_{\gamma}}{\sqrt{n}} + 1 < \alpha \bar{X} < \frac{z_{1-\gamma}}{\sqrt{n}} + 1 \right) \\ = \left(\frac{z_{\gamma} + \sqrt{n}}{\sqrt{n} \bar{X}} < \alpha < \frac{z_{1-\gamma} + \sqrt{n}}{\sqrt{n} \bar{X}} \right), \text{ где } z_{\gamma} \text{ и } z_{1-\gamma} \text{ квантили уровня } \gamma \text{ и } 1-\gamma \text{ у } N(0, 1)$$

Часто мы дисперсию заменяем на её состоятельную оценку, так можно делать по 1. Слуцкого.

МСПС | Байесовская статистика

(1)

Отличие байесов. статистики от традиционной: в традицион. мы оцениваем параметр, ^{в байесовской} ~~узнаем~~ находим его распределение.

Т.е. результатом действий традиц. статистики ^{не всегда} ~~явл.~~ число, байесовская - вероятностная мера.

Плюсы

- 1) Заранее известно, что надо делать (есть алгоритм)
- 2) Для сложных (сильно параметрич.) моделей гораздо лучше работает
- 3) Интуитив. ясна: позволяет отвечать на вопросы типа $P(\theta > 0 | y)$, где y - наши данные.

При проверке гипотез $P(H_0 | y) = 1 - P(H_1 | y)$, что неверно для традиц. статистики.

Минусы

- 1) Использование априор. распр. для пар-ра θ . Конечно, с какого толчка оно берется.
- 2) Разные люди могут прийти к разн. выводам на одних и тех же данных.
- 3) Большие объ. сложн.

В общем, можно исп. и те, и другие методы для получения ясной картины о данных.

Итак, байесовский анализ состоит из 3 шагов:

- 1) Установл. совместного распр. $p(y, \theta) = p(y | \theta) \cdot p(\theta)$

распр. выбора θ т.е. априорное распр., которое мы выбираем

- 2) Поиск услов. распр. (наш результат!)

$$p(\theta | y) = \frac{p(y, \theta)}{p(y)} = \frac{p(y | \theta) \cdot p(\theta)}{\int p(y | \theta) \cdot p(\theta) d\theta}$$

вообще говоря, нормиров. константа:

Байес. оценка

$$\hat{\theta} = E \int \theta p(\theta | y) d\theta$$

«Байесовская теорема» $p(\theta | y)$ содержит всю инф. из выбора + априор. распр.

- 3) Проверка модели (тестирование гипотез байес. статистики обходит этот шаг)

Есть 3 величины, кот. нам очень интересны:

- 1) $p(y) = \int p(y, \theta) d\theta$ «prior predictive» (априор. ^{прогноз} предсказание?)

- 2) Пусть $\theta = (\theta_1, \dots, \theta_m)$.

Нам будет интересен $p(\theta_i | y) = \int p(\theta_i | \hat{\theta}_i, y) \cdot p(\hat{\theta}_i | y) d\hat{\theta}_i (= \int p(\theta_i, \hat{\theta}_i | y) d\hat{\theta}_i)$

где $\hat{\theta}_i = \theta \setminus \theta_i$ - вектор $(\theta_1, \theta_2, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_m)$.

- 3) Апостериорная прогноз

Пусть \tilde{y} - будущие данные, сверх с нашими экспериментальными (т.е. будущие наблюдения)

$$\text{Тогда } p(\tilde{y} | y) = \int p(\tilde{y} | \theta, y) \cdot p(\theta | y) d\theta = \int p(\tilde{y} | \theta) \cdot p(\theta | y) d\theta,$$

поскольку при выбранном θ \tilde{y} и y независимы (также понятно, что $p(\tilde{y} | y)$ зависима, потому что y и \tilde{y} одной модель описаны).