

# Скрытые Марковские модели переменного порядка для анализа данных ChIP-seq

Атаманова Анна, кафедра системного программирования СПбГУ, [anne.atamanova@gmail.com](mailto:anne.atamanova@gmail.com)

31 марта 2015 г.

## Аннотация

Здесь нужно кратко описать суть работы и результаты.

## 1 Введение

Дезоксирибонуклеиновая кислота (ДНК) является своеобразным кодом жизни. Эта молекула хранит и передает генетическую программу развития и функционирования живого организма. В то же время, все функции ДНК зависят от ее соединений с белками. Поэтому, изучение ДНК-белковых взаимодействий актуально и привлекательно. Chip-seq (chromatin immunoprecipitation - sequencing)[1] является одним из современных методов, позволяющим выделить участки ДНК связанные с конкретным белком (одинаково применим к разным белкам). Однако, по понятным причинам (сложный биологический эксперимент), погрешность данного метода не может быть нулевой, и безрассудная вера ему лишена смысла. По этому, обычно, к результатам подобных методов накладывается вероятностная модель. Конечно, это добавляет ряд существенных ограничений. Однако, в качестве неоспоримого плюса можно привести тот факт, что хорошо подобранная модель позволяет понять природу данных и изучить их свойства.

Итого, нашей задачей является нахождение модели по последовательности чисел выдаваемых Chip-seq.

В настоящее время, в качестве семейства искомых моделей, активное применение находит НММ (Hidden Markov Model)[2] второго порядка с Пуассоновским испусканием. Данное семейство допускает предположение о том, что каждое состояние (наличие/отсутствие белка в заданной части генома) зависит только от одного предыдущего. Можно ограничиться и более лояльным допущением о том, что состояние зависит от  $n$  предыдущих состояний, однако такое допущение резко увеличивает сложность модели ( $O(2^n)$  параметров). Также, сложность заключается в подборе этого  $n$  и переобучении в случае, если не все состояния имеют одинаковые длины контекстов зависимости. Последнее замечание подводит к идее использования VОНММ (Variable Order Hidden Markov Model)[3]

## 2 Скрытые марковские модели переменного порядка

### Список литературы

- [1] David S Johnson, Ali Mortazavi, Richard M Myers, and Barbara Wold. Genome-wide mapping of in vivo protein-dna interactions. *Science*, 316(5830):1497–1502, 2007.
- [2] Lawrence Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [3] Y Wang, Lizhu Zhou, and Jianhua Feng. Mining complex time-series data by learning Markovian models. In *International Conference on Data Mining*, 2006.