

NATIONAL RESEARCH UNIVERSITY
HIGHER SCHOOL OF ECONOMICS
FACULTY OF MATHEMATICS

Project Proposal

Clustering of Multidimensional Random Variables to Improve HMM Sequence Alignment Accuracy Кластеризация многомерной случайно величины для поднятия точности выравнивания строк.

Student: Denis Grachev, BMT181

Scientific advisor:
PhD candidate,
Timofey Prodanov

Moscow 2022

Contents

1	Abstract	3
2	Introduction	3
3	Preliminaries	4
3.1	Strings	4
3.2	Task	5
3.3	Additive scoring functions	6
3.4	Pair Hidden Markov Model	6
3.5	Clustering	6
3.6	Stepwise iterative maximum likelihood algorithm	6
3.7	Likelihood	7
3.8	Algorithm	8
4	Main results	8
4.1	Clustering	8
4.2	Longshot	8

1 Abstract

Whole-genome sequencing is an complex and important task of bioinformatics. Different technologies can generate data of different nature. Most popular technologies, such as illumina, have low error rate, but give limited information about genome. Newer techologies, for example Pacific Biosciences and Oxford Nanopore, can give more details about genome, but have higher error rate. Combination of these methods can help improve accuracy of genome sequencing. During this work a new method of string clustering was developed, to identify different data profiles and new functionality to existing tools were added, to work with multiple datasets.

2 Introduction

Bioinformatics is an interdisciplinary science that aims to develops methods and software tools for understanding biological data. One of the ways to model haploid genome is to present it as pair of sequences or strings over the alphabet $\{A, C, G, T\}$. Modern technologies of reading genome do not sequence it as one continious string, but a number of random overlapping substrings that are called reads. As within one biological species genomes coincide almost completely, it is convenient to determine one reference genome for one species and identify for every individual it's deviations from the reference. These single nucleotide variants are called SNVs. Taking into cosideration these facts and the fact that various errors happen during all stages of the process, a number of problems appears for example:

- **Genome assembly** is process of deciphering genome using reads obtained from it.
- **Sequence alignment** is process of arranging sequences to spot similarities between them.
- **Variant calling** is process of identifying SNVs of an individual based on reads aligned on the reference genome.

The most commonly used technologies nowadays, such as illumina, allow to sequence reads of length 200-500 bp. Sequencing human genome using such short length reads has many limitations. First, due to diploidy of humans, it is important to obtain long-range haplotype information. This might be difficult with short-reads provided by illumina. Secondly 3.6% of human genome consists of long highly repetative duplicated regeions that can not be uniquely aligned, which lowers accuracy of SNVs. Third-generation single-molecule sequencing (SMS) techologies, such as Oxford Nanopore Technologies allow to genereate longer reads of length 10-30 kb. This technology might help overcome limitations that short-reads have. Variant calling tools that were developed for short-reads do not show high accuracy when applied to long-reads sequenced, due to significant difference in error rates and type of errors between these two types of reads. Also these tools process data in short windows of a few hundred bases lengths and not designed to agregate haplotype information present in long-reads, which might be crucial for distinguishing true deviations from errors.

A number of new methods of variant calling were developed to work with long-reads, such as Deep Variant and Longshot that uses deep learning and pair-Hidden Markov model algorithms to effectively work with such data. Accuracy of these methods can

be improved by agrigating information from multiple sequence data sets together. An example of effectiveness of such method is Spades genome assembly tool, which uses short and long read data together.

Longshot works effectively with one type of reads. To improve it's sencetivity and accuracy various methods of collecting data from diiferent sources can be applied. Previous research show that sequencing reads of different profiles and clustering them by their profiles can improve variant calling accuracy.

During this work we are going to extend functionality of Longshot tool to work effectively with multiple datasets. For that we have to develop clustering algorithm for read profiles and implement such feature into Longshot.

3 Preliminaries

3.1 Strings

Definition 3.1. String is a sequence of letters of finite size.

String of length l over alphabet $A = \{1 \dots m\}$ is a map $s : \{1 \dots l\} \rightarrow A$. Usually elements of A are denoted as characters for convenience.

Definition 3.2. Alignment of strings is a way of representing them to spot similarities. Alignment of strings s_1 and s_2 of lengths l_1 and l_2 respectively, over alphabet A is a pair of strings \hat{s}_1 and \hat{s}_2 of length l over alphabet $A \sqcup \{'-\'}$, such that there exists increasing functions $f_i : \{1 \dots l_i\} \rightarrow \{1 \dots l\}, i \in \{1, 2\}$ such that $\hat{s}_i \circ f_i = s_i$ and $\text{Im}(f_i) = \hat{s}_i^{-1}(A)$.

Letter $'-'$ represents gap in a string. String \hat{s}_i represents string s_i with inserted letter $'-'$, that was not present in alphabet A , into some places, and function f_i maps indexes of letters in s_i to corresponding indexes in \hat{s}_i .

Example 3.1. Alignment of strings $s_1 = CABC AABA$ and $s_2 = ABADBBAD$ over alphabet $\{A, B, C, D\}$.

$$\begin{array}{c} \text{Initial strings.} \\ \left\| \begin{array}{c} s_1 \\ s_2 \end{array} \right| \begin{array}{cccccccc} C & A & B & C & A & A & B & A \\ A & B & A & D & B & B & A & \end{array} \right\| \end{array}$$

$$\begin{array}{c} \text{Aligned strings.} \\ \left\| \begin{array}{c} \hat{s}_1 \\ \hat{s}_2 \end{array} \right| \begin{array}{cccccccc} C & A & B & C & - & A & A & B & A \\ - & A & B & - & A & D & B & B & A \end{array} \right\| \end{array}$$

Definition 3.3. Define score function for an alignment $S(\hat{s}_1, \hat{s}_2)$. The *better* alignment the more score it gets. Different score functions will be introduced later.

Definition 3.4. We are interested in the alignment with the highest possible score. So we define score between two strings as

$$S(s_1, s_2) = \max_{\hat{s}_1, \hat{s}_2} S(\hat{s}_1, \hat{s}_2)$$

Definition 3.5. Substring is continuous piece of a given string.

For a string s of length l , sub-string s_s is a string of length l_s , such that there exists a function $f : \{1 \dots l_s\} \rightarrow \{1 \dots l\}$ such that

$$f(i) = i + d$$

$$s \circ f = s_s.$$

Definition 3.6. We want to find similarities between reference string and reads. For this we define local alignment and score for a string and a substring.

For a string s_1 and s_2 of lengths l_1, l_2 correspondingly, define string-sub-string score as

$$S_s(s_1, s_2) = \max\{S(s_s, s_2) | s_s \text{ is a sub-string of } s_1\}$$

and corresponding alignment \hat{s}_1, \hat{s}_2 .

Definition 3.7. For a string s of length l and set of strings $R = \{s_1 \dots s_n\}$ of lengths $\{l_1 \dots l_n\}$ correspondingly, multiple alignment is tuple $\hat{s}, \hat{s}_1 \dots \hat{s}_n$, of strings of length l over alphabet $A \sqcup \{'-\'}$, such that $\sum_{i=1}^n S(\hat{s}, \hat{s}_i)$ is maximal.

Definition 3.8. Set of reads R for string s of length l and rate r is

$$R = \{s_s | \text{length of } s_s > l, S_s(s, s_s) < r\}$$

3.2 Task

Given reference string s_r and reads R for an unknown target string s_t , we know that $S(s_r, s_t) < D$ and want to find s_t .

Plan:

1. Make multiple alignment of R over s_r .
2. Estimate most likely difference between s_r and s_t .

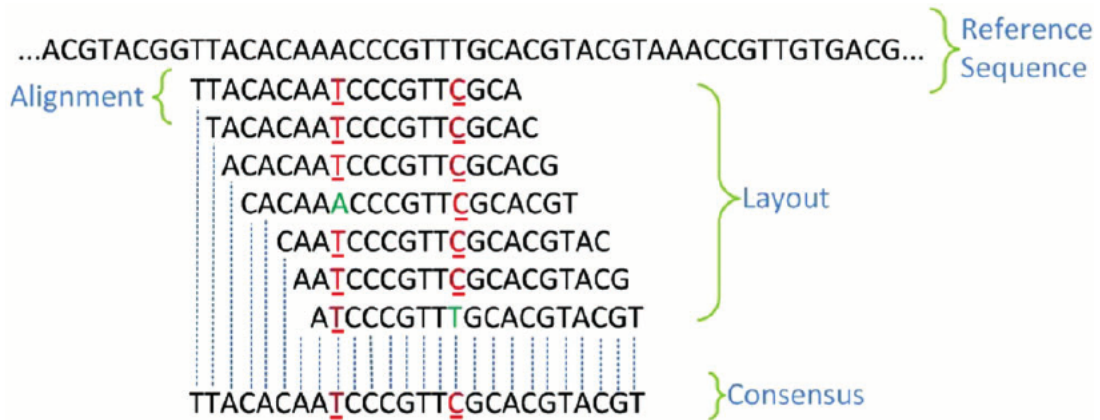


Figure 1: Example of reference string, target string and reads.

3.3 Additive scoring functions

One of the ways to score an alignment is to add award for matches and subtract penalty for mismatches and gaps.

Definition 3.9. For a given matrix $G \in \mathbb{R}^{|A \cup \{-'\}} \times |A \cup \{-'\}|$ and $p \in \mathbb{R}$ score of alignment \hat{s}_1, \hat{s}_2 is

$$S(\hat{s}_1, \hat{s}_2) = \sum_{i=1}^l G_{\hat{s}_1(i), \hat{s}_2(i)}.$$

Matrix G stores predefined penalties for mismatches and gaps and encouragement for matches.

Sometimes fines penalty for first gap is higher than prolonging a continious gap, because one continious gap is more likely to appear than several small gaps.

Best alignment for such score function can be found using Needleman-Wunsch algorithm.

3.4 Pair Hidden Markov Model

Each step of pairwise alignment can be assigned to one of the three states $\{M, X, Y\}$, where M is a match, X is a gap in s_1 , Y is a gap in s_2 . Given transition probabilities each alignment can be assigned a probability. The most likely alignment can be found using Viterbi algorithm.

Add about begin and end states?

3.5 Clustering

Definition 3.10. Clustering algorithm aims to group points together into predefined number of sets.

Clustering algorithm is a map

$$\text{cluster}(X, m) \rightarrow C$$

$$X = \{x_i | x_i \in \mathbb{R}^d, i \in (1 \dots n)\}, m \in \mathbb{N}$$

$$C = \{c_i | c_i \in (1 \dots m), i \in (1 \dots n)\}$$

where m is number of clusters and n is number of points.



Figure 2: Example of clustering for $d = 2$, $m = 2$, color represents class.

3.6 Stepwise iterative maximum likelihood algorithm

To distinguish different read profiles a new method of clustering was developed.

3.7 Likelihood

Assume that the data consists of different multidimensional normal distributions. Let's denote data as

$$X = \{x_i | x_i \in \mathbb{R}^d, i \in (1 \dots n)\}, m \in \mathbb{N}$$

where m is number of clusters and n number of points. Assume that we have some clustering C , we will try to increase it's likelihood.

Denote i th cluster as ω_i and it's estimated parameters as θ_i

$$\Theta = \{\theta_i | i \in (1 \dots n)\}.$$

Then probability of x in i th cluster

$$p(x|\Theta) = p(x|\omega_i, \theta_i)P(\omega_i).$$

Denote clusters as $\chi_1 \dots \chi_l \subset X$, than logarithm of probability for all points in cluster i is

$$\begin{aligned} L_i &= \sum_{x \in \chi_i} \log(p(x|\omega_i, \theta_i)P(\omega_i)) \\ &= \log \left(\frac{\exp \left(\frac{-1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right)}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \right) + n_i \log(P(\omega_i)) \\ &= -\frac{1}{2} n_i d - \frac{n_i d}{2} \log(2\pi) - \frac{n_i}{2} \log |\Sigma_i| + n_i \log \frac{n_i}{n}. \end{aligned}$$

Where μ_i is mean value and Σ_i - covariation of i th cluster. Overall likelihood is

$$L = \sum_{i=1}^l L_i.$$

Move \hat{x} from χ_i to χ_j , then

$$\begin{aligned} \Delta L_i &= -\frac{1}{2} \log |\Sigma_i| + \frac{n_i - 1}{2} \log \left(1 - \frac{(\hat{x} - \mu_i)^T \Sigma_i^{-1} (\hat{x} - \mu_i)}{n_i - 1} \right) + \\ &\quad + \log \frac{n_i}{n} - (n_i - 1) \left(\frac{d}{2} + 1 \right) \log \frac{n_i - 1}{n_i} \\ \Delta L_j &= -\frac{1}{2} \log |\Sigma_j| - \frac{n_j + 1}{2} \log \left(1 + \frac{(\hat{x} - \mu_j)^T \Sigma_j^{-1} (\hat{x} - \mu_j)}{n_j + 1} \right) + \\ &\quad + \log \frac{n_j}{n} + (n_j + 1) \left(\frac{d}{2} + 1 \right) \log \frac{n_j + 1}{n_j}. \end{aligned}$$

$$\Delta L = \Delta L_i + \Delta L_j$$

3.8 Algorithm

Main idea

1. Initialize clusters (randomly or using another algorithm)
2. Iterate over all points
 - (a) Move point to a cluster, such that overall likelihood increases the most. (With most ΔL_j)
 - (b) Update clusters and their parameters.
3. Repeat step 2 while it makes changes.

Advantages

- After every step overall likelihood increases.
- This implies that the cycle will end.

Problems

- Updating parameters after every step is very slow.

Transition of one point does not change estimated parameters significantly, so to reduce iteration complexity, we will update parameters once in k points. Also to avoid possible loops we will iterate over points in different orders.

Final algorithm

1. Initialize clusters and estimate their parameters
2. Devide X into p random disjoined groups $g_1 \dots g_p$.
3. Loop c from 1 to p .
 - (a) Loop x over g_c .
 - (b) Let x currently be in cluster i .
 - i. If $n_i \leq 1$, then pass to next point.
 - ii. Calculate $\delta_j = \begin{cases} \Delta L_j, & j \neq i \\ \Delta L_i, & j = i \end{cases}$
 - iii. Trasfer x to $\text{argmax}(\delta_j)$ cluster.
 - (c) Update parameters.
 - (d) If overall likelihood decreased, revert changes.
4. If any changes were made, repeat step 2.

4 Main results

4.1 Clustering

Examples of clustering reads with likelyhood information

4.2 Longshot

Added longshot features

References

- [1] Longshot is variant calling tool for long reads based on pair-Hidden Markov Model.
<https://www.nature.com/articles/s41467-019-12493-y>