

Анализ подходов к моделированию данных метилирования ДНК

Сергей Лебедев

Руководитель: Олег Шпынов

JetBrains

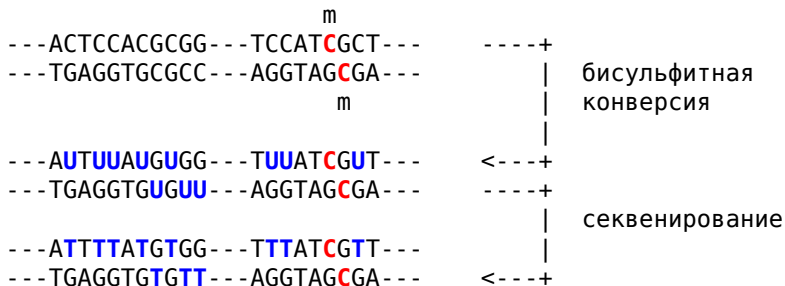
21 сентября, 2013

Эпигенетический минимум

Метилирование ДНК —

химическая модификация, добавляющая метильную группу к цитозину.

Бисульфитное секвенирование ДНК



0

5

4 кол-во метилированных цитозинов

8 покрытие

Долины метилирования

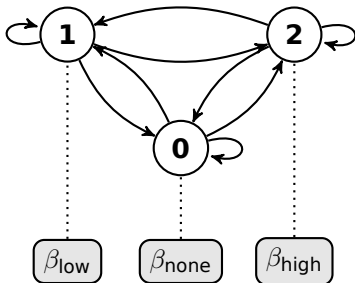
- DNA Methylation Valleys [**XSL+13**] — длинные слабо метилированные участки генома.
- Что особенного в **DMV**:
 - ассоциированы с транскрипционными факторами и генами развития,
 - пересекают некоторые промоторы lncRNA,
 - консервативны между различными видами,
 - остаются слабо метилированными при дифференциации клетки и
 - становятся сильно метилированными в раковых клетках.
- Как локализовать DMV, используя данные бисульфитного секвенирования?
- Как сегментировать метилом на слабо- и сильно- метилированные участки?

Сегментирование метилома

Уровень метилирования

$$\beta = \frac{mC}{C} \in [0, 1]$$

Скрытая марковская модель с Гауссовыми вероятностями испусканий [SMB⁺11]



$$\beta_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

Сегментирование метилома с GHMM

- С помощью модели мы можем приписать каждому цитозину одну из трех меток:
 - 0 метилирование отсутствует,
 - 1 низкий уровень метилирования,
 - 2 высокий уровень метилирования.
- Всем остальным нуклеотидам припишем метку, соответствующую ближайшему цитозину.
- Пример:

```
      m
---ACTCCACGCGG---TCCATCGCT---
      000000000000      0000111100
              ^
              |
            почти DMV
```

- Уровень метилирования – это величина из $[0, 1]$, в то время как GHMM описывает величину из \mathbb{R} . Можно:
 - моделировать испускания с помощью бета-распределения,
 - преобразовать уровень метилирования:

$$M = \log \frac{\beta}{1 - \beta} \in \mathbb{R}$$

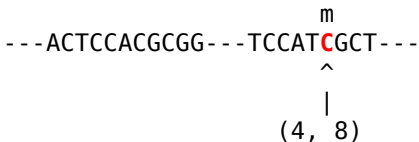
$$\Gamma = -\log(1 - \beta) \in [0, +\infty)$$

- При переходе от пары (mC, C) к уровню метилирования теряется информация о покрытии, например:

$$\beta_i = 1/2 = \beta_j = 10/20 = 0.5$$

Альт. подход к сегментированию метилома

- Будем моделировать не уровень метилирования, а пары (m_C, C) напрямую:



- Таким образом, значение m_C в каждой позиции определяется покрытием в этой позиции C и некоторой вероятностью “метилирования” p_i :

$$m_C \sim \mathcal{B}(C, p_i) \quad \text{for } i \in \{\text{high, low, none}\}$$

- Покрытие в каждой позиции C мы считаем известным с точностью до распределения.

Как моделировать покрытие?

- Категориальным распределением на множестве $\{1, \dots, k\}$, где k — максимальное покрытие.
- Смесь распределений Пуассона:
 - Если предположить, что каждый нуклеотид в геноме длины N покрыт одинаковым числом фрагментов n , то вероятность получить k ридов в некоторой позиции:

$$B(k; n, p = 1/N) \underset{N \rightarrow \infty}{\approx} \text{Pois}(k; \frac{n}{N})$$

- В реальном мире покрытие почти никогда не равномерно, поэтому разумно использовать смесь распределений Пуассона.

Биномиальная смесь для парных наблюдений

- Моделирует независимые пары (mC, C) , полученные в результате бисульфитного секвенирования.
- Предполагает распределение покрытия в каждой позиции известным: категориальное или смесь распределений Пуассона.
- Использует набор компонент, аналогичный GHMM [SMB+11]: {high, low, none}.

- Провести сравнительный анализ результатов РВММ и GHMM на данных [SMB⁺11, XSL⁺13].
- Обобщить модель для поиска различно метилированных регионов в результатах нескольких экспериментов.



¹http://www.flickr.com/photos/maggie_oc/9210581429



Michael B Stadler, Rabih Murr, Lukas Burger, Robert Ivanek, Florian Lienert, Anne Schöler, Erik van Nimwegen, Christiane Wirbelauer, Edward J Oakeley, Dimos Gaidatzis, et al.

Dna-binding factors shape the mouse methylome at distal regulatory regions.

[Nature](#), 2011.



Wei Xie, Matthew D Schultz, Ryan Lister, Zhonggang Hou, Nisha Rajagopal, Pradipta Ray, John W Whitaker, Shulan Tian, R David Hawkins, Danny Leung, et al.

Epigenomic analysis of multilineage differentiation of human embryonic stem cells.

[Cell](#), 2013.