

Rival-Penalized Competitive Clustering: A Study and Comparison

Alberto Borghese and Wiliam Capraro

Applied Intelligent Systems Laboratory, Dept. of Computer Science, University of Milano
borghese@di.unimi.it, wiliam.capraro@studenti.unimi.it

Abstract. A major recurring problem in exploratory phases of data mining is the task of finding the number of clusters in a dataset. In this paper we illustrate a variant of the competitive clustering method which introduces a rival penalization mechanism, and show how it can be used to solve such problem. Additionally, we present some tests aimed at comparing the performance of our rival-penalized technique with other classical procedures.

1 Introduction

The term central clustering refers to a family of clustering algorithms that are based on moving a set of points, referred to as prototypes, inside the data space until their position minimizes a certain cost function, representing a measure of the goodness by which the prototypes represent the data points.

In the literature, approaches based on soft-clustering, like fuzzy c-means [14], Neural-Gas [13,4] or Self-Organizing Maps (SOM) [6,16], and competitive learning [5,15,12], have been proposed to partition a given dataset into a predefined number of clusters, each one represented by a prototype usually corresponding to the cluster centroid. All these algorithms suffer from several issues, most notably, the optimal value of the cost function is rarely reached. Only stochastic optimization [10], that is extremely costly, or a careful initialization of the prototypes allow escaping local minima. Although a few attempts have been proposed to derive a robust initialization (e.g. [11]), there seem to be no universal and reliable way to proceed, and some prototypes typically get stuck during the clustering process. These prototypes are referred to as “dead units” [13] and affect the proper operation of the algorithm and the quality of the result. As a consequence, the general approach is to repeat the clusterization process several times with different, random initializations of the prototypes, so as to allow the algorithm to escape from local minima from time to time.

A slightly different task is to find the number of clusters in a dataset (e.g. [3]). This is indeed a frequent problem in exploratory phases of data mining, and a straightforward approach is to adopt a parameterized version of a clustering algorithm using the desired number of data clusters K as parameter and to try a different clusterization for each possible value of the parameter. Subsequently, the best result would be chosen based on some validity measure or index.

However, it is impractical to run several random initializations of one algorithm for each possible values of its parameter, especially when the latter spans a wide interval of

values. Alternatively, in some cases it is possible to exploit the intrinsic characteristics of some algorithms to produce dead units (e.g. [6]), to facilitate the search for a good solution.

Recently, in the framework of central clustering, a different approach has been proposed for moving the prototypes. The idea is that, while the winning prototype is attracted by the closest data point, other prototypes are moved in the opposite direction. This mechanism, known as rival-penalization [5,15], is somehow similar to the BCM model proposed by Bienenstock, Cooper and Munro [9] in a typical Hebbian learning fashion.

Rival-penalization clustering has been overlooked in the past. In this paper we present the results of some tests we performed, aimed at comparing the rival-penalization approach with classical clustering techniques, namely standard competitive learning and SOM. The ability of rival-penalization of discovering the proper number of clusters in a given dataset is analysed and discussed. Specifically, we show that, by introducing the rival penalization mechanism into a competitive learning setting, results comparable with soft-clustering can be achieved. Moreover, the number of clusters can be discovered in a robust and reliable way.

2 Algorithms

We'll consider a collection of N d -dimensional observations, $\{\xi_j\}$. Goal of clustering algorithms is to assign each observation to one of K clusters Ψ_i , according to a similarity measure with the other elements in the same cluster. Each cluster is represented by its centroid, ψ_i , which is also a point in \mathbb{R}^d .

The following subsections give a quick coverage of the algorithms we employed.

2.1 Competitive Learning and Rival Penalization

Competitive learning (CL) is an effective tool for data clustering, widely applied in a variety of signal processing problems such as data compression, classification, adaptive noise cancelation, image retrieval and image processing [2].

For the purposes of this contribution, a feed-forward neural network with a single layer consisting of K output units is used to achieve a K -cluster data partitioning. Each unit represents a cluster centroid ψ_i .

The training of the network proceeds as follows. At each iteration, each data point ξ is presented in turn to the network and a winning unit, w , is elected. This is the prototype whose Euclidean distance from the point is minimum:

$$w = \underset{i}{\operatorname{arg\,min}} \|\xi - \psi_i\|. \quad (1)$$

Subsequently, the position of the winning unit is updated towards the data point using the following updating rule

$$\psi_{wj} = \psi_{wj} + \eta_{(t)}(\xi - \psi_w) \quad (2)$$

where j denotes a component of the prototype vector and $\eta_{(t)}$ is a learning rate parameter whose value decays as a function of the time t .

In a pure competitive learning setting, only the winning unit is updated. The procedure is repeated multiple times for each data point, until the prototypes converge to their final position—i.e. when the maximum difference in the position of any centroid in two successive iterations is smaller than a fixed tolerance ε , or when a maximum number of iterations is reached.

The prototypes are initialized using the "Forgy" approach [1]—i.e. K of the available data points are randomly chosen to serve as cluster prototypes. In this context, this is enough to guarantee that no dead unit will ever appear, as every prototype shall win the competition for at least one data point, that is, the prototype itself.

The rival-penalized competitive learning (RPCL) algorithm improves on the pure competitive learning approach by introducing a rival penalization mechanism, as proposed in [5] and [15]. With this approach, not only the position of the winning unit is updated towards the input vector, but additionally the position of its rival unit is updated in the opposite direction.

In order to find the winning unit and its rival, a relative winning frequency is introduced, which keeps track of how many times each unit happens to win a competition for some input vector. The relative winning frequency for unit i is defined as

$$\gamma_i = \frac{s_i}{\sum_{j=1}^K s_j} \quad (3)$$

where s_i is the number of times unit i was declared winner in the past. When $\sum_{j=1}^K s_j = 0$ —i.e. initially, then $\gamma_i = 1$ in order to give every prototype a fair chance to win.

The winning unit w for an input vector ξ is now given by

$$w = \arg \min_i \gamma_i \|\xi - \psi_i\|. \quad (4)$$

Notice how the parameter γ_i acts as a "conscience" for the unit—if the unit has won too often in the past, its chances to win the competition for the current data point are reduced accordingly. Moreover, for each input vector ξ , the rival penalized competitive learning algorithm computes not only the winning unit w , but also a second winning unit, referred to as the rival, defined by

$$r = \arg \min_i \gamma_i \|\xi - \psi_i\|, \quad i \neq w. \quad (5)$$

Equation 2 is used to update both the winner and its rival. The latter, however, moves away its centroid from the input point with a de-learning rate β , which is related to η by

$$\beta_{(t)} = -c\eta_{(t)}\gamma_r \quad (6)$$

where γ_r is the relative winning frequency of the rival and $c = 1/10$ is a predefined constant. Unlike the implementation of [5], here β depends on both the learning rate η and the winning frequency γ_r , so that the rival is dynamically penalized according to γ_r even for constant η (which is not the case anyway).

In contrast to the CL algorithm, here a "Forgy" initialization of the prototypes is not enough to guarantee dead unit avoidance. In fact, even if the prototypes are initialized using the input data points, depending on the de-learning rate β , a rival unit may incur

considerable modification in the value of its prototype, and thus it can fail to win the competition even for the input data point to which it had been initialized.

What is interesting with this approach is that, as reported in [5], if the learning rate η is chosen to be at least one order of magnitude larger than β , then the adequate number of output clusters will be automatically found. In other words, assuming that the actual number of clusters is unknown and that the number of units K is chosen greater than the cluster number, the prototype vectors will converge towards the centroids of the actual clusters with few of them overlapping in space. In our implementation, this condition holds in each iteration as $c = 1/10$. In each iteration, the RPCL algorithm pushes away the rival, thus allowing for faster convergence, and invalidates extra prototypes by eventually making their cluster empty. Hence, the RPCL algorithm is believed to be able to perform appropriate clustering without knowing the cluster number.

2.2 SOM

The limitation of considering only one data point at a time in competitive learning has been overcome by soft-clustering approaches [2] in which the position of all the prototypes is updated for each data point. Among these approaches, Self-Organizing Maps (SOMs) represent an excellent tool in exploratory phases of data mining. They project the input space onto prototypes in a low-dimensional regular grid that can be effectively used to visualize and explore properties of the data. The SOM consists of a regular, one- or two-dimensional grid of units, with each unit i represented by its prototype vector ψ_i . Additionally, each unit i is assigned a place in the output grid, represented by its coordinates $r_i = (x_i, y_i)$, and the units are logically linked to adjacent ones by a neighborhood relation. During training, data points lying near each other in the input space are mapped to nearby units in the output hyperplane. Thus, the SOM can be regarded as a topology-preserving tool for mapping the input space onto the output grid.

The SOM is trained iteratively. At each training step, a data point ξ is randomly chosen from the input data set, and the distance between ξ and all the prototype vectors is computed. Subsequently, all prototype vectors are updated, each proportionally to the distance of the corresponding unit from the winning unit in the output grid:

$$\psi_{ij} = \psi_{ij} + \eta_{(t)} \Lambda(i, w) (\xi - \psi_i). \quad (7)$$

In the hereabove equation $\Lambda(i, w)$ denotes the value of the neighborhood function between unit i and the winning unit w , as given by

$$\Lambda(i, w) = \exp\left(-\frac{\|r_i - r_w\|^2}{2\sigma^2}\right) \quad (8)$$

where the parameter σ defines the radius of the neighborhood. $\Lambda(i, w)$ therefore defines a region of influence for the prototype w . Notice that the value of Λ is exactly 1 when $i = w$, and decreases as the distance of the prototype from the other data points increases. Also, it is useful to adjust the radius as well as the learning rate at each iteration, so that the influence region of a prototype decays with time as a function of σ and η .

In this work, we are mainly concerned with the SOM’s ability to perform appropriate clustering of a given data set. Thus, only SOMs with one-dimensional output arrays are actually used. As stated in [16], this configuration is expected to produce better results as compared to the 2-dimensional grid configuration. This is due to the fact that the “tension” exerted in each unit by the neighboring units is much higher in the second configuration, and such a tension limits the plasticity of the SOM to adapt to the particular distributions of the dataset.

3 Experimental Setting and Test Results

In order to test the algorithms, we have generated a specific dataset containing 250 3- d data points distributed over 5 non-overlapping clusters. It has been generated by perturbing the centroid of each cluster with a Gaussian distribution with mean value 0 and variance 1.

Since the resulting clusterization depends strongly on the initialization of prototypes, it is essential that each algorithm be tested several times with different initializations. For our tests, 60 “Forgy” initializations (which we’ll refer to as trials) have been generated and evaluated for each algorithm. This should be enough to overcome random fluctuations.

As to the tests we conducted, they can be divided into two types. In a first type thereof—we call it type-A experiment—we focused on a specific algorithm and tried to partition our dataset varying the cluster number from $K = 2$ to $K = 10$. (And for each K , 60 trials have been performed as described before.) On the other hand, in type-B experiments we tested 60 trials of an algorithm with a fixed value of K but varying the parameters of the algorithm instead.

In our tests, the validity of the resulting clusterization is evaluated by means of quality indexes, and the number of iterations required by the algorithm to converge was also measured. The quality indexes used are the Davies-Bouldin (DB) index and the mean quadratic error (MQE). The mean quadratic error is simply the ratio of the sum of all the squared distances of each data point from its cluster prototype to the total number of data points:

$$MQE = \frac{\sum_{i=1}^K \sum_{\xi \in \Psi_i} \|\xi - \psi_i\|^2}{\sum_{i=1}^K |\Psi_i|}. \quad (9)$$

The Davies-Bouldin index is defined as

$$DB = \frac{1}{K} \sum_{i=1}^K \max_{j \neq i} \left\{ \frac{S_i + S_j}{\|\psi_i - \psi_j\|} \right\} \quad (10)$$

where S_i is the within i -th cluster scatter, as given by

$$S_i = \sqrt{\sum_{\xi \in \Psi_i} \frac{\|\xi - \psi_i\|^2}{|\Psi_i|}}. \quad (11)$$

For a detailed review of these and other cluster validity measures see [8]. Notice that it is geometrically plausible to seek clusters that have minimum within-cluster scatter

and maximum between-class separation, so the number of clusters \bar{K} that minimizes the Davies-Bouldin index can be reasonably taken as the optimal value of K . As reported in [8], for well-separated clusters, the Davies-Bouldin index is expected to decrease monotonically as K increases until the correct number of clusters is achieved.

In all the tests conducted, we fixed a tolerance of $\varepsilon = 0.001$ and the maximum number of iterations was set to 500. It should be enough for the algorithms to produce good clusterizations given the time-decay rule for η adopted, which is

$$\eta_{(t)} = \exp\left(-\frac{t}{50}\right) \cdot \eta_0 \quad (12)$$

where $\eta_0 = 0.1$ is the initial value and $t = 0, 1, \dots$ is the iteration number.

In the remaining of this section we illustrate the results of our tests.

3.1 Competitive Learning

To begin with, we measured the effectiveness of the standard CL algorithm in partitioning our dataset by running a type-A test. The results can be used throughout the rest of this work as a reference for the other algorithms. For each value of K , Table 1 reports the Davies-Bouldin index and the quadratic error for the best outcome out of the 60 trials of the algorithm, along with the number of iterations performed.

As Table 1 shows, the algorithm succeeds in discovering the correct number of clusters: the DB index takes on its optimal value for $K = 5$.

As expected, the mean quadratic error is a decreasing function of the number of clusters (indeed one expects the within-cluster variance to decrease in this case), and hence it does not convey any useful information on the goodness of the result.

As a downside, the CL algorithm takes a considerable number of iterations to converge, as in each iteration only the winning unit is moved towards the current data point by a small, η -dependent, fraction of the distance. Moreover, in order to discover the optimal value of the parameter K , every possible value has to be investigated and the result evaluated. The average number of iterations is a decreasing function of K , which is rather obvious since, for small K , we expect the amount of modification in the position of each centroid as a function of the data points to be higher in each iteration as compared to when K is large.

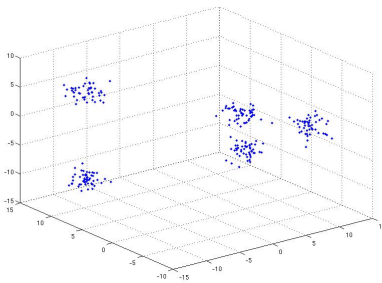


Fig. 1. Scatterplot of the dataset

Table 1. Type-A test results for CL

K	DB	MQE	it
2	0.712	49.642	344
3	0.396	27.360	344
4	0.429	10.583	328
5	0.298	2.816	298
6	0.690	2.648	298
7	0.751	2.603	298
8	0.768	2.540	298
9	0.746	2.435	298
10	0.809	10.197	328

Table 2. Type-A test results for SOM

(a) $\sigma = 0.5$					(b) $\sigma = 1$					(c) $\sigma = 1.5$				
K	nD	DB	MQE	it	K	nD	DB	MQE	it	K	nD	DB	MQE	it
2	0	0.7123	49.6421	346	2	0	0.7123	49.6421	346	2	0	0.7123	49.6421	346
3	0	0.5779	32.8641	345	3	0	0.5779	32.8641	345	3	0	0.5779	32.8641	345
4	0	0.4291	10.5827	330	4	0	0.4291	10.5827	330	4	0	0.4291	10.5827	330
5	0	0.2983	2.8161	300	5	0	0.2983	2.8161	300	5	0	0.2983	2.8161	300
6	1	0.2983	2.8161	300	6	1	0.2983	2.8161	300	6	1	0.2983	2.8161	300
7	2	0.2983	2.8161	300	7	2	0.2983	2.8161	300	7	2	0.2983	2.8161	300
8	3	0.2983	2.8161	300	8	3	0.2983	2.8161	300	8	3	0.2983	2.8161	300
9	4	0.2983	2.8161	300	9	3	0.7244	2.6374	294	9	2	1.0505	2.5142	294
10	5	0.2983	2.8161	300	10	3	1.0228	2.4724	294	10	3	1.0710	2.5077	294

We have also investigated the role of the learning rate η in the learning process. To this end, a type-B test has been performed in which a value of $K = 5$ has been fixed and η takes on some values in the range (0.2-0.02). As before, 60 trials of the algorithm have been tested for each value of η , and the best outcome is considered. We do not report the results for space issues. We report, however, that the initial learning rate seems to play no crucial role in the learning process, since for every value of η the best value obtained for the DB index is the same as that of Table 1.

3.2 SOM

As our second test, we have investigated the performance of a SOM in achieving proper clusterizations of the dataset. As before we ran a type-A test using a 1-dimensional output array of units, as we are not interested in the spatial organization of the resulting cluster centroids. The distance of each prototype from its neighbor prototypes on the output array has been set arbitrarily to 1, which is also the initial value for the radius of the neighborhood σ . The general idea is that, by exploiting the SOM's inherent ability to produce dead units, it is possible to avoid testing every possible value of the parameter K provided it is chosen larger than the actual number of clusters. Results are reported in Table 2b, where nD represents the number of dead units and again each line represents the best outcome for all the 60 initializations, according to the Davies-Bouldin index.

As the table implies, the minimum of the DB index is obtained for values of K in the range between $K = 5$ and $K = 8$. The values reported confirm the ability of the SOM to produce good clusterizations of the dataset, with values comparable with that of the competitive learning approach for all values of the parameter K . The results also advocate the thesis that the SOM is able to invalidate extra clusters and discover the correct number of clusters if the parameter K is chosen in a neighborhood of its optimal value.

We did not expect the quality of the resulting clusterization to change considerably as a function of the initial learning rate η , so we did not conduct any test in this respect. It is interesting, however, to observe the behavior of the algorithm when the initial radius is enlarged or restricted. Tables 2a and 2c also show the result of type-A tests

Table 3. Type-B test results for RPCL

(a) $\eta = 0.1$					(b) $\eta = 0.3$					(c) $\eta = 0.5$				
#	K	DB	MQE	it	#	K	DB	MQE	it	#	K	DB	MQE	it
1	5	0.2973	2.8166	298	1	5	0.3312	3.6632	366	1	6	0.9961	1743.7830	499
2	8	14.1375	436.7141	499	2	5	0.2955	2.8196	353	2	7	1.1205	529.9871	499
3	5	0.2983	2.8161	298	3	5	0.2978	2.8162	353	3	1	n.a.	n.a.	499
4	7	0.8392	3.4301	337	4	6	1.0019	594.3351	499	4	9	1.0335	8.1133	440
5	5	0.2977	2.8162	298	5	10	3.2478	12092.0839	499	5	8	2.0086	1834.3202	499
6	5	0.2972	2.8165	298	6	9	2.2782	27149.0575	499	6	6	0.6528	2.9607	411
7	5	0.2977	2.8162	298	7	7	2.2396	407.3565	499	7	7	1.7469	27172.1130	499
8	5	0.2972	2.8168	298	8	1	n.a.	n.a.	499	8	1	n.a.	n.a.	499
9	8	1.2113	127.1700	485	9	9	2.0219	14926.2644	499	9	1	n.a.	n.a.	499
10	5	0.2980	2.8161	298	10	8	1.0686	111.9457	499	10	1	n.a.	n.a.	499
11	9	1.9186	283.6699	486	11	1	n.a.	n.a.	494	11	3	1.3013	84.3499	454
12	5	0.2981	2.8161	298	12	1	n.a.	n.a.	496	12	7	1.0529	659.0174	499
13	9	1.5104	2.2541	285	13	8	3.3781	736.2310	499	13	8	1.5026	353.7500	499
14	10	1.4525	32.4955	438	14	6	0.6007	2.6212	353	14	5	0.2983	2.8161	379
15	8	1.8814	119.3174	478	15	9	1.1331	2.4473	363	15	7	1.3130	663.7360	499
16	9	3.2711	1289.8495	499	16	8	1.1660	254.0439	499	16	1	n.a.	n.a.	499
17	5	0.2981	2.8161	298	17	9	4.7230	68607.6355	499	17	4	0.9470	45.7815	432
18	10	1.3259	2.0547	283	18	1	n.a.	n.a.	499	18	9	8.7524	200.3968	499
19	5	0.2964	2.8175	298	19	1	n.a.	n.a.	486	19	8	1.3188	35309.9092	499
20	7	0.8711	2.4896	298	20	1	n.a.	n.a.	459	20	8	2.2503	712.0748	499
21	5	0.2961	2.8180	298	21	5	0.2948	2.8272	353	21	5	0.2983	2.8161	379
22	8	1.2890	187.3231	478	22	1	n.a.	n.a.	499	22	1	n.a.	n.a.	499
23	5	0.2983	2.8161	298	23	8	1.4005	22.1895	463	23	1	n.a.	n.a.	499
24	10	1.4116	2.0938	287	24	7	0.8376	2.4683	345	24	1	n.a.	n.a.	499
25	8	1.7381	87.3881	429	25	9	1.1667	2.7230	374	25	7	0.7805	3.0149	401
26	5	0.2959	2.8205	298	26	1	n.a.	n.a.	495	26	6	1.0100	270.6201	499
27	8	1.4844	786.6143	499	27	9	1.1134	6.0580	410	27	1	n.a.	n.a.	454
28	9	2.2929	931.2280	499	28	8	1.1568	355.1591	499	28	7	2.1072	12438.7022	499
29	5	0.2983	2.8161	298	29	6	1.4096	91.4646	443	29	1	n.a.	n.a.	499
30	5	0.2967	2.8174	298	30	7	1.1727	7076.8766	499	30	7	1.0933	3285294.0383	499
31	5	0.2970	2.8167	298	31	9	1.2331	22.9518	479	31	8	0.9186	3.4217	410
32	5	0.2975	2.8164	298	32	9	1.4621	599.8796	499	32	9	2.0005	4802.1818	499
33	5	0.2983	2.8161	298	33	9	1.2037	2.8699	384	33	7	0.9568	2.8659	402
34	5	0.2959	2.8183	298	34	8	0.9067	2.8137	374	34	8	2.0441	15709.4822	499
35	7	0.8889	2.5653	292	35	1	n.a.	n.a.	470	35	9	8.6214	631.4390	499
36	5	0.2983	2.8161	298	36	10	3.6183	406.4338	499	36	7	1.0028	502.4863	499
37	5	0.2973	2.8165	298	37	7	2.0509	729.5399	499	37	1	n.a.	n.a.	499
38	5	0.2979	2.8161	298	38	8	0.9922	2.6372	352	38	1	n.a.	n.a.	499
39	5	0.2973	2.8164	298	39	9	1.2467	32.5807	487	39	7	2.0152	147.8561	482
40	7	1.4684	206.7596	491	40	1	n.a.	n.a.	467	40	2	0.7452	106.3408	445
41	5	0.2975	2.8164	298	41	1	n.a.	n.a.	499	41	8	2.7086	503.6955	499
42	5	0.2959	2.8188	298	42	8	1.8052	5702.3187	499	42	6	2.6311	41958.7468	499
43	9	1.0727	2.5370	297	43	9	1.0551	2.6401	363	43	8	2.3926	889.5130	499
44	5	0.2975	2.8163	298	44	1	n.a.	n.a.	499	44	6	0.9915	904.8396	499
45	9	1.2583	69.7734	465	45	9	1.6941	22.7070	461	45	5	1.6410	564.5119	499
46	9	1.2287	150.8042	494	46	1	n.a.	n.a.	499	46	1	n.a.	n.a.	499
47	8	1.3926	246.4092	499	47	7	1.7933	317.9108	499	47	1	n.a.	n.a.	499
48	5	0.2959	2.8183	298	48	9	1.1090	3.3827	386	48	7	1.1248	208.1925	499
49	7	2.1666	222.3385	487	49	9	1.1030	3.8625	396	49	1	n.a.	n.a.	499
50	5	0.2983	2.8161	298	50	7	1.0535	1192.2263	499	50	1	n.a.	n.a.	499
51	5	0.2958	2.8192	298	51	10	1.9694	1978.7778	499	51	2	0.8228	102.4013	450
52	9	1.4211	51.5298	445	52	5	0.2982	2.8161	353	52	7	0.8945	2.7744	378
53	5	0.2976	2.8162	298	53	8	5.0827	783.8008	499	53	7	1.0473	1132.3572	499
54	5	0.2970	2.8168	298	54	7	1.9763	3376261.9184	499	54	1	n.a.	n.a.	499
55	9	1.3163	988.7258	499	55	9	1.2069	39.0237	497	55	6	1.0058	582.4626	499
56	5	0.2982	2.8161	298	56	5	0.2983	2.8161	353	56	1	n.a.	n.a.	499
57	9	1.2233	2.2424	292	57	5	0.2894	2.9133	354	57	10	6.9798	282791.2331	499
58	8	1.7596	428.0865	499	58	7	1.0936	701.7895	499	58	1	n.a.	n.a.	469
59	5	0.2978	2.8161	298	59	1	n.a.	n.a.	499	59	6	0.9993	464.6136	499
60	5	0.2982	2.8161	298	60	8	2.7499	133.5437	499	60	1	n.a.	n.a.	499

for $\sigma = 1.5$ and $\sigma = 0.5$: results show a general tendency of the radius to influence the ability of the SOM to kill extra units—this ability seems to increase as the radius of the neighborhood narrows.

3.3 Competitive Clustering with Rival Penalization

Lastly, we have analysed the performance of our RPCL implementation in discovering the correct number of clusters. As suggested in [5] we have chosen a number of clusters, $K = 10$, larger than the true number of clusters. Recall that the de-learning rate β is always at least one order of magnitude smaller than η . Once again we considered 60 “Forgy” initialization of the algorithm, and ran three type-B tests using different learning rates, namely $\eta = 0.1$, $\eta = 0.3$ and $\eta = 0.5$. Results are reported in Tables 3a through 3c, where we have indicated with k the number of partitions in the resulting clusterization, and with $\#$ the trial number. For each value of η , we have highlighted the best result according to the Davies-Bouldin index.

Results reveal the following aspects. As expected, the algorithm exhibits a strong ability to invalidate extra units. Such ability appears to be stronger compared to the SOM, as suggested by the fact that the algorithm has always been able to obtain correct clusterizations of the dataset—i.e. five clusters, associated with extremely good values for the Davies-Bouldin index.

Moreover, this ability is only partially affected by the choice of the initial learning rate η —as Table 3 implies, the RPCL algorithm has been able to obtain correct partitionings of the dataset in the 56.67% of the trials for $\eta = 0.1$, 11.67% for $\eta = 0.3$ and 5% for $\eta = 0.5$. In this respect, a higher learning rate does augment the ability of the rival units to move in the data space, and hence the ability of the algorithm to invalidate extra clusters¹, but the success or failure of the algorithm is ultimately due to the goodness of the initialization of the prototypes. As a consequence, we expect the algorithm to succeed independently of the learning rate as long as the number of initializations tested is large enough.

4 Discussion and Conclusion

In conclusion, all the algorithms tested work reasonably well and produce good clusterizations of the dataset. However, if the main task is to make use of one such methods to discover the number of clusters in a given dataset, the rival-penalized competitive learning approach appears to be more robust and practical, since it exhibits a remarkable ability to invalidate extra units—i.e. clusters—depending on the prototype initialization, provided the number of initializations tested is large enough. Hence, this method comes in handy when the number of clusters of a dataset is unknown.

If the number of clusters is not known exactly but it is known to belong to a range of a few possible values, then the self-organizing map can also guess the correct number of clusters and yield a good clusterization, providing multiple, random initializations are tested and the prototypes are drawn from the input dataset. However, the performance

¹ Note that, in some cases, this ability has reached a point in which the algorithm produced a 1-cluster partitioning, for which the Davies-Bouldin index is structurally not defined and the quadratic error loses its significance.

of the SOM exhibits a strong dependency on the value of the parameters, and finding the optimal values for the radius and the step-length can be a challenging task. In this respect, the rival-penalized competitive learning approach is to be preferred over the SOM. Moreover, the SOM involves greater workload compared to the rival-penalized method or the standard competitive learning method, and hence its use in a context where the spatial organization of the output units is of little or no interest appears to be questionable.

Lastly, if the number of clusters is known in advance and the goal is simply to produce a clusterization of the dataset, the basic competitive learning algorithm works fairly well and is less susceptible to the initialization of the prototypes and does not suffer from the dead unit problem. Additionally, it has the highest performance-to-cost ratio, although the number of clusters and the learning rate can play a crucial role in this respect.

References

1. Forgy, E.W.: Cluster analysis of multivariate data: efficiency vs interpretability of classifications. *Biometrics* 21, 768–769 (1965)
2. Xu, R., Wunsch, D.: Survey of Clustering Algorithms. *IEEE Transactions on Neural Networks* 16, 645–678 (2005)
3. Sugar, C.A., James, G.M.: Finding the number of clusters in a dataset: an information theoretic approach. *Journal of the American Statistical Association* 98, 750–763 (2003)
4. Martinetz, T.M., Berkovich, S.G., Schulten, K.J.: “Neural Gas” network for vector quantization and its application to time-series prediction. *IEEE Transactions on Neural Networks* 4, 558–569 (1993)
5. Budura, G., Botoca, C., Miclău, N.: Competitive Learning Algorithms for Data Clustering. *Electronics and Energetics* 19, 261–269 (2006)
6. Kohonen, T.: Self-organized formation of topologically correct feature maps. *Biological Cybernetics* 43, 59–69 (1982)
7. Erwin, E., Obermayer, A., Schulten, K.: Self-organizing maps: ordering, convergence, properties and energy functions. *Biological Cybernetics* 67, 47–55 (1992)
8. Bezdek, J.C., Pal, N.R.: Some new indexes of cluster validity. *IEEE Transactions on Systems, Man, and Cybernetics* 28, 301–315 (1998)
9. Bienenstock, E.L., Cooper, L.N., Munro, P.W.: *Neuroscience* 2, 32–48 (1982)
10. Kirkpatrick, S., Gelatt, C., Vecchi, M.: Optimization by Simulated Annealing. *Science* 220, 671–680 (1983)
11. Ferrari, S., Ferrigno, G., Piuri, V., Borghese, N.A.: Reducing and Filtering Point Clouds with Enhanced Vector Quantization. *IEEE Transactions on Neural Networks* 18, 161–177 (2007)
12. Uchiyama, T., Arbib, M.A.: An algorithm for competitive learning in clustering problems. *Pattern Recognition* 27, 1415–1421 (1994)
13. Fritzke, B.: A growing Neural Gas network learns topologies. *Advances in Neural Information Processing Systems* 7, 625–632 (1995)
14. Bezdek, J.: *Pattern recognition with fuzzy objective function algorithms*. Plenum Press, New York (1981)
15. King, I., Lau, T.-K.: Non-hierarchical Clustering with Rival Penalized Competitive Learning for Information Retrieval. In: Perner, P., Petrou, M. (eds.) *MLDM 1999*. LNCS (LNAI), vol. 1715, pp. 116–130. Springer, Heidelberg (1999)
16. Bação, F., Lobo, V., Painho, M.: Self-organizing Maps as Substitutes for K-Means Clustering. In: Sunderam, V.S., van Albada, G.D., Sloot, P.M.A., Dongarra, J. (eds.) *ICCS 2005*. LNCS, vol. 3516, pp. 476–483. Springer, Heidelberg (2005)