

Hierarchical clustering

BioInformatic: Report Module 3
Groupe 2

Thomas FEUILLETIN

thomas.feuilletin@etudiant.univ-rennes1.fr

Florian EPAIN

florian.epain@etudiant.univ-rennes1.fr

May 20, 2022

Abstract

We are close to get something nice but we put our hearts into it.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 2 |
| 2 | Compare sequences | 3 |
| 2.1 | Step 1 | 3 |
| 2.2 | Step 2 | 5 |
| 3 | Hierarchical clustering | 5 |
| 3.1 | Initialization | 5 |
| 3.2 | Visualization | 8 |
| 3.3 | Clustering | 12 |
| 4 | Application to hemoglobin | 15 |
| 4.1 | Read the gene sequence from a fasta file | 15 |
| 4.2 | Cluster of gene sequences | 16 |
| 5 | Tutorial | 16 |
| 5.1 | Install Rust | 16 |
| 5.2 | Install Foam | 16 |

1 Introduction

This project's goal is to compare a gene and its associated proteins across multiple species. We will process some basic algorithm with basic data sample. We're using a method called clusterize by agglomerative.

First, we have to compare each sequence. Using the [Levenshtein distance](#), we will know how far -or close- is a sequence from another. With this well-known distance, we group the closest sequences and repeat the process with the remaining and this new cluster.

Tests are performed on a hemoglobin dataset. Algorithms coded in rust. Results observed in terminal and/or in *.md* by foam vc-code extension. The section 5 is a brief tutorial on how to install rust and run our program and how to install foam and visualize the result graph.

Our project outline, gave by subject:

We want to verify whether similar species do have similar genes, and whether these genes produce similar proteins. The first part of the project consists in implementing a generic analysis workflow on any kind of sequence. The second part is an application of the workflow on hemoglobin. Section 2 focuses on determining the distance between a pair of genes (resp. of proteins). Section 3 groups similar sequences by performing hierarchical clustering according to their distances. Eventually, section 4 focuses on hemoglobin genes and proteins

This record is written with the [v0.9.0](#) of our program.
Some final results revealed:

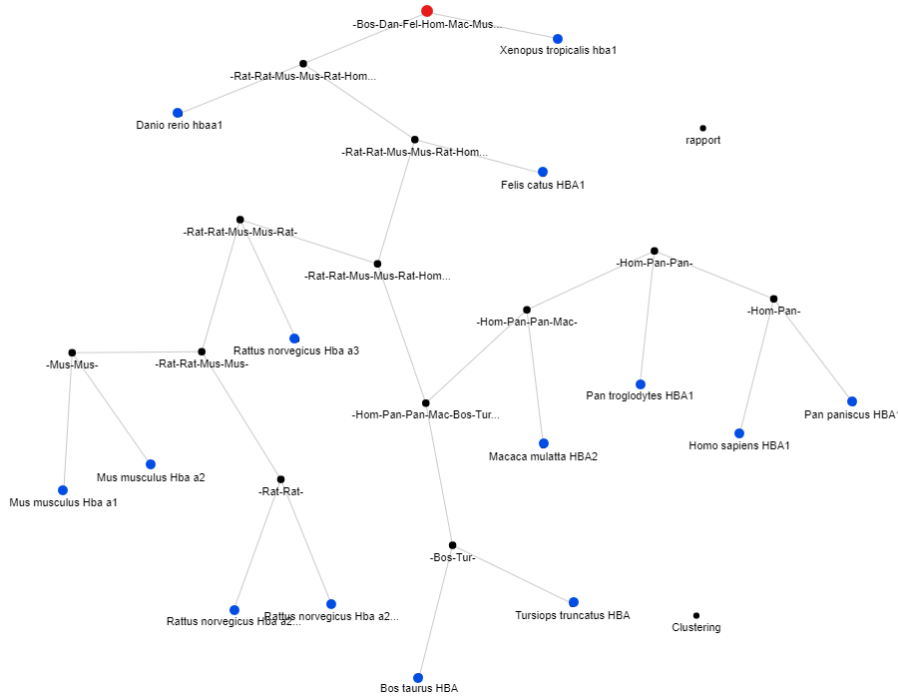


Figure 1: Hierarchically clustered hemoglobin dataset

2 Compare sequences

This section aims at computing the distance between two sequences.

2.1 Step 1

To do so, let's create a structure to represent Sequence.

```
1  #[derive(Clone)]
2  pub struct Sequence {
3      name: String,
4      seq: String,
5  }
```

Let's implement this structure with constructor

```
1  impl Sequence {
2
3      pub fn new() -> Sequence{
4          Sequence {
5              name: String::new(),
6              seq: String::new(),
7          }
8      }
9
10     pub fn new_with_string(name: String, seq: String) -> Sequence {
11         Sequence {
12             name,
13             seq,
14         }
15     }
16
17     pub fn new_with_sequence(name: String, s: Sequence) -> Sequence {
18         Sequence {
19             name,
20             seq: s.seq,
21         }
22     }
23
24     pub fn get_name(&self) -> String {
25         self.clone().name
26     }
27
28     ...
29
30 }
```

We now code a pretty close Levenshtein distance, as described in wikipedia.

```
1  impl Sequence {
2      ...
3
4      /**
5       * -----
6       * ^ ----- ^
7       * ^ ~ATCCCTGCA~ ^
8       * ^ ~ATCCCTGCT~ ^
9       * ^ ----- ^
10      * ^  AGCCCTGCT  ^
11      * -----
12      *
13      * Levenshtein distance
14      * return the distance between the two sequences
15      */
16      pub fn levenshtein_distance(&self, other_seq: &Sequence) -> usize {
17          // convert our seq to Vec<Char>
18          let s1: String = self.seq.to_string();
19          let s2: String = other_seq.seq.to_string();
20          let seq_1: Vec<char>= s1.chars().collect();
21          let seq_2: Vec<char>= s2.chars().collect();
22
23          // initialization of result matrix
24          let mut d: Vec<Vec<usize>> = vec![vec![0;s2.len()+1]; s1.len()+1];
25          // println!("d.len(): {}", d.len());
26
27          // write a meter
28          for i in 1..d.len() {
29              d[i][0] = i
30          }
31          for j in 1..d[0].len() {
32              // println!("j: {}, d[0].len(): {}", j, d[0].len());
33              d[0][j] = j
34          }
35
36          // idk what it determine
37          let mut substitution_cost: usize;
38          for i in 1..d.len(){
39              for j in 1..d[i].len() {
40                  // let substitution_cost: i8;
41                  // maybe to expensive to recreate every case
42                  if seq_1[i-1] == seq_2[j-1] {
43                      substitution_cost = 0;
44                  }else {
45                      substitution_cost = 1;
46                  }
47
48                  d[i][j] = lib::minimum_3(d[i-1][j] + 1,                // deletion
49                                          d[i][j-1] + 1,                // insertion
50                                          d[i-1][j-1] + substitution_cost)// substitution
51              }
52          }
53
54          d[seq_1.len()][seq_2.len()]
55      }
56 }
```

2.2 Step 2

A simple example of these methods

```
1 let seq_1 = Sequence::new_with_string("seq_1".to_string(), "ATTACG".to_string());
2 println!("{}", seq_1.levenshtein_distance(&seq_1));
```

Here the distance calculated by our `levenshtein_distance()` method:

| | levenshtein_distance() | | | | | | |
|--------|------------------------|--------|--------|--------|--------|-------|-------|
| | ATTACG | ATATCG | ACCCCG | GGGGAA | TTTACG | ATTAC | ATATC |
| ATTACG | 0 | 2 | 3 | 6 | 1 | 1 | 3 |
| ATATCG | 2 | 0 | 3 | 6 | 3 | 3 | 1 |
| ACCCCG | 3 | 3 | 0 | 6 | 4 | 4 | 4 |
| GGGGAA | 6 | 6 | 6 | 0 | 6 | 5 | 6 |
| TTTACG | 1 | 3 | 4 | 6 | 0 | 2 | 4 |
| ATTAC | 1 | 3 | 4 | 5 | 2 | 0 | 2 |
| ATATC | 3 | 1 | 4 | 6 | 4 | 2 | 0 |

3 Hierarchical clustering

Hierarchical clustering is a kind of recursive classification that consists in organizing the set of elements into subsets included in to each others in a tree-like structure. There are two main approaches for determining this organization:

- Agglomerative approach.
- Divisive approach.

We will focus on the Agglomerative approach.

3.1 Initialization

For the creation of this structure, we choose to use vector.

Note: I added a name attribute which slow down the whole clustering process but helps with the foam representation

```
1 /**
2  * we're able to create a pure copy of a cluster with a new ref by the Clone derive
3  *
4  * sub_clusters correspond to all the clusters contained in this clusters
5  * elements contains all element in this cluster, eventually from the sub_clusters
6  */
7 #[derive(Clone)]
8 pub struct ClusterOfSequence {
9     name: String,
10    sub_clusters: Vec<ClusterOfSequence>,
11    elements: Vec<Sequence>,
12 }
```

The cluster's name depends on what it contains

- One sequence -> Sequence's name
- Many sequences -> The three first letter of all sequence name

```
1  impl ClusterOfSequence
2  {
3      pub fn new(element: Sequence) -> ClusterOfSequence
4      {
5          ClusterOfSequence {
6              name: element.clone().name,
7              sub_clusters: Vec::new(),
8              elements: vec![element],
9          }
10     }
11
12     pub fn new_with_sequences(elements:Vec<Sequence>) -> ClusterOfSequence
13     {
14         ClusterOfSequence {
15             name:
16             {
17                 let mut res = String::new();
18                 res.push('-');
19                 for elem in &elements {
20                     // to avoid Err :
21                     if elem.name.len() <= 3 {
22                         res.push_str(&elem.name[0..elem.name.len()]);
23                     }else {
24                         res.push_str(&elem.name[0..3]);
25                     }
26                 }
27                 res.push('-');
28             }
29             res
30             sub_clusters:
31             {
32                 let mut res = Vec::new();
33                 for e in &elements {
34                     res.push(ClusterOfSequence::new(e.clone()))
35                 }
36                 res
37             },
38             elements,
39         }
40     }
41
42     ...
43 }
```

Iterating through all the sequence names of a cluster is time consuming, one way to counter this is to give the cluster a name when needed (Foam Graph Representation)

```
1  impl ClusterOfSequence
2  {
3      ...
4
5      /**
6       * ClusterOfSequence without reference &?
7       */
8      pub fn new_with_clusters(clusters_1: ClusterOfSequence,
9                             clusters_2: ClusterOfSequence) -> ClusterOfSequence
10     {
11         ClusterOfSequence {
12             name:
13             {
14                 // we have to make sure there is no doublon
15                 // to avoid conflict and wrong relation in foam
16                 let mut foo= String::new();
17                 foo.push('-');
18                 for elem in &clusters_1.elements {
19                     if elem.name.len() <= 3 {
20                         foo.push_str(&elem.name[0..elem.name.len()]);
21                     }else {
22                         foo.push_str(&elem.name[0..3]);
23                     }
24                     foo.push('-');
25                 }
26                 for elem in &clusters_2.elements {
27                     if elem.name.len() <= 3 {
28                         foo.push_str(&elem.name[0..elem.name.len()]);
29                     }else {
30                         foo.push_str(&elem.name[0..3]);
31                     }
32                     foo.push('-');
33                 }
34                 foo
35             },
36             sub_clusters: vec![clusters_1.clone(), clusters_2.clone()],
37
38             // get all seq of clust_1 and clust_2 to elements
39             elements:
40             {
41                 let mut n = clusters_1.clone().elements;
42                 n.extend(clusters_2.clone().elements);
43                 n
44             }
45         }
46     }
47
48     ...
49 }
```

3.2 Visualization

Step 5:

```
1 let seq_0_test = Sequence::new_with_string("ATTACG".to_string(), "ATTACG".to_string());
2 let seq_1_test = Sequence::new_with_string("ATATCG".to_string(), "ATATCG".to_string());
3 let seq_2_test = Sequence::new_with_string("GCCGAG".to_string(), "GCCGAG".to_string());
4 let seq_3_test = Sequence::new_with_string("ACCCCG".to_string(), "ACCCCG".to_string());
5 let seq_4_test = Sequence::new_with_string("TCCCCG".to_string(), "TCCCCG".to_string());
6 let sequences_test = vec![seq_0_test, seq_1_test, seq_2_test, seq_3_test, seq_4_test];
7 let mut bio_cluster = ClusterOfSequence::new_with_sequences(sequences_test);
```

Here the simple newick method

```
1 pub fn get_newick_old(&self) -> String
2 {
3     let mut res = String::from("(");
4     // we could set "!= 0" to "!= 1"
5     // to avoid displaying clusters containing a single sequence
6     if self.sub_clusters.len() != 0 {
7         for i in 0..self.sub_clusters.len() {
8             res.push_str(&self.sub_clusters[i].get_newick_old().as_str());
9             // e != self.elements.last().unwrap() -> != not implement
10            if i != self.sub_clusters.len()-1 {
11                res.push(',');
12            }
13        }
14    }
15    }else {
16        // this else occurs when a cluster contains only sequences
17        for e in 0..self.elements.len() {
18            res.push_str(&self.elements[e].name.as_str());
19            if e != self.elements.len()-1 {
20                res.push(',');
21            }
22        }
23    }
24    res.push(')');
25    res
26 }
27 }
```

Step 7:

```
1 println!("bio_cluster: \n{}", bio_cluster.get_newick_old());
```

results: bio_cluster: ((ATTACG),(ATATCG),(GCCGAG),(ACCCCG),(TCCCCG))

We tried to have a cool looking newick rep in the terminal but it not goes like intended...

```
1  /**
2   * TODO recheck step 6 and 7
3   * if you want to get a cleaner look change :
4   * .elements[e].seq to .elements[e].name
5   * if you want to get the precise sequence in the representation :
6   * .elements[e].name to .elements[e].seq
7   */
8  pub fn get_newick(&self) -> String
9  {
10     self.get_newick_with_space(0)
11 }
12
13 fn get_newick_with_space(&self, space: i32) -> String
14 {
15     let mut res = String::new();
16     if self.sub_clusters.len() != 0 {
17         for i in 0..self.sub_clusters.len() {
18             // res.push_str("/");
19             res.push_str(
20                 &self.sub_clusters[i].get_newick_with_space(space+1)
21                     .as_str());
22
23             // e != self.elements.last().unwrap()
24             // TODO : check this if else usefulness
25             if i != self.sub_clusters.len()-1 {
26                 res.push_str("\n");
27             }else {
28                 res.push_str("\n");
29             }
30
31         }
32     }else {
33         for e in 0..self.elements.len() {
34             if e ==0 || e == &self.elements.len() -1{
35                 for _count in 0..space {
36                     res.push_str("-");
37                 }
38             }
39             // .elements[e].seq or .elements[e].name depending
40             res.push_str(&self.elements[e].name.as_str());
41
42             if e != self.elements.len()-1 {
43                 res.push_str("\n");
44             }
45
46         }
47     }
48     // res.push_str("\n");
49     res
50 }
```

```

1 println("bio_cluster: \n{}", bio_cluster.get_newick());
2 bio_cluster.clusterize_agglomerative();
3 println("bio_cluster: \n{}", bio_cluster.get_newick());

```

imagining that clusterize method exists, it goes:

```

1 bio_cluster:
2 -ATTACG
3 -ATATCG
4 -GCCGAG
5 -ACCCCG
6 -TCCCCG
7
8 clusterize...
9 ended in 1.1177ms
10
11 bio_cluster:
12 ---ACCCCG
13 ---TCCCCG
14
15 --GCCGAG
16
17 --ATTACG
18 --ATATCG

```

And now for the visualization by graph (not [professional graph](#) but np). We opted for the FOAM vs-code extension which reads links between markdown and create a cool constellation. Go to Section 5, for more information on how to install foam in vs-code.

What does the create_foam_rep() method ?

Create a serie of markdown file and each file contains:

- ▶ `[[wikilinks]]` to subclusters
- ▶ The list of elements contained

Make sure you install foam vscode, and have theses folder :

- ▶ `_layout`
- ▶ `.vscode`
- ▶ `.foam`
- ▶ `assets`

At the end, you will prompt the foam graph in vscode by : `ctrl+p ">Foam: Show Graph"`

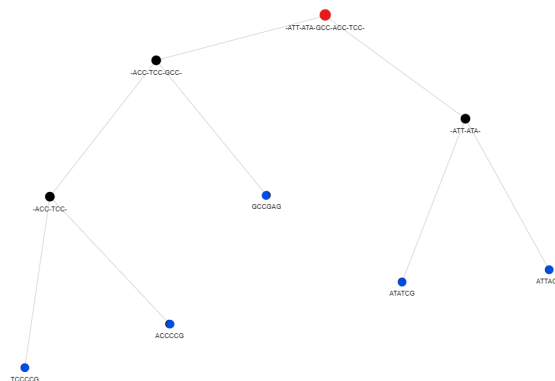


Figure 2: Hierarchically clustered random small dataset

```

1  pub fn create_foam_rep(&self, folder_name: &str) // -> Result<()>
2  {
3
4      let mut foo =
5      "----
6      \n title: name
7      \n type: bio_inf_results
8      \n ----
9      \n".to_string();
10
11     // Add up all sub_cluster links
12     for counter in 0..self.sub_clusters.len() {
13         let mut bar = "\n [[wikilink]]"
14             .replace("wikilink", &self.sub_clusters[counter].clone().name);
15
16         if counter != self.sub_clusters.len()-1 {
17             bar.push_str(", ");
18         }
19         foo.push_str(&bar);
20
21         // now create the md file associated with this sub cluster
22         self.sub_clusters[counter].create_foam_rep(folder_name);
23     }
24
25     // separator
26     foo.push_str("\n");
27
28     // if you prefer have the sequence prompted
29     // change element.name to element.seq
30     for element in &self.elements {
31         let mut bar = "element"
32             .replace("element", &element.name);
33         if element.name != self.elements.last().unwrap().name {
34             bar.push_str(", ");
35         }
36         foo.push_str(&bar);
37     }
38
39     // write this foobar to a .md file
40
41     let name = self.clone().name;
42     // change title
43     foo = foo.replace("name", &name);
44
45     // create docs folder if it did not exist
46     let folder_path = "./foam_rep/docs/folder_name/"
47         .replace("folder_name", folder_name);
48     fs::create_dir_all(&folder_path);
49     let md_path = "./foam_rep/docs/folder_name/name.md"
50         .replace("folder_name", folder_name)
51         .replace("name", &name);
52
53     fs::write(md_path, foo)
54         .expect("Unable to write file");
55
56     // Ok(())
57 }

```

3.3 Clustering

Linkage method creation

```
1 pub fn linkage(&self, a_cluster: ClusterOfSequence) -> Option<f32>
2 {
3     // covert the division by 0
4     if self.elements.len() == 0 || a_cluster.elements.len() == 0 {
5         None
6         // f32::MAX;
7     }else {
8         let mut result = 0.00;
9         let mut counter = 0.00;
10        /*
11         * compare every e1: Sequence (in the self called Cluster)
12         * with every e2: Sequence (in a_cluster)
13         * by the levenshtein_distance
14         * return the average distance between all sequence
15         */
16        for e1 in &self.elements {
17            for e2 in &a_cluster.elements {
18                result += e1.levenshtein_distance(&e2) as f32;
19                counter += 1 as f32;
20            }
21        }
22
23        result = result/counter;
24        Some(result)
25    }
26 }
```

With this cluster : ((ATTACG,ATATCG),(GCCGAG,(ACCCCG,TCCCCG)))
We have : (cl1, cl2) and cl2 = (GCCGAG, cl3)
So, we should have $cl3.linkage(cl2) < cl1.linkage(cl2)$. The average distance between cl1 and cl2 should be pretty high > 2 . Same for cl1 and cl3. But cl2 and cl3 should be quite close (being in the same cluster must bring us closer)

| linkage() | | | |
|-----------|-----|------|------|
| | cl1 | cl2 | cl3 |
| cl1 | 1 | 4 | 3.5 |
| cl2 | 4 | 1.55 | 1.33 |
| cl3 | 3.5 | 1.33 | 0.5 |

To create the `clusterize_agglomerative()` method we imagined some approach :
First approach :

- ▶ Compare each sequence.
- ▶ Group two close sequences in new cluster
- ▶ Group other sequences in another cluster
- ▶ Clusterize the second cluster

Second approach :

- ▶ Calculate the distance for each sequence two by two
- ▶ Compare this distance with ALL of the other sequence
- ▶ If the first distance is the smallest
 - ▶ Create the sub-cluster containing the two sequence
- ▶ Verify if the members are coherent
 - ▶ Distance between all of them is smaller than any distance of a another sequence (in the universe)
- ▶ If this cluster is good to go, create it for good.
- ▶ Compare two cluster

Third approach :

state 1. Create clusters for all sequence indiv (ok when initialization)

state 2. Compare all sub-clusters one by one by distance

state 2.1 Get the smallest distance

state 2.2 Check cluster condition if it's correct confirm this cluster else switch to the next

state 3. Back to state 2 until the total of sub-cluster ≤ 2

The problem is this method used a while loop that contains two nested for loop. The first for loop determine the smaller distance between all clusters. The second verify the coherence of the potential new cluster (with `i` from the while 'main_loop' and `keep` from the first for loop). Once the cluster is verified (if not we're incrementing `i` and get back to the 'main_loop'), we create the new cluster (`i` and `keep`), remove the sub-cluster `i` and `keep` from the general list of sub-clusters. The new cluster is inserted in place 0. And we reiterate the while loop ('main_loop') from 0, to compare the fresh cluster to others.

```

1 pub fn clusterize_agglomerative(&mut self)
2 { // start timer
3     let st = SystemTime::now();
4
5     let mut i = 0;
6     'main_loop:
7     while i < self.sub_clusters.len() && self.sub_clusters.len() > 2 {
8
9         let mut min: f32 = f32::MAX;
10        let mut keep = 0;
11        // state 2
12        for j in 0.. self.sub_clusters.len() {
13            if i != j {
14                // state 2.1
15                let dist = self.sub_clusters[i]
16                    .linkage(self.sub_clusters[j].clone())
17                    .unwrap();
18                if dist < min {
19                    min = dist;
20                    keep = j;
21                }
22            }
23            // we tried here to verify at the same time the cluster's coherence
24            // but we need to keep track of a smallest distance BEFORE checking
25            // the coherence (in the case if the closest is after j)
26
27        }
28        // state 2.2
29        for j in 0..self.sub_clusters.len() {
30            if j != keep && j != i {
31                if self.sub_clusters[keep]
32                    .linkage(self.sub_clusters[j].clone())
33                    .unwrap() < min {
34                    i+=1;
35                    continue 'main_loop;
36                    // the while exit is ONLY here
37                }
38            }
39        }
40
41        // create a new subCluster which contains two coherent element
42        let sub_1 = self.sub_clusters[i].clone();
43        let sub_2 = self.sub_clusters[keep].clone();
44
45        // commit: from here I remove all the ref & to self call
46
47        // remove the two clusters from the first cluster
48        self.sub_clusters.remove(i);
49        if keep > i {
50            self.sub_clusters.remove(keep-1);
51        }else {
52            self.sub_clusters.remove(keep);
53        }
54
55        // add the new cluster which contains the two last clusters
56        self.sub_clusters.insert(
57            0, ClusterOfSequence::new_with_clusters(sub_1, sub_2));
58
59        i = 0;
60    }
61    // check timer
62    let ed = SystemTime::now();
63    println!("{:?}", ed.duration_since(st).unwrap());
64 }

```

4 Application to hemoglobin

The data directory contains the nucleotide sequences of hemoglobin for several species in fasta format. They were downloaded from [Ensembl](#).

4.1 Read the gene sequence from a fasta file

After reforming all () fasta to be :

```
1 >some text that is skipped
2 ATGGTGTCTGTCTGCCGCCGACAAGGGCAATGTCAAGGCCGCTGGGGCAAGGTTGG
```

utils.rs

```
1 use std::{
2     fs,
3     time::SystemTime
4 };
5 use crate::sequence::*;
6
7 fn read_fasta(path: &str) -> String {
8     let start = SystemTime::now();
9     let content = fs::read_to_string(path).expect("Can't touch it");
10    let mut res:String = String::new();
11    //println!("{:?}", content);
12    for line in content.split("\n") {
13        if !line.starts_with(">") {
14            res.push_str(&line.replace("\r", ""));
15        }
16    }
17    let end = SystemTime::now();
18    println!("read content: {:?}", end.duration_since(start).unwrap());
19    res
20 }
```

4.2 Cluster of gene sequences

We created a method that read automatically all .fa file and return a Cluster containing all raw sequence.

```
1  /**
2   * it's meant to simplify the main method by just one method
3   * calling as many time the read_fasta method for all the
4   * fasta files in the path (in param).
5   * @param a specific folder name in the res folder
6   */
7  pub fn analyze_from_ressource_folder() -> ClusterOfSequence{
8      let paths = fs::read_dir("./data_reformed").unwrap();
9
10     let mut sequences: Vec<Sequence> = Vec::new();
11
12     for path in paths {
13
14         let mut name = path.as_ref()
15             .unwrap()
16             .path()
17             .display()
18             .to_string();
19
20         name = name.replace("./data_reformed\\", "")
21             .replace("_sequence.fa", "")
22             .replace("_", " ");
23
24         sequences.push(Sequence::new_with_string(
25             name.to_string(),
26             read_fasta(&path.unwrap().path().display().to_string())
27             ));
28         println!("name of the sequence: {}", name);
29     }
30
31     let bio_cluster = ClusterOfSequence::new_with_sequences(sequences);
32
33     bio_cluster
34
35
36
37 }
```

5 Tutorial

5.1 Install Rust

For your own good: [Install Rust](#)

Open a terminal and check the cmd : \$ rustup -V and \$ cargo -V

If it doesn't work GL (internet is your friend).

In terminal, go to the root of the Clustering project, and type \$ cargo run

Here we go the program runs !

5.2 Install Foam

Foam is a developing tools for thought. But we can use it for bioinfo ! Launch [Visual Studio Code](#), install the [Foam extension](#).

Open the Clustering project by >File >Open Folder...

Now, type *ctrl + p* and > *Foam : ShowGraph* The Graph will normally display.