

# Метод Монте-Карло по схеме марковской цепи для оценки вероятности редких событий в задачах биоинформатики

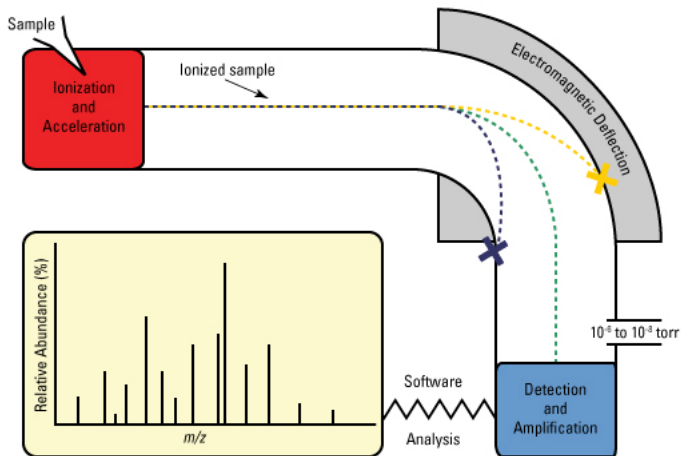
Абрамова Анастасия Николаевна, гр. 15.M03-мм

Санкт-Петербургский государственный университет  
Математико-механический факультет  
Кафедра статистического моделирования

Научный руководитель: к.ф.-м.н., д. Коробейников А. И.  
Рецензент: разработчик ПО Тарасов А. Л.



# Методы масс-спектрометрии



Пептид  $P$  с массой  $M$ , состоящий из  $k$  аминокислот

- вектором масс аминокислот  $m = (m_1, \dots, m_k)$
- матрицей структуры  $\mathbb{H}$  размера  $k \times r$ .

Ожидаемый масс-спектр пептида  $P$  имеет вид  $S = \mathbb{H}m$ .

Определим множество  $\mathcal{M}$ :

$$\mathcal{M} = \{m = (m_1, \dots, m_k) \mid m_i > 0, \sum_{i=1}^k m_i = M\}.$$

Тогда  $\mathcal{M}$  и фиксированная матрица  $\mathbb{H}$  описывают множество пептидов одинаковой массы и химической структуры.

Рассмотрим экспериментальный спектр:

$$\tilde{S} = (\tilde{s}_1, \dots, \tilde{s}_\ell), \quad \tilde{s}_i > 0.$$

Мерой похожести экспериментального спектра  $\tilde{S}$  и масс-спектра  $\mathbb{H}m^*$  назовем значение  $t = \overline{\text{Score}}(\tilde{S}, \mathbb{H}m^*)$ .

## Определение

*Статистической значимостью значения  $t$  будем называть вероятность*

$$p = \mathbb{P}(\overline{\text{Score}}(\tilde{S}, \mathbb{H}m) \geq t) = \mathbb{P}(\text{Score}(m) \geq t) = \mathbb{P}(m \in \mathcal{S}),$$

*предполагая что  $m$  имеет равномерное распределение на  $\mathcal{M}$ .*

В приложениях  $p \approx 10^{-20}$ .

В работе (Mohimani *et al.*, 2013) предложен алгоритм MS-DPR, вычисляющий оценку статистической значимости.

Недостатки данного алгоритма состоят в том, что

- его точность неизвестна,
- размер выборки, используемой для построения оценки, задается заранее,
- оценки являются смещенными вниз.

Таким образом, целью работы является:

- построение оценки  $\hat{p}$ ,
- вычисление ее дисперсии  $\sigma_{\hat{p}}^2$  и построение доверительного интервала для  $\hat{p}$ ,
- оценка достаточного размера выборки для получения  $\hat{p}$  заданной точности.

# Оценка по методу существенной выборки

Рассмотрим выборку  $m_1, \dots, m_N \sim f$ , где  $f$  — плотность  $U(\mathcal{M})$ . Пусть  $g$  — плотность некоторого распределения  $\mathcal{G}$ .

Оценка по методу существенной выборки для вероятности  $p = \mathbb{P}(m \in \mathcal{S})$ :

$$\hat{p}_{IS} = \frac{1}{N} \sum_{i=1}^N \frac{f(m_i)}{g(m_i)} \mathbb{I}_{\mathcal{S}}(m_i).$$

Пусть  $g(m) \propto w(\text{Score}(m))f(m)$ , тогда:

$$\hat{p}_{IS} = \frac{\sum_{i=1}^N \mathbb{I}_{\mathcal{S}}(m_i) / w(\text{Score}(m_i))}{\sum_{i=1}^N 1 / w(\text{Score}(m_i))}.$$

Для построения  $\hat{p}_{IS}$  необходимо уметь моделировать случайные величины из распределения  $\mathcal{G}$ .

Алгоритм Метрополиса–Гастингса позволяет построить марковскую цепь  $m_1, \dots, m_N$  со стационарным распределением  $\mathcal{G}$ .

## Предложение

*Если для марковской цепи выполняется закон больших чисел, то оценка  $\hat{p}_{IS}$  является состоятельной оценкой  $p$ .*



## Выбор весов $w$

Когда  $g(m) \propto w(\text{Score}(m))f(m)$ , выбор плотности  $g$  сводится к выбору весов  $w$ .

Если значения функции  $\text{Score}$  дискретны, то выбор

$$w(x) \propto \frac{1}{\mathbb{P}(\text{Score}(m) = x)}$$

уменьшит дисперсию  $\hat{p}$ .

Алгоритм Ванга-Ландау (Iba *et al.*, 2014) является модификацией метода Метрополиса-Гастингса: в процессе моделирования цепи строится оценка весов.

- 1 Построение оценки  $\hat{w}$  алгоритмом Ванга–Ландау.
- 2 Построение марковской цепи со стационарным распределением  $g(m) \propto \hat{w}(\text{Score}(m))f(m)$  алгоритмом Метрополиса–Гастингса.
- 3 Построение оценки

$$\hat{p}_{IS} = \frac{\sum_{i=1}^N \mathbb{I}_{\mathcal{S}}(m_i) / \hat{w}(\text{Score}(m_i))}{\sum_{i=1}^N 1 / \hat{w}(\text{Score}(m_i))}.$$

На практике, вместо множества  $\mathcal{M}$  можно рассматривать его дискретное подмножество: зафиксируем некоторое  $m \in \mathcal{M}$  и рассмотрим вектора

$$\mathcal{M}_{d,m} = \{\tilde{m} = m + dv \mid \tilde{m} \in \mathcal{M}, \sum_{i=1}^k v_i = 0, v_i \in \mathbb{Z}\}.$$

## Теорема

- ❶ *Марковская цепь, построенная методом Метрополиса-Гастингса на множестве  $\mathcal{M}_{d,m}$ , является эргодической.*
- ❷ *Марковская цепь, построенная методом Метрополиса-Гастингса на множестве  $\mathcal{M}$ , является Харрис-эргодической.*

## Теорема (Flegal et al., 2013)

Пусть  $\hat{\lambda}_N$  – оценка дисперсии вдоль траектории  $\lambda_p$ :

$\hat{\lambda}_N \xrightarrow{\text{п. н.}} \lambda_p$ ,  $\hat{\sigma}_N^2 \xrightarrow{\text{п. н.}} \sigma_p^2$ . Предположим, что

$\sqrt{N}(\hat{p}_N - p) \xrightarrow{d} \mathcal{N}(0, \sigma_p^2)$ ,  $N \rightarrow \infty$ .

Обозначим

- $N_\epsilon = \inf \left\{ N > 0 : 2z_{\delta/2}\hat{\sigma}_N/\sqrt{N} \leq \epsilon\hat{\lambda}_N \right\}$
- $C_N = (\hat{p}_N - z_{\delta/2}\hat{\sigma}_p^2/\sqrt{N}; \hat{p}_N + z_{\delta/2}\hat{\sigma}_p^2/\sqrt{N})$

Тогда при  $N \rightarrow \infty$  и  $\epsilon \rightarrow 0$  моделирование прекратится с вероятностью 1 и  $\mathbb{P}(p \in C_{N_\epsilon}) \rightarrow 1 - \delta$  при  $N \rightarrow \infty$ .

Оценка  $\hat{p}_{IS}$ , полученная по траектории марковской цепи  $m_1, \dots, m_N$ , имеет вид

$$\hat{p}_{IS} = \frac{1}{N} \sum_{i=1}^N h(m_i).$$

Дисперсия такой оценки вычисляется как

$$\sigma_p^2 = \frac{1}{N} \sum_{k=-(N-1)}^{N-1} \left(1 - \frac{|k|}{N}\right) \text{cov}(h(m_i), h(m_{i+k})).$$

Методы оценки дисперсии: метод перекрывающихся средних (Jones *et al.*, 2006), спектральные оценки (Hobert *et al.*, 2002).

В предложенном алгоритме длина траектории  $N$  марковской цепи увеличивается последовательно, поэтому при использовании классических методов возникают проблемы

- хранения траектории,
- высокая трудоемкость.

Используется рекурсивная оценка по «правилу выбора трапеции» (Chan and Yau, 2014).

**Свойства:**

- вычислительная сложность рекурсивного пересчета  $O(1)$ ,
- затрачивает  $O(1)$  памяти.

Были построены:

- полученные оценки  $\hat{p}_{IS}$ ,
- оценки по методу Монте-Карло  $\hat{p}_{MC}$ ,
- оценки по алгоритму MS-DPR  $\hat{p}_{DPR}$ ,

для пептидов:

- простой структуры: линейных и циклических (Mohimani *et al.*, 2013),
- сложной, циклической с разветвлениями структуры (*Surfactin*).

Также для оценок  $\hat{p}_{IS}$  и  $\hat{p}_{MC}$  были сосчитаны оценки дисперсий  $\hat{\sigma}_{IS}^2$ ,  $\hat{\sigma}_{MC}^2$  и построены 95% доверительные интервалы.

Для оценок по методу Монте-Карло, построенных с достаточно большим  $N$ :

- 1 оценки  $\hat{p}_{IS}$  лежат в доверительных границах оценок по методу Монте-Карло.
- 2 оценки  $\hat{p}_{DPR}$  в большинстве случаев выходят за границы доверительного интервала по методу Монте-Карло.



Для оценок, построенных по одинаковому размеру выборки:

Пептид	$\hat{p}_{IS}$	$\hat{\sigma}_{IS}^2$	$\hat{\sigma}_{MC}^2$	$\hat{\sigma}_{MC}^2/\hat{\sigma}_{IS}^2$
PPAEDSQK	$4.87 \cdot 10^{-7}$	$2.09 \cdot 10^{-10}$	$4.94 \cdot 10^{-7}$	2358.98
GQGDPGSNPKNK	$4.70 \cdot 10^{-7}$	$2.33 \cdot 10^{-10}$	$1.49 \cdot 10^{-7}$	639.49
GEEEPSQGQK	$1.03 \cdot 10^{-6}$	$1.23 \cdot 10^{-9}$	$7.89 \cdot 10^{-7}$	642.19
<i>Surfactin</i>	$1.18 \cdot 10^{-5}$	$1.15 \cdot 10^{-7}$	$1.00 \cdot 10^{-5}$	86.96
(10, 20, 40)	$1.84 \cdot 10^{-3}$	$5.47 \cdot 10^{-4}$	$1.88 \cdot 10^{-3}$	3.43

Результаты демонстрируют, что отношение  $\hat{\sigma}_{MC}^2/\hat{\sigma}_{IS}^2$  увеличивается с уменьшением значения оцениваемой вероятности  $p$ .

- 1 Был предложен способ оценки статистической значимости меры схожести двух спектров.
- 2 Проведено сравнение полученного алгоритма с MS-DPR.
- 3 Его код был написан на C++ и интегрирован в Dereplicator — метод идентификации пептидных спектров.
- 4 Написана статья и подана на конференцию по алгоритмической биоинформатике (WABI, 2017).