

You will implement algorithms for clustering vectors in the Euclidean space \mathbb{R}^K and of polygonal curves in Euclidean space \mathbb{R}^2 using the 8 combinations of variations that follow. The metric L_1 (Manhattan) will be used for the vectors and the pseudometric Dynamic Time Warping (DTW) for curves.

Initialization

1. Random selection of K points / K curves (simplest)
2. K-means++

Assignment

1. Lloyd's assignment (simplest approach)
2. Assignment by Range search with LSH for vectors / curves (inverse assignment)

Update

1. Partitioning Around Medoids (PAM) instead of Lloyd's
2. Update of Mean Vector / Update of DTW centroid Curve

ENTRANCE

1) A text file `input.dat` separated by spaces, colons or commas (e.g. `100 100 100` or `"100:100:100"`), which will have the following notation in the case of polygonal curves:

curves

curve `id` `mL` (`x11,y11`) (`x12,y12`) ... (`xlm1i,ylm1i`)
curve `id` `mN` (`xN1,yN1`) (`xN2,yN2`) ... (`xNmN,yNmN`)

where (`S11`, `c11`) the coordinates of the j -th point of the curve, where $j \leq p_i$ and p_i is the number of points of the curve i .

In `input.dat` the following notation:
vectors

item `id` `X11` `X12`

item `id` `XN1` `XN2`

where `Xi` the coordinates of the infinite vector that corresponds to item i

2) A configuration file of the following format (lines where there is a value may not be given so the default value is used):

`number_of_clusters:` `<int>` // K of K-means

`number_of_grids:` `<int>` //default: 2

`number_of_vector_hash_tables:` `<int>` //default: L=3

`number_of_vector_hash_functions:` `<int>` // k of LSH for vectors

Then, in `input.dat`, `cluster.conf` and `output.dat` will be done through the command:

```
./cluster -i <input file> -c <configuration file> -o <output file> -
```

complete <optional>

exit

A text file containing the clusters of data produced by each variant of the algorithm, the execution time in each case as well as the internal index

ag  ioAdynons tn  ovotadoroinons Silhouette.

In the case of vectors, the output file must follow the following pattern, the which is repeated for each variant:

Algorithm: IxAxUx

CLUSTER-1 {size: <int>, centroid: <item _id> or array pe TL  suvtTetayp veg tou centroid

in the case of k-means Update}

CLUSTER-K {size: <int>, centroid: <item _id> or array pe TL  suvtTetayp veg tou centroid

otnv meptutwon K-means Update }

clustering time: <double> //in seconds

Silhouette: [s1,...,si,...,sK, stotal]

/* si=average s(p) of points in cluster i, stotal=average s(p) of points in dataset */

/* Optionally with command line parameter -complete */

CLUSTER-1 {item idA, item idB, ..., item idC}

CLUSTER-K {item idR, item idT, ..., item idZ}

In the case of polygonal curves, the output file must follow the following template, which is repeated for each variation:

Algorithm: IxAXUx

CLUSTER-1 {size: <int>, centroid: <curve_id> or array of onpefa tou centroid in meptutwon DTW medoid curve update}

CLUSTER-K {size: <int>, centroid: <curve id>}

clustering time: <double> //in seconds

Silhouette: [s1,...,si,...,sK, stotal]

/* si=average s(c) of curves in cluster i, stotal=average s(c) of curves in dataset */

/* Optionally with command line parameter -complete */

CLUSTER-1 {curve idA, curve idB, ..., curve 140}

CLUSTER-K {curve idR, curve idT, ..., curve 146}

Additional requirements

1. File (or section in VeDaPITH) that compares algorithms based on results.
2. The program must be well organized with separation of declarations / definitions functions, structures, and data types into logical groups that correspond to separate header and source code files. The program must be compiled with the use of the tool piãKke and the existence of a suitable MDKePI6. Its quality is also rated code (e.g. avoiding PIEPIOI |6ã(s)).
3. The deliverable must be adequately documented with full code commentary and the existence of file (6b(η)e which includes at least: a) title and description of the program, b) list of code / header files and their description, c) its compilation instructions of the program, 6) instructions for using the program and e) full details of the students who use it they developed.
4. The implementation of the program should be done using a version management system software and collaboration ((ii or 5NN) [groups of 2 people].
5. Using an appropriate library and performing unit testing.