

FEDERAL STATE AUTONOMOUS EDUCATIONAL INSTITUTION  
FOR THE HIGHER EDUCATION  
NATIONAL RESEARCH UNIVERSITY “HIGHER SCHOOL OF ECONOMICS”  
FACULTY OF MATHEMATICS

**Grachev Denis Vadimovich**

# **Clustering of Multidimensional Random Variables to Improve HMM Sequence Alignment Accuracy**

**Project proposal**

Scientific supervisor:  
Prodanov Timofey Petrovich

Moscow 2022

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Clustering . . . . .	3
1.2	Strings . . . . .	3
<b>2</b>	<b>Task</b>	<b>4</b>

# 1 Introduction

## 1.1 Clustering

Given  $X = \{x_i | x_i \in \mathbb{R}^d, i \in (1 \dots n)\}$  and  $m \in \mathbb{N}$ , where  $n$  is the number of points,  $m$  - number of clusters.

Clustering algorithm takes  $X$  and  $m$  and outputs  $C = \{c_i | c_i \in (1 \dots m), i \in (1 \dots n)\}$ .



Figure 1: Example of clustering for  $d = 2$ ,  $m = 2$ , color represents class.

## 1.2 Strings

**Definition 1.1.** 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 1.2.** 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  $f_i : \{1 \dots l_i\} \rightarrow \{1 \dots l\}$  such that  $\hat{s}_i|_{\hat{s}_i^{-1}(A)} \circ f_i = s_i$ .

*Remark.*  $\text{Im}(f_i) = \hat{s}_i^{-1}(A)$

**Example 1.1.** Alignment of strings  $s_1 = CABC AABA$  and  $s_2 = ABADBBAD$  over alphabet  $\{A, B, C, D\}$ .

$$\left\| \begin{array}{c|cccccccc} s_1 & C & A & B & C & A & A & B & A \\ s_2 & A & B & A & D & B & B & A & \\ \hline \hat{s}_1 & C & A & B & C & - & A & A & B & A \\ \hat{s}_2 & - & A & B & - & A & D & B & B & A \end{array} \right\|$$

**Definition 1.3.** For given matrix  $G \in \mathbb{R}^{|A| \times |A|}$  and  $p \in \mathbb{R}$  score of alignment  $\hat{s}_1, \hat{s}_2$  is

$$S(\hat{s}_1, \hat{s}_2) = \sum_{i=1}^l \delta_i, \text{ where } \delta_i = \begin{cases} g_{\hat{s}_1(i)\hat{s}_2(i)}, & \hat{s}_1(i) \neq - \text{ and } \hat{s}_2(i) \neq - \\ p, & \end{cases}$$

**Theorem 1.** If  $G$  is symmetric and  $g_{ij} = \begin{cases} 0, & i = j \\ > 0, & \end{cases}$  and  $p > 0$ , then we can define metric for strings over alphabet  $A$  as

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

*Proof.*

□

**Definition 1.4.** For a string  $s$  of length  $l$ , sub-string  $s_s$  is a string of length  $l_s$ , such that there exists an function

$$\begin{aligned} f &: \{1 \dots l_s\} \rightarrow \{1 \dots l\} \\ f(i) &= i + d \\ s \circ f &= s_s \end{aligned}$$

**Definition 1.5.** For a string  $s_1$  and  $s_2$  of lengths  $l_1, l_2$  correspondingly, define string-sub-string score as

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

and corresponding alignment  $\hat{s}_1, \hat{s}_2$  are pair of strings of lengths  $l$  over alphabet  $A \sqcup \{-\}$  such that there exists increasing functions  $f_1 : \{1 \dots l_1\} \rightarrow \{1 \dots l\}$

**Definition 1.6.** 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 minimal.

**Definition 1.7.** Set of reads  $R$  for string  $s$  of length  $l$  and rate  $r$  is

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

## 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  $s_t$ .

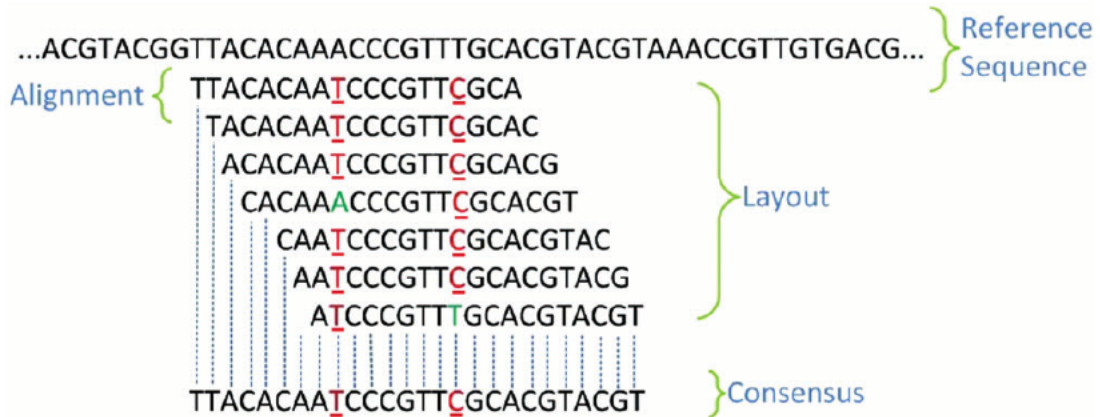


Figure 2: Example of reference string, target string and reads.