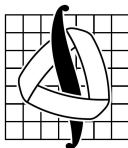


Центральные предельные теоремы для перестановочных случайных величин

Александр Ракитько, МГУ



Международная научная конференция
"Ломоносов-2014"
Москва, 9 сентября

1. Основные определения и понятия.

Определение

Набор случайных величин $\{X_i\}_{i=1}^n$ называется **перестановочным**, если

$$\text{Law}(X_1, \dots, X_n) = \text{Law}(X_{\sigma(1)}, \dots, X_{\sigma(n)}),$$

где $\sigma \in \mathbb{S}(n)$.

Определение

Бесконечная последовательность X_1, X_2, \dots называется **перестановочной**, если для каждого $n \in \mathbb{N}$ набор $\{X_i\}_{i=1}^n$ перестановочен.

Примеры:

- 1 Независимые одинаково распределенные случайные величины.
- 2 Пусть $\{\varepsilon_i\}$ – н.о.р.с.в. Положим $X_i = Y + \varepsilon_i$ или $X_i = Y \cdot \varepsilon_i$.
- 3 Урновые схемы: выбор без возвращения, схема Пойя. Здесь X_i – цвет шара, извлечённого на i -ом шаге.
- 4 Пусть $n > 1$ и $\sigma : \Omega \rightarrow \mathbb{S}(n)$ распределена равномерно. Положим $X_i = \mathbb{I}\{\sigma(i) = i\}$.

Классическая теорема **де Финетти** утверждает:

Теорема

Каждой перестановочной последовательности $\{X_i\}_{i \geq 1}$ с $X_1 \in \{0, 1\}$ соответствует единственное распределение μ на $[0, 1]$ такое, что

$$\begin{aligned} \mathbb{P}(X_1 = 1, \dots, X_k = 1, X_{k+1} = 0, \dots, X_n = 0) = \\ = \int_0^1 p^k (1 - p)^{n-k} \mu(dp). \end{aligned}$$

В случае конечного набора перестановочных случайных величин теорема де Финетти неверна.

Контрпример:

$$\begin{aligned}\mathbb{P}(X_1 = 0, X_2 = 1) &= \mathbb{P}(X_1 = 1, X_2 = 0) = 1/2, \\ \mathbb{P}(X_1 = 0, X_2 = 0) &= \mathbb{P}(X_1 = 1, X_2 = 1) = 0.\end{aligned}$$

Набор $\{X_1, X_2\}$ перестановочен. Но

$$0 = \int_0^1 p^2 \mu(dp) = \int_0^1 (1-p)^2 \mu(dp).$$

Следовательно, подходящей меры μ не существует.

2. ЦПТ для перестановочного массива.

Будем рассматривать массив $\{X_{n,i}, i = 1, \dots, k_n\}_{n \geq 1}$ случайных величин.

Пусть выполняются следующие условия:

- Набор $\{X_{n,i}, i = 1, \dots, k_n\}$ перестановочен для каждого $n \in \mathbb{N}$.
- $k_n \in \mathbb{N}$ и $k_n \rightarrow \infty$ при $n \rightarrow \infty$.
- $\sup_{n \in \mathbb{N}} \mathbb{E} X_{n,1}^4 < \infty$.

Theorem

Пусть выполнены следующие условия

- 1 $\mathbb{E}X_{n,1}^2 - \mathbb{E}X_{n,1}X_{n,2} \rightarrow \sigma^2 > 0$ при $n \rightarrow \infty$,
- 2 $\text{cov}(X_{n,1}^2, X_{n,2}^2) + \text{cov}(X_{n,1}X_{n,2}, X_{n,3}X_{n,4}) - 2\text{cov}(X_{n,1}^2, X_{n,1}X_{n,2}) \rightarrow 0$ при $n \rightarrow \infty$.

Тогда для любой возрастающей последовательности $\{m_n\}_{n \geq 1}$ натуральных чисел такой, что $m_n/k_n \rightarrow \alpha < 1$ при $n \rightarrow \infty$, справедливо

$$m_n^{-1/2} \sum_{i=1}^{m_n} (X_{n,i} - \hat{\mu}_{k_n}) \rightarrow Z_{0, (1-\alpha)\sigma^2} \sim \mathcal{N}(0, (1-\alpha)\sigma^2),$$

где $\hat{\mu}_{k_n} := k_n^{-1} \sum_{i=1}^{k_n} X_i$.

Рассмотрим $\{Y_i\}_{i=1}^m$ – перестановочный набор, причем

$$\mathbb{E} Y_1 = 0, \quad \mathbb{E} |Y_1|^3 < \infty, \quad \sum_{i=1}^m Y_i = 0 \text{ п.н.} \quad (1)$$

Пусть $\Sigma = (\sigma_{i,j})_{1 \leq i,j \leq m}$ – ковариационная матрица вектора (Y_1, \dots, Y_m) .

Для функции $h : \mathbb{R}^d \rightarrow \mathbb{R}$ и $k \in \mathbb{N}$ положим

$$C_h^{(k)} := \max_{i_1, \dots, i_d \geq 0, \sum_{j=1}^d i_j = k} \left\| \frac{\partial^k h}{\partial x_1^{i_1} \dots \partial x_d^{i_d}} \right\|_{\infty}.$$

Сформулируем вспомогательный результат из [1] (Rollin, 2013):

Лемма

Пусть Y – вектор, состоящий из перестановочных случайных величин, с ковариационной матрицей Σ . Предположим, что условия (1) выполнены. Тогда

$$\begin{aligned} |\mathbb{E}h(Y) - \mathbb{E}h(Z)| &\leq \\ &\leq C_h^{(2)} \left[\text{var} \left(\sum_{i=1}^m Y_i^2 \right) \right]^{\frac{1}{2}} + 16m C_h^{(3)} \mathbb{E}|Y_1|^3, \end{aligned}$$

где $Z \sim \mathcal{N}(0, \Sigma)$.

Схема доказательства:

- Переходим к перестановочному набору $\{Y_i := m_n^{-1/2}(X_{n,i} - \hat{\mu}_{k_n}), i = 1, \dots, k_n\}$.
- Метод Линдеберга.
- С помощью леммы оцениваем

$$|\mathbb{E}h(Y) - \mathbb{E}h(Z)|,$$

где $h(x_1, \dots, x_{m_n}) = f(x_1 + \dots + x_{m_n})$,
а Z – гауссовский вектор с той же
ковариационной структурой, что и Y .

Замечание

Введем $\hat{\sigma}_{k_n}^2 := \frac{1}{k_n} \sum_{i=1}^{k_n} (X_{n,i} - \hat{\mu}_{k_n})^2$. Предположим, что

$$\sup_{n \in \mathbb{N}} \mathbb{E}((X_{n,1} - \hat{\mu}_{k_n})/\hat{\sigma}_{k_n})^4 < \infty.$$

Тогда, для последовательности $(m_n)_{n \in \mathbb{N}}$, определенной ранее, можно установить следующий вариант центральной предельной теоремы

$$\frac{1}{\sqrt{m_n}} \sum_{i=1}^{m_n} \left(\frac{X_{n,i} - \hat{\mu}_{k_n}}{\hat{\sigma}_{k_n}} \right) \xrightarrow{law} Z_{0,1-\alpha} \sim \mathcal{N}(0, 1 - \alpha).$$

Интересные результаты можно получать с помощью **мартингального подхода**, предложенного Вебером. Введём фильтрацию

$$\mathfrak{F}_{n,j} := \sigma\{X_{n,1}, \dots, X_{n,j}, \sum_{i=j+1}^{k_n} X_{n,i}\}.$$

Очевидно, можно переписать

$$\begin{aligned} \frac{1}{\sqrt{m_n}} \sum_{i=1}^{m_n} X_{n,i} &= \frac{1}{\sqrt{m_n}} \sum_{i=1}^{m_n} (X_{n,i} - \mathbb{E}(X_{n,i} | \mathfrak{F}_{n,(i-1)})) + \\ &+ \frac{1}{\sqrt{m_n}} \sum_{i=1}^{m_n} \mathbb{E}(X_{n,i} | \mathfrak{F}_{n,(i-1)}). \end{aligned}$$

Действительно, в силу измеримости и перестановочности

$$\begin{aligned}\sum_{i=j}^{k_n} X_{n,i} &= \mathbb{E}\left(\sum_{i=j}^{k_n} X_{n,i} \mid \mathfrak{F}_{n,(j-1)}\right) = \\ &= \sum_{i=j}^{k_n} \mathbb{E}(X_{n,i} \mid \mathfrak{F}_{n,(j-1)}) = (k_n - j + 1) \mathbb{E}(X_{n,j} \mid \mathfrak{F}_{n,(j-1)}).\end{aligned}$$

Таким образом, получаем

$$\mathbb{E}(X_{n,j} \mid \mathfrak{F}_{n,(j-1)}) = \frac{1}{k_n - j + 1} \sum_{i=j}^{k_n} X_{n,i}.$$

Пусть $X = (X_1, \dots, X_n)$ – случайный вектор с $X_k : \Omega \rightarrow \{0, 1, \dots, s\}$, $k = 1, \dots, n$. Положим

$$\mathbb{X} = \{0, \dots, s\}^n,$$

$$\mathbb{Y} = \{-m, \dots, 0, \dots, m\},$$

здесь $m \in \mathbb{N}$. Предположим, что $Y : \Omega \rightarrow \mathbb{Y}$, $f : \mathbb{X} \rightarrow \mathbb{Y}$ и задана штрафная функция $\psi : \mathbb{Y} \rightarrow \mathbb{R}_+$. Введем функционал ошибки

$$Err(f) := \mathbb{E}|Y - f(X)|\psi(Y).$$

- $Err(f) \rightarrow \min_f \implies f_{opt}$
- Распределение (X, Y) неизвестно:

$$\begin{aligned} f_{opt}(\cdot) &\implies \widehat{f}_{PA}(\cdot), \\ \psi(\cdot) &\implies \widehat{\psi}(\cdot), \\ Err(\cdot) &\implies \widehat{Err}(\cdot). \end{aligned}$$

- **K-кроссвалидация.** Для $K \in \mathbb{N}$, $(K > 1)$ и $k = 1, \dots, K$ введем разбиение

$$\begin{aligned} S_k(N) := \{ & (k-1)[N/K] + 1, \dots, \\ & k[N/K] \mathbb{I}\{k < K\} + N \mathbb{I}\{k = K\} \}. \end{aligned}$$

Удается преобразовать

$$Err(f) = \sum_{i=0}^{2m-1} \sum_{i-m < |y| \leq m} \psi(y) P(Y = y, |f(X) - y| > i).$$

Тогда в качестве оценки будем использовать

$$\widehat{Err}_K(f_{PA}, \xi_N) := \sum_{i=0}^{2m-1} \sum_{i-m < |y| \leq m} \frac{1}{K} \sum_{k=1}^K \widehat{\psi}(y, \xi_N(\overline{S_k(N)})) \times \\ \sum_{j \in S_k(N)} \frac{1}{\#S_k(N)} \mathbb{I}\{Y^j = y, |f_{PA}(X^j, \xi_N(\overline{S_k(N)})) - y| > i\}.$$

В случае строго состоятельной оценки $\hat{\psi}(\cdot)$ в [2] (Булинский, Ракитько, 2014) был установлен критерий, гарантирующий

$$\widehat{Err}_K(f_{PA}, \xi_N) \rightarrow Err(f) \text{ п.н., } N \rightarrow \infty.$$

Пусть также для оценки $\hat{\psi}(\cdot)$ справедливо

$$\sqrt{\#S_k(N)}(\hat{\psi}(y, \xi_N(S_k(N))) - \psi(y)) = O_P(1), \quad N \rightarrow \infty,$$

и

$$\sup_{y \in \mathbb{Y}, N \in \mathbb{N}, k \in \{1, \dots, K\}} \mathbb{E} \left(\hat{\psi}(y, \xi_N(S_k(N))) \right)^4 < \infty.$$

Теорема

Пусть $(m_N)_{N \in \mathbb{N}}$ – последовательность натуральных чисел такая, что $m_N \leq q$ для $q = \lfloor N/K \rfloor$ и

$$m_N \rightarrow \infty, \quad m_N/N \rightarrow 0, \quad \text{as } N \rightarrow \infty.$$

Выберем такую последовательность положительных чисел $\varepsilon = (\varepsilon_N)_{N \in \mathbb{N}}$, что $\varepsilon_N \rightarrow 0$ и $m_N^{1/2} \varepsilon_N \rightarrow \infty$ as $N \rightarrow \infty$. Тогда при $N \rightarrow \infty$

$$\sqrt{m_N} \left(\widehat{Err}_K(f_{PA,\varepsilon}, \xi_N) - Err(f) \right) \xrightarrow{law} Z_{0,\sigma^2} \sim \mathcal{N}(0, \sigma^2).$$

- [1] Булинский А.В., Ракитько А.С. (2014). Оценивание небинарного случайного отклика. *Доклады РАН*, **455**(6), 623–627.
- [2] Bulinski, A., Butkovsky, O., Sadovnichy, V., Shashkin, A., Yaskov, P., Balatskiy, A., Samokhodskaya, L., Tkachuk, V. (2012). Statistical methods of SNP data analysis and applications. *Open Journal of Statistics*. **2**(1), 73–87.
- [3] Röllin, A. (2013). Stein's method in high dimensions with applications. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques* **49**(2), 529–549.