

Web 2.0 and Mobile Interaction - Assignment 1

Olivier Jais-Nielsen

October 7, 2009

1 Approach

I decided to implement different kind of clustering on movies information.

More precisely, my program can download information about the current top rated movies from IMDb, using the API IMDbPy and perform clustering on these films based on the data provided by IMDb : the international titles, several plot summaries and the genres attributed to the films. One can easily add any other data provided by IMDb. I choosed the previous because I think any other data (like actors, directors, languages, ...) doesn't describe the films themselves but their environment and thus wouldn't be relevant. The selection of the words to take into account was inspired by the script given in lecture 3 : only the words present in at least 10% and at most 50% of the movies are kept.

Both the hierarchical and the k-means clusterings are implemented. The Tanimoto and the Pearson similarity measures can be used as well as the inverse Euclidian distance, defined by : $IE(v_1, v_2) = \exp^{-\|v_1 - v_2\|}$

The program produces a box diagram representation of the clusters, including for each movie some information (title, year, rating) and a thumbnail.

2 Python code

2.1 API

In addition to IMDbPy, several classes and functions have been defined to handle the data and the algorithms involved in clustering. A complete documentation of all the modules is given in the document *modules.pdf*. The most important classes are *WordVector* which defines a vector of word occurrences, *Tree* which defines a tree used for the hierarchical clustering and *ClustersSet* which defines a set of clusters for k-mean clustering.

All the source files are in the *Sources* folder.

2.2 Main algorithms

The main algorithms are contained in the module *clustering.py*.

2.2.1 Function *hierarchicalCluster(vectors, simMeasure)*

```
1 """
2 This function computes the hierarchical clustering from the WordVector instances in the "
  vectors" list using the similarity measure "simMeasure". It returns a Tree with
  WordVector labels.
```

```

3  """
4
5  # We create an empty Tree with the appropriate label merger. It will act as the top tree.
6  tree = Tree(mergeWordVectorsLabels)
7
8  vectPairSimilarities = {}
9
10 # For each vector, we create a leaf with the appropriate label and we add it as a branch of
    the top Tree. These are the initial clusters.
11 for vectID in vectors.keys():
12     branch = Tree();
13     branch.setLeaf((vectID, vectors[vectID]))
14     tree.addBranch(branch)
15
16 # While there is more than one cluster at the top level, we iterate.
17 while tree.countBranches() > 1:
18
19     # We find the highest similarity among the cluster pairs (we store the similarities
        in vectPairSimilarities in order to avoid calculating several times the same
        thing)
20     maxSimilarity = -1;
21     for i in range(tree.countBranches()):
22         vectIID, vectI = tree.getBranchLabel(i)
23         for j in range(i + 1, tree.countBranches()):
24             vectJID, vectJ = tree.getBranchLabel(j)
25             sim = None
26             vectIID = str(vectIID)
27             vectJID = str(vectJID)
28             if (vectIID, vectJID) in vectPairSimilarities.keys():
29                 sim = vectPairSimilarities[(vectIID, vectJID)]
30             elif (vectJID, vectIID) in vectPairSimilarities.keys():
31                 sim = vectPairSimilarities[(vectJID, vectIID)]
32             else:
33                 sim = simMeasure(vectJ, vectI)
34                 vectPairSimilarities[(vectIID, vectJID)] = sim
35                 vectPairSimilarities[(vectJID, vectIID)] = sim
36             if sim > maxSimilarity:
37                 maxSimilarity = sim
38                 bestI = i
39                 bestJ = j
40                 bestVectIID = vectIID
41                 bestVectJID = vectJID
42
43     # We merge the clusters of the corresponding pair
44     tree.mergeBranches(bestI, bestJ)
45     del vectPairSimilarities[(bestVectIID, bestVectJID)]

```

```

46         del vectPairSimilarities[(bestVectJID, bestVectIID)]
47
48 # Once there is only one cluster at the top level, we return it
49 return tree

```

2.2.2 Function *kMeansClustering*(vectors, simMeasure, k)

```

1  """
2  This function computes the k-mean clustering from the WordVector instances in the "vectors"
   list using the similarity measure "simMeasure" and "k". It returns a ClusterSet
   containing WordVector instances.
3  """
4
5  # We create a ClusterSet with no clusters and with the set of vectors      clusters =
   ClusterSet(vectors)
6  # We create k clusters and fill them arbitrarily : the i-th vector belongs to the (i mod[k])
   -th cluster.
7  i = 0
8  for vectID in vectors.keys():
9      clusters.attributeVector(vectID, i % k)
10     i += 1
11
12 # While the composition of the cluster differs from the previous iteration, we iterate.
13 while clusters.hasChanged():
14     # We compute the centroid of each cluster for their current composition.
       clusters.computeCentroids()
15
16 # For each vector :
17     for vectID in vectors.keys():
18
19         # We find the highest similarity between the vector and a centroid
20         maxSimilarity = -1
21         for clusterID in range(k):
22             currentSim = simMeasure(vectors[vectID], clusters.getCenter(
               clusterID))
23             if currentSim >= maxSimilarity:
24                 maxSimilarity = currentSim
25                 bestClusterID = clusterID
26
27         # We assign the vector to the cluster corresponding to the found centroid
               clusters.attributeVector(vectID, bestClusterID)
28
29 # Once the composition of the clusters stops changing, we return the ClusterSet
30 return clusters

```

2.2.3 Example application

```

1 # We create an interface with IMDb
2 interface = IMDbInterface("data.db")
3
4 # We get the IMDb movie IDs corresponding to the Top 250
5 movieIDs = getIMDbTop250()
6
7 # We prepare dictionaries intended to contain, for each movie ID, an IMDb movie object and a
  WordVector
8 movies = {}
9 wordVectors = {}
10
11 # For each movie ID
12 for i in range(0, 50):
13     id = movieIDs[i]
14
15     # We get the IMDb movie object
16     movies[id] = interface.getMovie(id)
17
18     # We select the texts concerning the movie that we want to use
19     list = [movies[id]["canonical_title"] + movies[id]["akas"] + movies[id]["genres"] +
20             movies[id]["plot"]]
21
22     # We remove the annotations from these texts : they can be in a text, after a double
23     colon      list = map (lambda text: text.split("::")[0], list)
24
25     # We compute a word histogram from the texts and a WordVector from it
26     wordVectors[id] = WordVector(wordHistogramFromTextList(list))
27
28 ## HIERARCHICAL CLUSTERING
29
30 # We create a PIL image object intended to contain the box diagram of the hierarchical
31 clustering
32 im = Image.new("RGB", (0, 0), (255, 255, 255))
33
34 # We compute the clusters for the choosen similarity measure
35 clusters = hierarchicalCluster(wordVectors, lambda vect1, vect2: similarity(vect1, vect2, "
36 euclidian_inverse"))
37
38 # We define an appropriate full label writer : it displays the canonical title, the year
39 and the rating of the corresponding movie
40 clusters.setFullLabelWriter(lambda (vectID, vect): (u"%s_(Year_:_%i_Rating_:_%i)" % (movies[
41 vectID]["canonical_title"], movies[vectID]["year"], movies[vectID]["rating"])).encode("
42 utf-8"))

```

```

39 # We define an appropriate thumbnail accessor : it displays the corresponding movie's
    thumbnail    clusters.setThumbnailAccessor(lambda (vectID, vect):
40 interface.getThumbnail(vectID))
41
42 # We resize the image to the appropriate size , we draw the diagram and we save the image
43 im = im.resize(clusters.boxDimensions())
44 clusters.drawDiagram(im)
45 im.save("hierarchical_eucl.jpg")
46
47
48 ## K-MEANS CLUSTERING
49
50
51 # We create a PIL image object intended to contain the box diagram of the hierarchical
    clustering
52 im = Image.new("RGB", (0, 0), (255, 255, 255))
53
54 # We compute the clusters for the choosen similarity measure and the choosen k (here 6)
55 clusters = kMeansClustering(wordVectors, lambda vect1, vect2: similarity(vect1, vect2, "
    euclidian_inverse"), 6)
56
57 # We define an appropriate full label writer : it displays the canonical title , the year
    and the rating of the corresponding movie
58 clusters.setFullLabelWriter(lambda vectID: (u"%s_(Year_:_%i_Rating_:_%i)" % (movies[vectID][
    "canonical_title"], movies[vectID]["year"], movies[vectID]["rating"])).encode("utf-8"))
59
60 # We define an appropriate thumbnail accessor : it displays the corresponding movie's
    thumbnail
61 clusters.setThumbnailAccessor(interface.getThumbnail)
62
63 # We resize the image to the appropriate size , we draw the diagram and we save the image
    im = im.resize(clusters.boxDimensions())
64 clusters.drawDiagram(im)
65 im.save("k-means_eucl.jpg")

```

3 Results

Both kinds of clustering and for all the three similarity measures have been tested on the 50 top rated movies. The results are in the files named according to the format *kind_of_clustering-similarity_measure.jpg* in the *Results* folder.

After performing a hierarchical clustering, we find the major subdivisions that appear and it can give an idea on how many clusters to use with the k-mean clustering. Here, we chose 6 clusters.

Although it is obvious the hierarchical clustering with the inverse Euclidian similarity is a failure, we can also see this similarity measure isn't satisfactory for the k-means clustering : the clusters sizes have no coherence in comparison to the other results. This is explainable as the inverse Euclidian measure quickly reaches extremely low values, when two vectors differ, thus the computational precision is low.

For the other results, it is difficult to evaluate the quality as it is very subjective. One criterium could be to check whether films episodes are in the same clusters. For that criterium, the hierarchical clustering with a Pearson similarity is very good while the k-means clustering

with the same similarity is not that good.

4 Conclusion

These clustering method are very flexible and therefore can give very accurate results in specific cases. However, in general cases such as this one, it is not easy to appriciate the quality of the results. Moreover, when one criterium is satisfied, like for instance in the hierarchical clustering with a Pearson similarity, there can still be strange aspects : like associating “Amélie Poulain” with “Psycho”...