

Master 1 Bioinformatique

Object-oriented programming

Hierarchical clustering

Emmanuelle Becker, Olivier Dameron, Marine Louarn

February 13, 2018

Version 1.4

1 Objective

This project's goal is to classify a set of students according to their grades, and to generate the corresponding dendrogram.

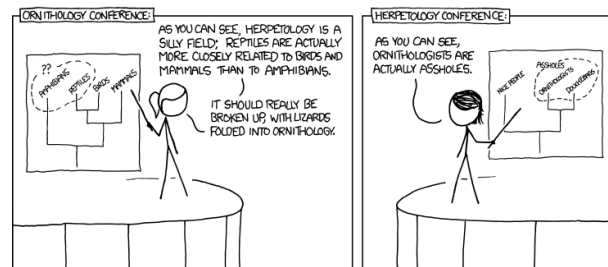


Figure 1: xkcd 867 (<http://xkcd.com/867/>)

2 Representing a class of students

2.1 Class Student: methods overload

Step 1 Create a class *Student* that represents the set of students. Each student has an (assumed) unique identifier (a string) and a grade (a double).

Step 2 Add a first constructor having for parameters an identifier and a grade. Then add a second constructor having an identifier as single parameter. This is an excellent opportunity to use method overload...

Step 3 Add the methods *getIdent()*, *getGrade()* et *setGrade(double newGrade)*.

Step 4 In the `main(...)` method, create the following instances and check that the methods from step 3 still work correctly:

```
1 Student riri = new Student("riri", 12.5);
2 Student fifi = new Student("fifi", 14.0);
3 Student loulou = new Student("loulou", 18.5);
4 Student geo = new Student("geo", 19.5);
5 Student donald = new Student("donald", 10.5);
```

2.2 Class Promo: inheritance and static methods

Step 5 Create a class `Promo` that represents a set of students. `Promo` is a sub-class of `java.util.ArrayList`¹. Make sure to read the documentation for `ArrayList`, you will need it soon. Please note that `ArrayList` is a generic class, whereas all the elements of a `Promo` instance are composed of instances of `Student`, so you will need to state that `Promo` is a subclass of `ArrayList<Student>`.

Step 6 In the `Promo`'s `main(...)` method, create an instance `m1bioinfo` of `Promo`, and add the members `riri`, `fifi`, `geo`, `donald` and `loulou` (respect this order so that the highest and lowest grades are in the middle of the list. The idea here is to avoid having the students almost sorted for the clustering).

Step 7 Add the methods `getMinGrade()`, `getMaxGrade()` et `getAverageGrade()` to the class `Promo`. The class `java.lang.Math`² has several useful methods. For iterating over all the students that compose a promotion, you can seek inspiration from the article "Traversing collections"³. Feel glad to have declared `Promo` as a subclass of `ArrayList`.

Step 8 In the `main(...)` method of `Promo`, add the code for printing the lowest grade, the highest and the average for `m1bioinfo`.

Step 9 The `PromoLoader` class (**FIXME: fournie sur l'ENT**) has a method `loadTsvFile(...)` that takes as argument a text file (one student per line; its identifier, a tabulation, its grade) for creating an instance of `Promo`. Why is the method `loadTsvFile(...)` declared as "static"? Draw the parallel with the methods `min(...)`, `max(...)` and `abs(...)` from `java.lang.Math`.

3 Hierarchical clustering

3.1 Principle

3.1.1 Agglomerative vs. divisive approaches

Classifying consists in organizing a set of elements in groups based on the elements' similarities or differences.

¹<http://docs.oracle.com/javase/7/docs/api/java/util/ArrayList.html>

²<http://docs.oracle.com/javase/7/docs/api/java/lang/Math.html>

³<http://docs.oracle.com/javase/tutorial/collections/interfaces/collection.html>

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

- the **agglomerative approach** (also called ascending) starts by creating one (atomic) cluster for each element, and then iteratively generates new clusters composed of the most similar two, until there only remains one cluster;
- the **divisive approach** (also called descending) starts by gathering all the elements into a single cluster, and then iteratively decompose the clusters into subclusters until each of them is only composed of a single element.

The divisive approach requires more operations than the agglomerative one and is therefore longer... except when we only need the most general clusters (e.g. to separate a sample into two groups).

3.1.2 Distance measures between elements and between clusters

For both the agglomerative and ther divisive approaches, clustering depends on two main parameters:

- a **distance measure between elements** (also simply called *distance*). There are several classical ones: euclidian distance, Manhattan distance... In our case, we will consider that the distance between two students is the absolute value of the difference of their grades;
- a **distance measure between clusters** (also called *linkage*) that relies on the *distance* between elements of the two clusters. There are several classical linkage measures: the average of the distances between all the combinations of elements, their maximum, their minimum... In our case, we will consider that the distance between two clusters of students is the average of the distances between all the elements of the first cluster and all the elements of the second cluster.

3.2 Classe ClusterOfStudents

This section aims at implementing the `ClusterOfStudents` class for representing a cluster of `Student` instances. A simple cluster is composed of a single instance of `Student`. A complex cluster is composed of several sub-clusters which can themselves be either simple or complex clusters. A complex cluster has a tree-like structure where all the leaves are simple clusters.

Initially, a complex cluster is only composed of simple clusters (Fig. 2). After clustering, a complex cluster is composed of sub-clusters that are intermediate complex clusters (Fig. 3)

3.2.1 Initialization

Step 10 Create a `ClusterOfStudents` with a `subclusters` attribute that represents the list of its sub-clusters. For simplifying the clustering step (when marshalling the students composing the cluster), add an attribute `students` that represents the set of students constituting the leaves of the cluster.

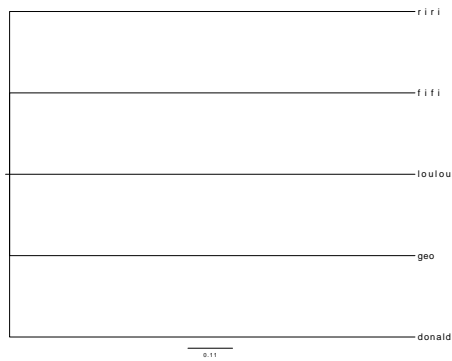


Figure 2: Complex cluster in its initial state: it is composed of five simple sub-clusters, each composed of a student.

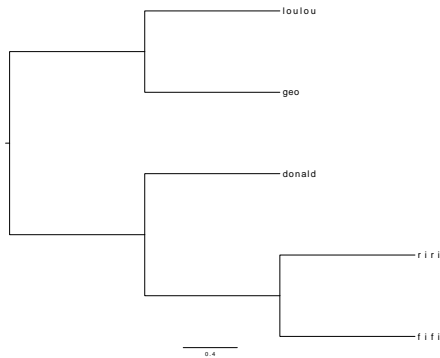


Figure 3: Complex cluster after clustering. It is composed of two intermediate complex sub-clusters. The first is itself composed of two simple clusters (geo and loulou). The second is composed of a simple sub-cluster (donald) and a complex sub-cluster composed of two simple sub-clusters (riri and fifi).

Step 11 Add the following constructors:

- a default constructor `ClusterOfStudents()` that creates an empty cluster (we do not really need it, but I find it cleaner to have a default constructor);
- a constructor for simple clusters `ClusterOfStudents(Student aStudent)`;
- a constructor for complex clusters before clustering `ClusterOfStudents(Promo aPromo)`.

Step 12 In `ClusterOfStudents`' `main(...)` main method, create an instance of a simple cluster `geoCluster` initialized with `geo`, and an instance of complex cluster `bioCluster` initialized with `m1bioinfo`.

3.2.2 Visualization

The Newick format⁴ provides a straightforward representation of trees and dendrograms, and is supported by most visualization tools. You can use the Tree Viewer web server⁵ or T-REX⁶ or the dedicated softwares FigTree⁷, dendroscope⁸ (free use in an academic context; getting a licence is not required for the basic functions). FigTree seems to give the best results.

The dendrogram from Fig. 2 can be represented by
`((loulou,geo),(donald,(riri,fifi)));`.

NB : for visualizing dendrograms, we could as well have used the R functions via the Java-R binding, but it is more complicated, and writting Newick files makes for an interesting exercise anyway.

Step 13 Add a `getNewick()` method to the class `ClusterOfStudents` that returns a string representing the dendrogram in the Newick format. Because of the final semicolon, you may need to introduce an intermediate function (aptly named `getNewickIntermediate()`). For marshalling the tree, you will make your life easier by considering a recursive approach (but this is not mandatory). Should these methods' visibility be public, protected or private?

Step 14 Generate a Newick representation of `bioCluster` and check (for example with T-REX or dendroscope) whether you get something similar to Fig. 2.

3.2.3 Clustering

Step 15 Add a method `linkage(ClusterOfStudents anotherCluster)` that returns the distance between the current cluster and `anotherCluster`. Choosing the average of the absolute value of the grade differences for each combination of students from each cluster is probably the easiest solution.

Step 16 In `ClusterOfStudents`'s `main(...)` method, create two simple clusters `loulouCluster` and `donaldCluster` and check whether the distance between `geoCluster`, `loulouCluster` and `donaldCluster` are what you expect them to be (check the six combinations).

Step 17 In `ClusterOfStudents`'s `main(...)` method, create the complex cluster `geoLoulouCluster` and check whether its distance with `donaldCluster` and `geoLoulouCluster` (and conversely).

Step 18 Add a method `clusterizeAgglomerative()`. Perform clustering on `bioCluster` and display the result as a Newick string.

Figure 4 shows the classification result for a set of students. Notice that because all the branches have the same length, the dendrogram seems to display two main clusters. Figure 5 shows that by making the length of each branch proportional to the distance separating the two clusters it joins, the dendrogram reveals three main clusters (cf. section 4.1).

⁴<http://evolution.genetics.washington.edu/phylip/newicktree.html>

⁵<http://www.proweb.org/treeviewer/>

⁶<http://www.trex.uqam.ca/>

⁷<http://tree.bio.ed.ac.uk/software/figtree/>

⁸<http://ab.inf.uni-tuebingen.de/software/dendroscope/>

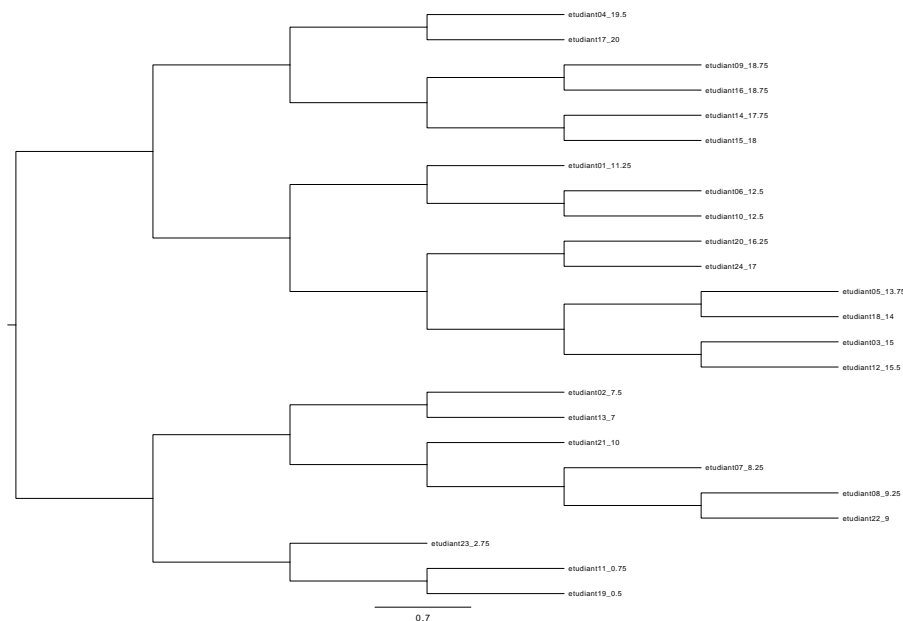


Figure 4: Complex cluster after clustering. The length of each branch is constant.

4 Extensions facultatives

4.1 Amélioration du dendrogramme

Le format Newick permet d'indiquer la longueur de chacun des branches. Utilisez plutôt l'application FigTree⁹ ou le site de Tree Viewer¹⁰ pour la visualisation, les autres outils semblent avoir des problèmes.

Step 19 Améliorez la méthode `getNewick()` pour que toutes les feuilles soient au même niveau.

Step 20 Améliorez la méthode `getNewick()` pour que toutes les feuilles soient au même niveau et que la longueur des branches soit proportionnelle à l'écart des notes.

4.2 Approche divisive

Step 21 Ajoutez une méthode `clusterizeDivisive()` à la classe `ClusterOfStudents`. Clusterisez `bigCluster` et comparez avec l'approche ascendante.

4.3 Modélisation de `ClusterOfStudents`

Dans la classe `ClusterOfStudents`, chaque instance d'étudiant apparaît deux fois :

⁹<http://tree.bio.ed.ac.uk/software/figtree/>

¹⁰<http://www.proweb.org/treewiewer/>

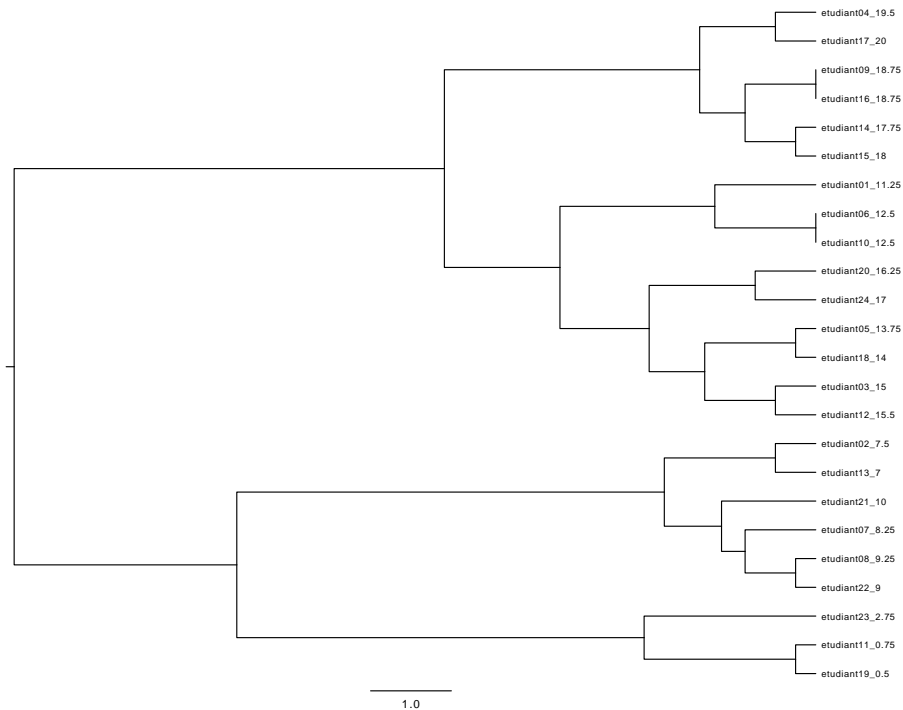


Figure 5: Complex cluster after clustering. The length of each branch is proportional to the distance between the two clusters it unites.

- dans l'attribut **subClusters** puisque le dendrogramme a autant de sous-clusters feuilles que d'étudiants ;
- dans l'attribut **students** qui permet un parcours plus simple de la liste des étudiants d'un cluster en évitant de devoir parcourir récursivement tout le dendrogramme à chaque fois.

On pourrait penser que cela occupe donc deux fois plus de place en mémoire que nécessaire, même si dans notre cas ce surcoût est acceptable dans la mesure où chaque instance occupe peu de place mémoire et qu'il y a relativement peu d'étudiants. Néanmoins, Java ne duplique évidemment pas les instances de **Student** dans les deux attributs. Chaque attribut ne contient en fait que les adresses des instances de **Student** (on appelle ça un *passage d'objet par référence*). Ainsi, si vous modifiez la note d'un étudiant dans **students**, cette modification apparaîtra également dans **subClusters**, et inversement.

Au final, le fait d'utiliser deux attributs qui semblent redondants parce qu'ils contiennent les mêmes objets :

- a l'avantage principal d'améliorer le temps de traitement en évitant un parcours de l'arborescence du dendrogramme chaque fois que l'on souhaite parcourir les étudiants (et cela arrive souvent durant l'étape de clustering) ;
- a l'avantage secondaire de simplifier l'écriture de la classe en vous dispensant justement d'écrire la fonction de parcours de l'arborescence du

dendrogramme ;

- a l'inconvénient d'augmenter légèrement la consommation de la mémoire.

Step 22 *Écrivez une classe `ClusterOfStudentsBis` qui ne contient que l'attribut `subClusters`. Comparez les temps de clusterisation de `ClusterOfStudents` et de `ClusterOfStudentsBis`.*