

МСПС | Расстояния в статистике |

(1)

Задача Есть выборка Y_1, \dots, Y_n из неизвестной плот. $p_0(x)$. Мы предполагаем, что $p_0(x)$ из семейства плотностей $p(x, \mu)$ (но это не обязательно так), $\mu \in \mathbb{R}$, ну и хотим оценить μ , т.е. найти такую плотность $p(x, \mu)$, которая вроде как должна равняться $p_0(x)$.

Цель: $\hat{\mu}(p_0) = \arg \min_{\mu} d[p_0, p(\cdot; \mu)]$, где d - некое расстояние (или мера близости) между плотностями.

Но эта штука не явл. оценкой, потому что $p_0(x)$ мы не знаем!

Знаем только выборку. Но тогда можно брать вкап. оценки $p_0(x)$

эмпирир. плотность $\delta(x - Y)$ ($x = (x_1, \dots, x_n)$, $Y = (Y_1, \dots, Y_n)$) (сравни с ядерными оценками плотности - там ведь тоже фактически δ -функцией пользуемся!)

Сначала разберемся с расстояниями:

1) Расстояние полной вариации - самое, пожалуй, естественное с вероятн. точки зрения. Пусть ζ, η - сл. в. с плотностями p_ζ и p_η .

$$\text{Тогда } D(p_\zeta, p_\eta) = \sup_{A \in \mathcal{F}} \left| \int_A p_\zeta(x) dx - \int_A p_\eta(x) dx \right|, \quad \mathcal{F} - \sigma\text{-алгебра.}$$

Вроде как от. неудобно вычислять, ведь супремум по всей σ -алгебре берется. Но есть такая теорема:

Теорема (Шерфе) Если существуют плотности (обобщенные) $p_\zeta(x), p_\eta(x)$, $x \in \mathbb{R}^n$ сл. векторов ζ, η , то $D(p_\zeta, p_\eta) = \frac{1}{2} \int_{\mathbb{R}^n} |p_\zeta(x) - p_\eta(x)| dx$

Проблема в том, что непонятно, как вычислить $\hat{\mu}(p_0) = \arg \min_{\mu} D[p_0, p(\cdot; \mu)]$

Кроме того, $D(\delta(\cdot - Y^n), p(\cdot; \mu)) = \frac{1}{2} \int_{\mathbb{R}^n} |\delta(x - Y^n) - p(x, \mu)| dx = 1$ (на 2 ант. надо разбить)

$$\approx \frac{1}{2} \int_{Y^n} |\delta(x - Y^n) - p(x, \mu)| dx + \frac{1}{2} \int_{\mathbb{R}^n / Y^n} |p(x, \mu)| dx.$$

2) Расстояние Кульбака-Лейблера

$$K(p_\zeta, p_\eta) = \int_{\mathbb{R}^n} p_\zeta(x) \log \frac{p_\zeta(x)}{p_\eta(x)} dx.$$

Прежде всего, $K(p_\zeta, p_\eta)$ расст. не является, т.к. $K(p_\zeta, p_\eta) \neq K(p_\eta, p_\zeta)$.

Но отсутствие симметрии - это реальная сит. в статистике.

Ведь, с одной стороны, есть реал. данные с неув. плотностью $p_0(x)$

(кладём $p_\zeta(x) = p_0(x)$), который мы управлять не можем, а есть семейства плотностей $p(x, \mu) (= p_\eta(x))$, где мы можем варьировать μ . Но тогда как свести расст. Кульбака-Лейблера с расст. полной вар.?

Теорема (Нер-во Пинскера)

Пусть сл. вектора ζ, η имеют плотности (обобщ.) $p_\zeta(x), p_\eta(x)$, $x \in \mathbb{R}^n$, тогда

$$K(p_\zeta, p_\eta) \geq 2 [D(p_\zeta, p_\eta)]^2$$

Крайне удобное св-во расст. К.-Л. таково:

Теорема Пусть $\zeta = (\zeta_1, \dots, \zeta_n)$ и $\eta = (\eta_1, \dots, \eta_n)$ - выборки. Тогда

$$K(p_\zeta, p_\eta) = n K(p_{\zeta_1}, p_{\eta_1})$$

Еще расстояние:

3) Расстояние Хеллингера: $H(p_1, p_2) := \left[\int |\sqrt{p_1(x)} - \sqrt{p_2(x)}|^2 dx \right]^{1/2}$

Из пер-ва Коши-Бун. следует, что $H(p_1, p_2) \geq D(p_1, p_2)$.

4) Расстояние хи-квадрат: $\chi^2(p_1, p_2) := \int_{\mathbb{R}^n} \frac{|p_1(x) - p_2(x)|^2}{p_1(x)} dx$.

Это первый неметрич. метр. в метр. метрике $K(p_1, p_2)$ (логарифмический метр).

Связь расст. Кульбака-Лейблера и ММП.

Пусть мы хотим использовать для проверки вер. модели расст. Кульбака-Лейблера. Если бы мы знали $p_0(z)$ - плотность Y_1, \dots, Y_n (всего n), то искали бы параметр μ семейства распредел. так:

$$\hat{\mu}(p_0) = \arg \min_{\mu} \int_{\mathbb{R}^n} p_0(z) \log \frac{p_0(z)}{p(z, \mu)} dz = \arg \max_{\mu} \int_{\mathbb{R}^n} p_0(z) \log p(z, \mu) dz.$$

Но, как и ранее, p_0 мы не знаем, и поэтому в качестве оценки p_0 берем $\delta(z - y)$, получаем:

$$\hat{\mu}(y) = \arg \max_{\mu} \int_{\mathbb{R}^n} \delta(z - y) \log p(z, \mu) dz = \arg \max_{\mu} p(y, \mu) - \text{так это } \mu$$

ищется максимум логарифмической φ -фн. правдоподобия! если это макс. если
 Но тогда $\log p(y, \mu) = \sum_{i=1}^n \log p(y_i, \mu) \Rightarrow \hat{\mu}(y) \rightarrow \arg \max_{\mu} E \log p(y, \mu)$
можно сказать, что так будет

$$\hat{\mu}(y) = \arg \min_{\mu} E \log \frac{p_0(y_i)}{p(y_i, \mu)} = \arg \min_{\mu} K(p_0, p(\cdot, \mu)).$$

Т.е. при больших n метод макс. правдоподобия вводит плотность из зад. параметриз. сем-ва, наход. на миним. расст. Кульбака-Лейблера от истинной, но неизвестной плотности наблюдений.

Основная книга - Вальд, "Последовательный анализ".

В процедуре проверки гипотез мы, как правило, действуем по след. алгоритму: по всем наблюдениям, которые у нас есть, строим статистику критерия, исходя из значений которой, принимаем или отвергаем основную гипотезу.

Пусть у нас наблюдения поступают в процессе проверки, тогда в последовательном анализе мы фактически делаем проверку для каждого k , начиная с 1, это позволяет достигнуть уровня значимости α и мощности β , используя в среднем меньшее количество наблюдений.

Итак, пусть $f_m(\theta) = \prod_{i=1}^m p(X_i, \theta)$ - ф-ция правдоподобия выборки (X_1, \dots, X_m)

Пусть проверим простую гипотезу $H_0: \theta = \theta_0$ против альтернативы $H_1: \theta = \theta_1$. Тогда на каждой стадии эксперимента (т.е. n) действуем по след. алгоритму:

- 1) Если $B < \frac{f_m(\theta_1)}{f_m(\theta_0)} < A$, то эксперим. продолжается и производится доп. наблюд.
- 2) Если $\frac{f_m(\theta_1)}{f_m(\theta_0)} \geq A$, то отклоним H_0 и законч. процесс проверки
- 3) Если $\frac{f_m(\theta_1)}{f_m(\theta_0)} \leq B$, то примем H_0 и законч. процесс проверки.

Эти вычисления проще проводить с помощью логарифма отнош. ф-ций:

$$\ln \frac{f_m(\theta_1)}{f_m(\theta_0)} = \sum_{i=1}^m \ln \frac{p(X_i, \theta_1)}{p(X_i, \theta_0)} =: \sum_{i=1}^m Z_i, \text{ тогда сравним } \sum_{i=1}^m Z_i \text{ с } \ln A \text{ и } \ln B.$$

Пример Пусть $X_1, \dots, X_n \sim \text{Bern}(p)$, $H_0: p = p_0$, $H_1: p = p_1$, тогда $Z_i = \begin{cases} \ln \frac{p_1}{p_0}, & \text{если } X_i = 1 \\ \ln \frac{1-p_1}{1-p_0}, & \text{если } X_i = 0 \end{cases}$, т.е. получается такое слуг. блуждание

Как выбирать константы A и B ?

Пусть процесс у нас закончился отклонением гипотезы H_0 , т.е. в какой-то момент $\frac{f_m(\theta_1)}{f_m(\theta_0)} \geq A$. Если гипотеза H_0 верна, то вероятность такого исхода (таких исходов) $\leq \alpha$, где α - ур. знач., если H_1 верна - то вероятность $1 - \beta$ (это нужно отдельно доказывать) - нужно док., что процесс окончится), т.е. $\frac{1-\beta}{\alpha} \geq A$ по идее должно быть.

А вероятность получения выборки, на кот. гипотеза H_0 отвергается, при невер. H_0 хотя бы в A раз больше, чем при вер. H_0 , поэтому $\frac{1-\beta}{\alpha} \geq A$. Аналогично, $B \geq \frac{\beta}{1-\alpha}$. И можно док. следующее утверждение:

Утв. Если ведь $B = \frac{\beta}{1-\alpha}$, $A = \frac{1-\beta}{\alpha}$, то вер. ошибки I рода неслучайно критерия будет $\leq \alpha$, а вер. ош. II рода $\leq \beta$.

* Это всё, конечно, прекрасно, но проверка сложных гипотез - гораздо более интересная и важная задача. Попробуем найти вероятность того, что последов. анализ закончится принятием гипотезы H_0 , когда истинное знач. параметра равно θ , $\theta \neq \theta_0, \neq \theta_1$, эта вероятность называется оперативной характеристикой, обозн. $L(\theta)$

МСПС (Последовательный анализ)

(2)

Если приобрести эффект от превышения функции при оконт. процедурах (об этом позже), то $L(\theta) = \frac{A^{h(\theta)} - 1}{A^{h(\theta)} - B^{h(\theta)}}$, где $h(\theta)$ находится (и $h(\theta) \neq 0$) как решение след. ур-ния: $E_{\theta} \left(\frac{p(x, \theta_1)}{p(x, \theta_0)} \right)^{h(\theta)} = 1$, т.е. (в нейр. случае) $\int \left(\frac{p(x, \theta_1)}{p(x, \theta_0)} \right)^{h(\theta)} p(x, \theta) dx = 1$.

Замечание. Что делать, если все наблюдения закончились, а выбора мы так и не сделали (т.е. не отвергли H_0 , не приняли её), то тогда решение следующее: если при n_{\max} $\ln B < \sum_{i=1}^{n_{\max}} z_i < \ln A$, то если $\sum_{i=1}^{n_{\max}} z_i > 0$, то отвергаем H_0 , если $\sum_{i=1}^{n_{\max}} z_i \leq 0$, то принимаем H_0 . Конечно, вер-от ошибок I и II рода уменьшается, но при большом n_{\max} эти величины становятся незначительными.

Замечание Так насколько же всё-таки последов. анализу Вальда лучше критерий Неймана-Пирсона? Ну, при фиксированных α и β последовательный анализ требует в среднем в 2 раза меньше наблюдений, чем стандартный критерий Неймана-Пирсона, подробнее см. книгу Вальда.

Как в итоге проверить сложную гипотезу? Дело в том, что специфич. х-ка, как правило, является монотон. ф-цией по θ .

Проверим гипотезу $H_0: \theta = \theta_0$ против $H_1: \theta > \theta_0$.

Можно ввести такое $\theta_1 > \theta_0$, что применение основной гипотезы в области $\theta > \theta_1$ приведёт к ошибке, имеющей практическое значение, т.е. область отклонения — $\theta \geq \theta_1$, и область безразличия $\theta_1 > \theta > \theta_0$.

Ну и строим обобщенный послед. критерий по θ_1 и θ_0 такой, что $L(\theta_0) = 1 - \alpha(\geq)$ и $L(\theta_1) \leq \beta$.

В случае двусторонней альтернативы $H_0: \theta = \theta_0$ и $H_1: \theta \neq \theta_0$ можно

ввести $f_m(\theta) = \frac{1}{2}(f_m(\theta_0 - \delta) + f_m(\theta_0 + \delta))$ вместо $f_m(\theta_1)$ и рассл.

послед. критерий $B \leq \frac{\hat{f}_m(\theta)}{f_m(\theta_0)} \leq A$, причём δ выбирается примерно так же, как и в одностороннем случае.