

# Статистические свойства некоторых процедур сжатия данных

Бзикадзе Андрей Важевич, гр. 15.M03-мм

Санкт-Петербургский государственный университет  
Кафедра статистического моделирования

Научный руководитель: к.ф.-м.н. Некруткин В.В.  
Рецензент: исследователь Советкин Е.А.



Санкт-Петербург  
06 июня 2017

Одна из рассматриваемых в ВКР процедур сжатия данных.

- $\mathbb{S} \stackrel{\text{def}}{=} \{1, 2, \dots, S\}$  — множество книг.
- Стопка книг.
- Начальный порядок  $\Xi_0$ .
- Из стопки случайная книга перекладывается наверх.

Одна из рассматриваемых в ВКР процедур сжатия данных.

- $\mathbb{S} \stackrel{\text{def}}{=} \{1, 2, \dots, S\}$  — множество книг.
- Стопка книг.
- Начальный порядок  $\Xi_0$ .
- Из стопки случайная книга перекладывается наверх.

Итеративная процедура.

- $\{\eta_i\}_{i=1}^{\infty}$  — последовательность **названий** случайных книг.
- $\{\Xi_i\}_{i=0}^{\infty}$  — последовательность состояний **стопки** случайных книг.
- $\{\xi_i\}_{i=1}^{\infty}$  — последовательность **положений** случайных книг.

Одна из рассматриваемых в ВКР процедур сжатия данных.

- $\mathbb{S} \stackrel{\text{def}}{=} \{1, 2, \dots, S\}$  — множество книг.
- Стопка книг.
- Начальный порядок  $\Xi_0$ .
- Из стопки случайная книга перекладывается наверх.

Итеративная процедура.

- $\{\eta_i\}_{i=1}^{\infty}$  — последовательность **названий** случайных книг.
- $\{\Xi_i\}_{i=0}^{\infty}$  — последовательность состояний **стопки** случайных книг.
- $\{\xi_i\}_{i=1}^{\infty}$  — последовательность **положений** случайных книг.

$$\eta_i = \Xi_{i-1}[\xi_i].$$

3
4
1
5
2

$\Xi_0$

$$\eta_1 = 4$$

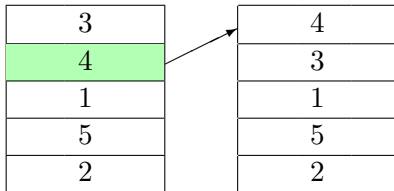
$$\xi_1 = 2$$

3
4
1
5
2

 $\Xi_0$

$$\eta_1 = 4$$

$$\xi_1 = 2$$



$\Xi_0$

$\Xi_1$

# Пример. $S = 5$

$$\begin{array}{lcl} \eta_1 & = & 4 \\ \xi_1 & = & 2 \end{array}$$

$$\begin{array}{lcl} \eta_2 & = & 4 \\ \xi_2 & = & 1 \end{array}$$

$$\begin{array}{lcl} \eta_3 & = & 2 \\ \xi_3 & = & 5 \end{array}$$

3
4
1
5
2

$\Xi_0$

4
3
1
5
2

$\Xi_1$

4
3
1
5
2

$\Xi_2$

2
4
3
1
5

$\Xi_3$



## Применение «Book Stack»-преобразования:

- Рябко Б.Я. (1980): алгоритм **сжатия** данных под названием «метод стопки книг».
- Bentley J.L. et al. (1986): тот же алгоритм под названием «**Move To Front**».
- Рябко Б.Я. et al. (2003–2004): тест для проверки свойств **генераторов псевдослучайных чисел** под названием «Book Stack».

Равносильны:

- $\mathbb{H}_0 : \eta_i$  независимы и равномерно распределены на  $\mathbb{S}$ .
- $\mathbb{H}_0^* : \xi_i$  независимы и равномерно распределены на  $\mathbb{S}$ .

Два критерия для проверки  $\mathbb{H}_0$ : при применении к «исходной» и к «преобразованной» выборке.

Альтернатива в (Бзикадзе А.В., Некруткин В.В, 2016) и бакалаврской ВКР:

$\mathbb{H}_1 : \{\eta_i\}_{i \geq 1} \text{ — н.о.р. с неравномерным распределением.}$

При больших  $n$ :

- Критерий  $\chi^2$  к  $\{\eta_i\}_{i=1}^n$  (как правило) мощнее, чем к  $\{\xi_i\}_{i \geq 1}$ .
- Критерий отношения правдоподобия к  $\{\eta_i\}_{i=1}^n$  мощнее, чем к  $\{\xi_i\}_{i \geq 1}$ .

Альтернатива в (Бзикадзе А.В., Некруткин В.В, 2016) и бакалаврской ВКР:

$\mathbb{H}_1 : \{\eta_i\}_{i \geq 1}$  — н.о.р. с **неравномерным** распределением.

При больших  $n$ :

- Критерий  $\chi^2$  к  $\{\eta_i\}_{i=1}^n$  (как правило) мощнее, чем к  $\{\xi_i\}_{i \geq 1}$ .
- Критерий отношения правдоподобия к  $\{\eta_i\}_{i=1}^n$  мощнее, чем к  $\{\xi_i\}_{i \geq 1}$ .

В магистерской ВКР:

$\mathbb{H}_1 : \{\eta_i\}_{i \geq 1}$  — **эргодическая** однородная марковская цепь (ЭОМЦ) со стационарным **равномерным** на  $\mathbb{S}$  распределением и переходной матрицей  $\mathbf{P}^{(\eta)}$ :  $\text{tr}(\mathbf{P}^{(\eta)}) \neq 1$ .

Пусть входная последовательность  $\{\eta_i\}_{i \geq 1}$  — ОМЦ с матрицей переходов  $\mathbf{P}^{(\eta)} = (p_{ij})$ .

## Предложение

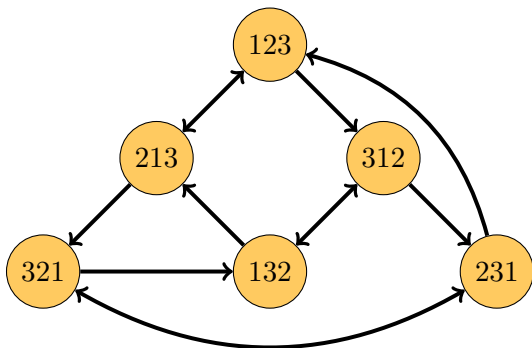
- 1 *Последовательность  $\{\Xi_n\}_{n \geq 1}$  образует ОМЦ.*
- 2 *Если входная ОМЦ  $\{\eta_n\}_{n \geq 1}$  — эргодическая, то последовательность  $\{\Xi_n\}_{n \geq 1}$  имеет ровно один непериодический эргодический класс и, быть может, несколько несущественных состояний.*
- 3 *Если же  $p_{ij} > 0$  при всех  $i, j$ , то несущественных состояний нет, т.е.  $\{\Xi_n\}_{n \geq 1}$  эргодическая.*

# Марковское свойство $\Xi_i$ . Пример

Рассмотрим графическое изображение ОМЦ  $\{\Xi_i\}_{i \geq 1}$ .

Входная ОМЦ  $\{\eta_i\}_{i \geq 1}$  — эргодическая с матрицей  $\mathbf{P}^{(\eta)} = (p_{ij})$ .

$S = 3$ ,  $p_{ij} > 0$  для всех  $i, j$ .

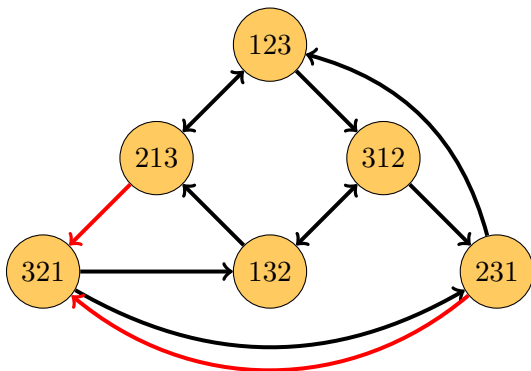


# Марковское свойство $\Xi_i$ . Пример

Рассмотрим графическое изображение ОМЦ  $\{\Xi_i\}_{i \geq 1}$ .

Входная ОМЦ  $\{\eta_i\}_{i \geq 1}$  — эргодическая с матрицей  $\mathbf{P}^{(\eta)} = (p_{ij})$ .

$S = 3$ ,  $p_{23} = 0$ ,  $p_{33} = 0$  и остальные  $p_{ij} > 0$ .



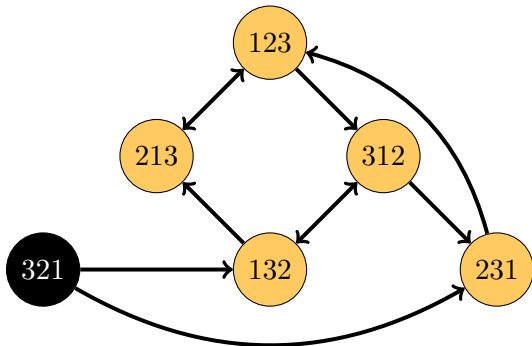
# Марковское свойство $\Xi_i$ . Пример

Рассмотрим графическое изображение ОМЦ  $\{\Xi_i\}_{i \geq 1}$ .

Входная ОМЦ  $\{\eta_i\}_{i \geq 1}$  — эргодическая с матрицей  $\mathbf{P}^{(\eta)} = (p_{ij})$ .

$S = 3$ ,  $p_{23} = 0$ ,  $p_{33} = 0$  и остальные  $p_{ij} > 0$ .

$(3, 2, 1)^T$  — несущественное состояние ОМЦ  $\{\Xi_i\}_{i \geq 1}$ .





Пусть  $\{\eta_i\}_{i \geq 1}$  — ЭОМЦ с м.п.  $\mathbf{P}^{(\eta)} = (p_{ij})$ .

Обозначим  $(\pi_\alpha, \alpha \in \mathfrak{S}_S)$  — стационарное распределение  $\Xi_i$ ,  
 $\tau_k = \tau_k(n) = \mathbb{I}_k(\xi_1) + \dots + \mathbb{I}_k(\xi_n)$  для  $1 \leq k \leq S$  и положим

$$s_k \stackrel{\text{def}}{=} \sum_{j=1}^S \sum_{\substack{\alpha \in \mathfrak{S}_S \\ \alpha_k = j}} \pi_\alpha p_{\alpha[1]j}.$$

## Теорема

Для всех  $k \in 1 : S$  при  $n \rightarrow \infty$

- $\mathbb{P}(\xi_n = k) \rightarrow s_k$ .
- $\mathbb{E}(\tau_k/n - s_k)^2 = O(1/n)$ .

Распределение с вероятностями  $s_j$  будем обозначать  $\mathcal{R}$ .

Обозначим  $\mathbf{s} = (s_1, \dots, s_S)^T$  и

$$a_{k,\ell} = \begin{cases} s_k(1 - s_k), & \text{при } k = \ell, \\ -s_k s_\ell, & \text{при } k \neq \ell. \end{cases}$$

## Теорема

$$\mathcal{L} \left( \sqrt{n} \left( \tau_n^{(\xi)} / n - \mathbf{s} \right) \right) \Rightarrow \mathcal{N}(0, \Sigma_\xi),$$

где матрица  $\Sigma_\xi$  размеров  $S \times S$  имеет компоненты

$$(\Sigma_\xi)_{k,\ell} = a_{k,\ell} + C_{k,\ell},$$

а  $C_{k,\ell}$  некоторые константы, зависящие только от  $\mathbf{P}^{(\eta)}$ .

## Теорема

*Пусть  $\{\eta_i\}_{i \geq 1}$  образуют ЭОМЦ со стационарным **равномерным** распределением и матрицей переходных вероятностей  $\mathbf{P}^{(\eta)}$ :  $\text{tr}(\mathbf{P}^{(\eta)}) \neq 1$ . Тогда предельное распределение  $\mathcal{R}$  не совпадает с равномерным.*

Без предположения, что  $\text{tr}(\mathbf{P}^{(\eta)}) \neq 1$ , утверждение Теоремы, вообще говоря, неверно.

Альтернатива:

$\mathbb{H}_1 : \{\eta_i\}_{i \geq 1}$  — эргодическая однородная марковская цепь со стационарным **равномерным** распределением и переходной матрицей  $\mathbf{P}^{(\eta)}$ :  $\text{tr}(\mathbf{P}^{(\eta)}) \neq 1$ .

## Теорема

*Критерий  $\chi^2$  против альтернативы  $\mathbb{H}_1$*

- При применении к «входным»  $\{\eta_i\}_{i \geq 1}$  — несостоятельный.
- При применении к «выходным»  $\{\xi_i\}_{i \geq 1}$  — состоятельный.

Модель марковской цепи:

- Задано:  $0 < \delta < 1$ .
- Матрица переходных вероятностей:  $\mathbf{P} = (p_{ij})$ , где  $p_{ii} = \delta$  и  $p_{ij} = (1 - \delta)/(S - 1)$  при  $i \neq j$ .

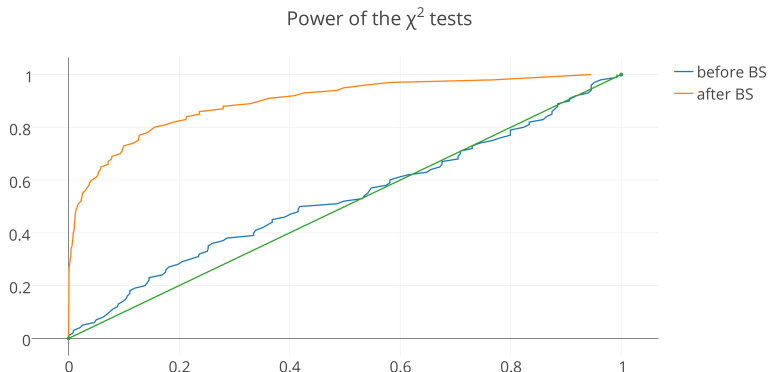
Моделирование:

- Вихрь Мерсенна.
- $n$  — размер выборки,  $m$  — количество выборок.
- Критерий  $\chi^2$  с  $S - 1$  степенью свободы:  $m$  штук  $P$ -значений.

Цель:

- Сравнение мощности критерия  $\chi^2$  «до» и «после» преобразования «Book Stack».

Параметры:  $S = 3$ ,  $n = 10^4$ ,  $\delta = 1/S + 0.01$ ,  $m = 100$ .



**Рис.:** Мощности критериев  $\chi^2$  до/после «Book Stack». Для  $\chi^2$  «до»  $P$ -значение критерия Колмогорова-Смирнова равно 0.48, «после» — меньше  $2.2 \cdot 10^{-16}$ .

- Получено обобщение теоретико вероятностных результатов бакалаврской ВКР на случай, когда  $\{\eta_i\}_{i \geq 1}$  образуют ЭОМЦ.
- При альтернативной гипотезе, что  $\{\eta_i\}_{i \geq 1}$  — ЭОМЦ с матрицей переходов  $\mathbf{P}^{(n)}$ :  $\text{tr}(\mathbf{P}^{(n)}) \neq 1$ , критерий  $\chi^2$  «до» преобразования несостоятельный, а такой же критерий «после» — состоятельный.

Также получены следующие результаты

- Рассмотрено **обобщение** «Book Stack»-преобразования, при котором **сохраняются** многие теоретико вероятностные результаты.
- Рассмотрен другой тест «Order» для проверки той же гипотезы  $\mathbb{H}_0$  и теоретически обоснована **бесперспективность** его применения.