

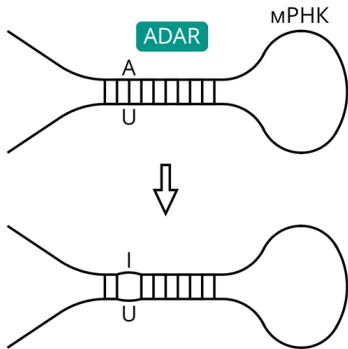
# Аналитические методы поиска сайтов редактирования РНК в данных RNA-seq

Щукина Ирина Алексеевна

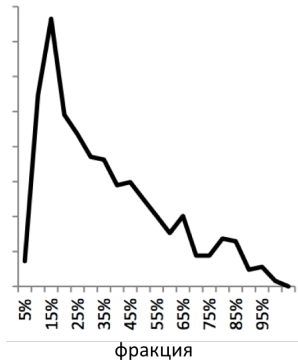
научные руководители: А. Канапин, А. Самсонова

СПб АУ НОЦНТ РАН

14 июня 2017 г.



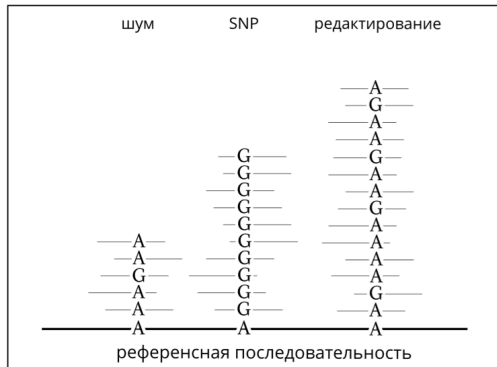
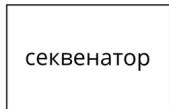
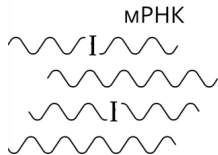
Клеточная машинерия и секвенатор интерпретируют инозин (I) как гуанин (G)



Ожидаемый вид распределения фракций отредактированных нуклеотидов

# Редактирование в RNA-seq

AG = редактирование + технический шум + SNPs



- REDITools
- RED
- GIREMI
- Фильтрация по фракции/покрытию/качеству...
- Использование DNA-seq или нескольких реплик для отбора сайтов

Принцип работы:

- 1 Оценивает вероятность замены каждого типа по всему эксперименту

$$H_0 : p_{AG} = \frac{n_G}{n_G+n_A}$$

$$H_1 : p_{AG} < \frac{n_G}{n_G+n_A}$$

- 2 Точный тест Фишера для каждой позиции
- 3 Поправка на множественное тестирование

Цель: разработать метод поиска сайтов редактирования РНК в данных RNA-seq

Задачи:

- 1 Изучить и модифицировать существующий алгоритм REDITools denovo
- 2 Применить байесовский подход к фильтрации технического шума
- 3 Разработать алгоритм для фильтрации SNP
- 4 Сравнить подходы с аналогами и произвести отбор

- Знаем, что редактированию соответствует только замена типа AG. Остальные типы – шум, будем оценивать вероятности ошибки только по ним
- Дополнительный тест для фильтрации гомозиготных SNPs (оригинальный метод никак не затрагивает эту проблему).

- 1 Оценим по шуму априорное распределение фракций

$$p_{AG} \sim B(\alpha_0, \beta_0)$$

- 2 Для всех сайтов типа AG найдём апостериорное распределение

$$p_{AG} | n_A, n_G \sim B(\alpha_0 + n_G, \beta_0 + n_A)$$

- 3 Оценим вероятность

$$P[p_{AG} | n_A, n_G \leq \text{thr}]$$

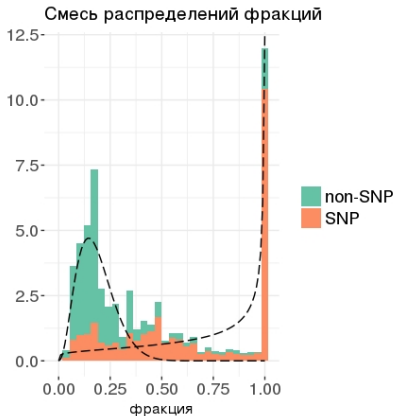
- 4 Установим FDR и отберём сайты в соответствии с ним

Данный метод позволяет избавиться от технического шума



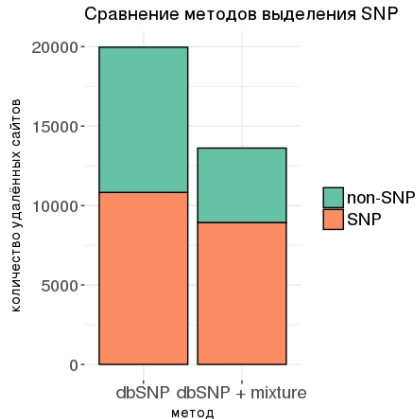
## Методы:

- Фильтрация против dbSNP — базы всех встречающихся SNP (избыточно)
- Смоделируем смесь распределений (EM алгоритм) и удалим позиции, из “SNP-компоненты”



## Методы:

- Фильтрация против dbSNP — базы всех встречающихся SNP (избыточно)
- Смоделируем смесь распределений (EM алгоритм) и удалим позиции, из “SNP-компоненты”



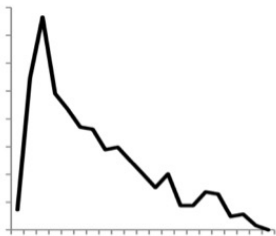
В роли baseline-метода использован биномиальный тест

Данные:

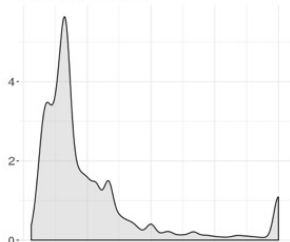
- GM12878 — достоверно известны SNP
- ICE-seq — экспериментальный протокол для поиска сайтов редактирования. Существует один датасет с набором (не полным) подтверждённых сайтов редактирования. Сайты разбиты на подгруппы в зависимости от степени валидированности.

# Получаемые распределения фракций

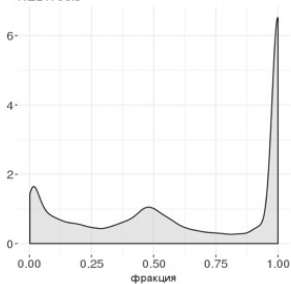
Ожидаемый вид распределения



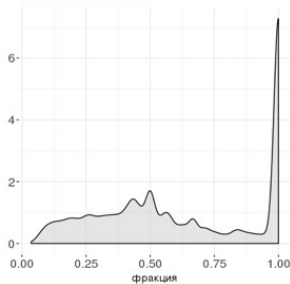
Байесовский подход



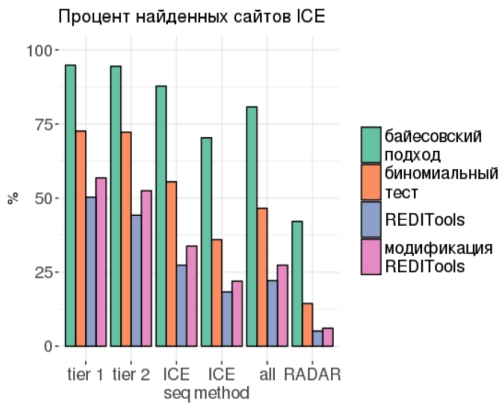
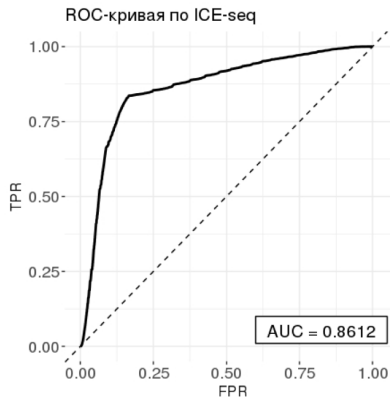
REDITools



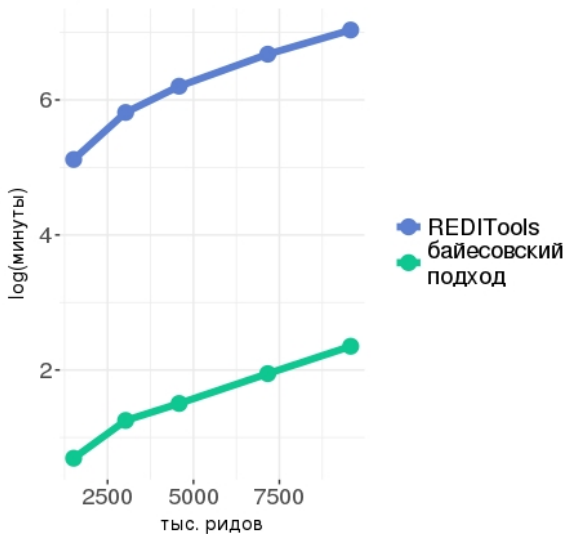
Биномиальный тест



# Сравнение методов



## Сравнение производительности методов



Помимо основного алгоритма реализованы:

- возможность фильтрации и определение стренда с помощью аннотации генома
- метрика для ранжирования сайтов на основе kpp (значимое различие для TP сайтов ICE:  $p < 2.2 \cdot 10^{-16}$ )

- модифицированы оценка вероятности и фильтрация SNP в существующем алгоритме REDITools
- применение байесовского подхода дало ощутимое улучшение результатов фильтрации технического шума
- разработан алгоритм для фильтрации SNP на основе базы dbSNP и разделения позиций по фракциям
- на основе отобранных методов реализован пайплайн для нахождения сайтов редактирования в данных RNA-seq



Доля SNP в найденных сайтах

