

EM v.s. DBSCAN

Miguel Alcón Doganoc
Universitat Politècnica de Catalunya
Barcelona, Spain
miguel.alcon@est.fib.upc.edu

1 INTRODUCTION

Clustering algorithms are really relevant in the scientific community. They are used in several fields like biology, medicine, business and marketing, computer science, etc.

Because of that, we want to analyze and compare two of the most well-known algorithms in cluster analysis, which are Expectation-maximization (EM) and Density-based spatial clustering of applications with noise (DBSCAN). We chose both algorithms because they have really different ideas behind its implementation, so we expect that they are significantly different to be compared.

To be more precise, the main objective of our work is to apply different datasets to both algorithms, analyze the results and their differences to select a winner, if possible.

2 ALGORITHMS

2.1 EM

Expectation-maximization algorithm (EM) [1, 6, 8] is an iterative method to find maximum likelihood or maximum a posteriori (MAP) estimates of parameters in statistical models, where the model depends on unobserved latent variables. The EM iteration alternates between performing an expectation (E) step, which creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters, and a maximization (M) step, which computes parameters maximizing the expected log-likelihood found on the E step. These parameter-estimates are then used to determine the distribution of the latent variables in the next E step.

As a clustering algorithm, given a fixed k clusters, EM computes probabilities of cluster memberships based on one or more probability distributions. Put another way, the EM algorithm attempts to approximate the observed distributions of values based on mixtures of different distributions in different clusters. The goal of the clustering algorithm then is to maximize the overall probability or likelihood of the data, given the (k) clusters.

The results of EM clustering are the classification probabilities of each observation of the data. In other words, each observation belongs to each cluster with a certain probability. Of course, as a final result you can assign observations to clusters, based on the (largest) classification probability.

2.2 DBSCAN

Density-based spatial clustering of applications with noise (DBSCAN) [2, 7] is a data clustering algorithm. Specifically, it is a density-based clustering non-parametric algorithm. So, given a set of points in some space, it groups together points that are closely packed together (points with many nearby neighbors), marking as outliers points that lie alone in low-density regions (whose nearest neighbors are too far away). This algorithm require two parameters:

- ϵ : the maximum distance between two points for them to be considered as in the same neighborhood.
- $minPts$: the minimum number of points required to form a dense region.

Given ϵ and $minPts$, the DBSCAN algorithm can be abstracted into the following steps:

- (1) Find the points in the ϵ neighborhood of every point, and identify the core points with more than $minPts$ neighbors.
- (2) Find the connected components of core points on the neighbor graph, ignoring all non-core points.
- (3) Assign each non-core point to a nearby cluster if the cluster is an ϵ neighbor, otherwise assign it to noise.

3 DATA

In order to perform the experimentation and to compare both clustering algorithms, we created different scenarios using the module dataset of the *scikit-learn* [5] package of Python. These scenarios consist in 2-D points located in different places among the space, forming clusters of different shapes (blob, circle and moon). They are shown in figure 1.

We also used the dataset of the previous work (Data mining techniques applied to medicine), which is the Heart Disease dataset [4], to see how the algorithms behave with realistic data.

4 IMPLEMENTATION

All the experiments were programmed with Python 3. We used the well-know *scikit-learn* [5] Python package to apply the clustering methods and to calculate their accuracy (with the adjusted random score), and *matplotlib* [3] to create the plots of the first and second experiments.

For more details of the implementation, you can always take a look at the Python script that we developed (inside 'src' folder).

5 EXPERIMENTATION

5.1 First experiment

The goal of the first experiment is to observe the differences of applying EM and DBSCAN to the generated data with blob shape. For this, we created an algorithm that given a number of blobs (n_b) and the space between them (s_b), it places the n_b isotropic Gaussian blobs in a well-distributed way, and with centers separated at a distance of s_b (see table 1). We tried EM and DBSCAN for several cases generated with the algorithm, but we selected 4 of them that show clearly the differences between both. These are represented in sub-figures 1a, 1b, 1c and 1d.

The selected parameters for EM (k) vary for each dataset, but for DBSCAN we fixed them to $\epsilon = 0.2$ and $minPts = 10$, with which we obtained better results after trying several possibilities. Notice that the colors of the clusters in figures 2 and 3 are not an identifier,

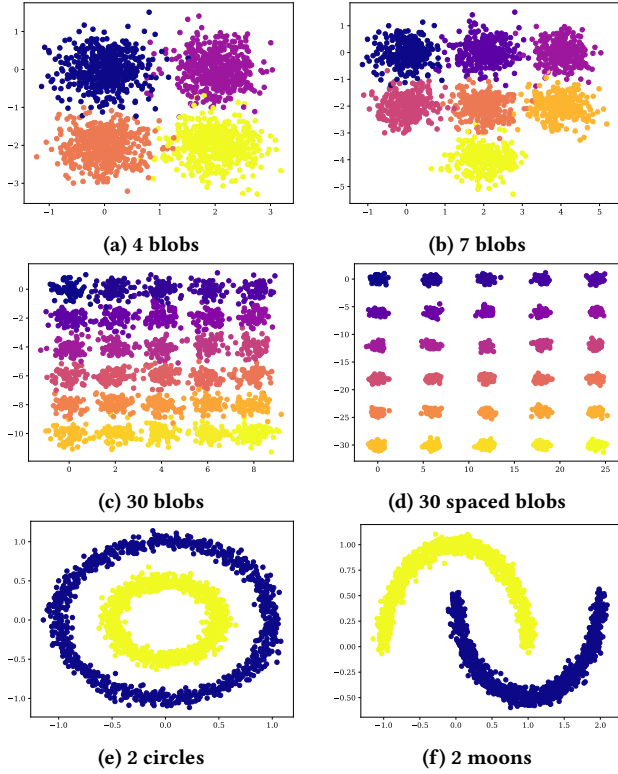


Figure 1: Generated data

Sub-figure	n_b	s_b
1a	4	2
1b	7	2
1c	30	2
1d	30	6

Table 1: Parameters for the generation of the the blob shape data

they are just there to ease the visualization of them. Also, the grey color in DBSCAN means noise. Next, we explain the results of this experiment.

n_b	Algorithm	Accuracy	Sub-figure
4	EM	0.9552	2a
4	DBSCAN	0.8650	2b
7	EM ($k = 3$)	-	2c
7	DBSCAN	0.7575	2d
30	EM	0.9356	2e
30	DBSCAN	0.0111	2f
30 ($s_b = 6$)	EM	1	3a
30 ($s_b = 6$)	DBSCAN	0.0111	3b

Table 2: Accuracy of the predictions of the first experiment

When $k = n_b$, EM is by far the algorithm with more accuracy when identifying the clusters, as we can see in sub-figures 2a, 2b, 2e, 2f, 3a, 3b, and in table 2. Also, when k is incorrect, EM predicts the clusters as better as possible, as seen in sub-figure 2c. But DBSCAN is not bad. In sub-figures 2b and 2d, it predicts well the clusters, identifying the outer points as noise, which is a property that can be desirable in some situations. It is true that in cases with lots of small clusters it fails, as in sub-figures 2f and 3b, but the bad choice of its parameters is not discarded, although we tried several combinations of them. In these sub-figures we can also see that the separation of the clusters does not affect to the predictions, DBSCAN detects roughly the same ones in both cases. This does not happen for EM, which with more space between clusters it predicts them much better.

5.2 Second experiment

The goal of the second experiment is to observe the differences of applying EM and DBSCAN to the generated data with circle and moon shapes. For this, we used the functions `dataset.make_circles` and `dataset.make_moons` of the *scikit-learn* package. The resultant data is shown in sub-figures 1e and 1f.

Unlike the first experiment, and as sub-figures 3c, 3d, 3e and 3f show, DBSCAN predicts the clusters perfectly, while EM is unable to identify them because of their shape.

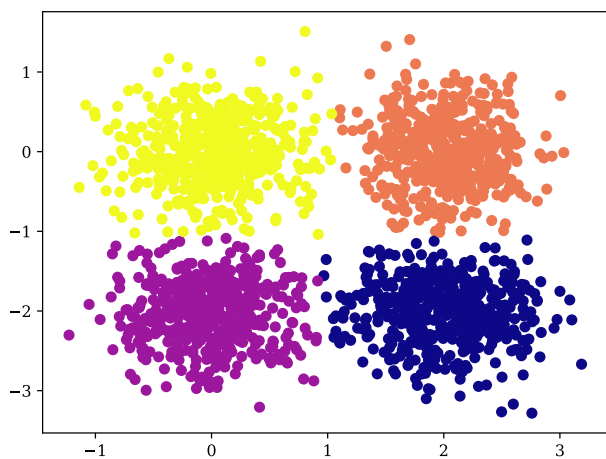
5.3 Third experiment

The goal of the third experiment is to observe the differences of applying EM and DBSCAN to a real dataset (Heart disease), with lots of data and dimensions, and also to compare their predictions with the results obtained in the previous work.

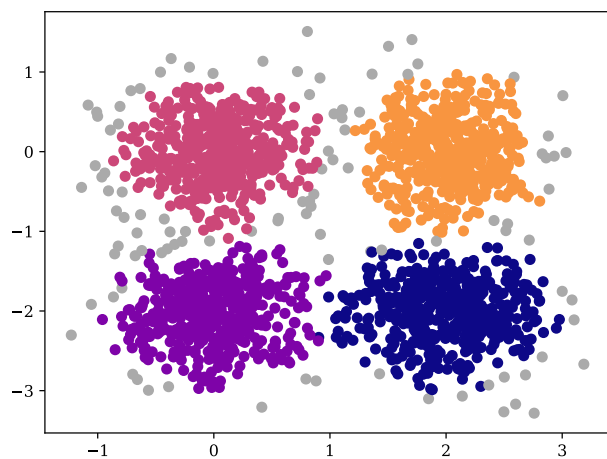
We tried to detect the degree of presence of heart disease in a patient using the provided data. As we did in the previous work, we made the detection first to predict if the patient has presence of heart disease or not (binary decision), and then to predict its degree. For DBSCAN, we selected the parameters in order to obtain the desired number of clusters (with noise being a cluster itself). The parameters are $\{\epsilon = 13, \text{minPts} = 8\}$ for the binary decision, and $\{\epsilon = 14.7, \text{minPts} = 8\}$ for detecting the degree. For EM we selected $k = n_b$. The final results are just too bad to comment them. None of the algorithms cannot be better than selecting the clusters at random. Hence, the classifier process, developed in the previous work, is by far better. If you want to take a look at the results of the experiment, or whichever of the other two, the output of the script is in the file `'experimentation_out.txt'` in the main folder, or you can just execute it.

6 CONCLUSIONS

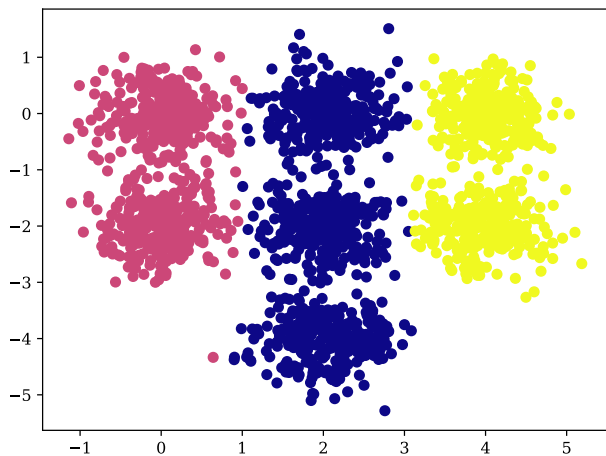
We cannot conclude this work with a clear winner. EM and DBSCAN are both good clustering algorithms, and one is better than the other only depending on the data. As we have seen in the first and second experiment in section 5, and with the appropriate value of k , EM predicts with high accuracy the data with blob shape, but it fails when predicting moons and circles. So, it suffers when points of the same cluster can be far one from the other, but 'connected' by the rest of the points.



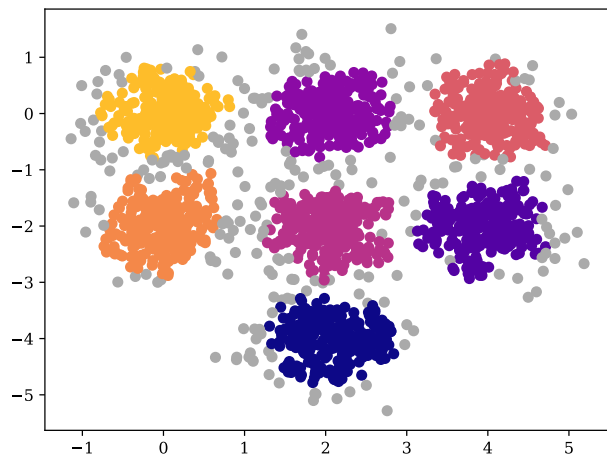
(a) EM $k = 4$



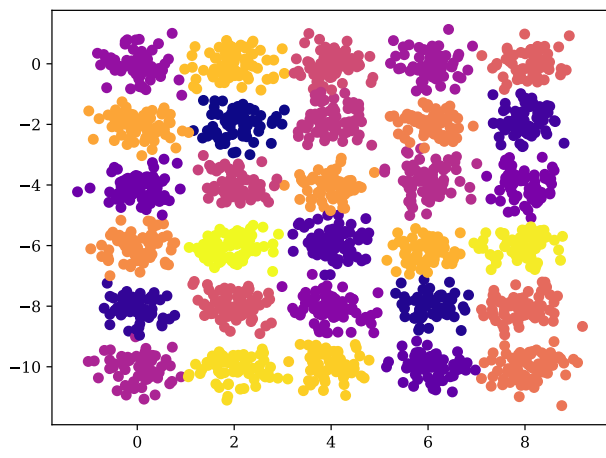
(b) DBSCAN



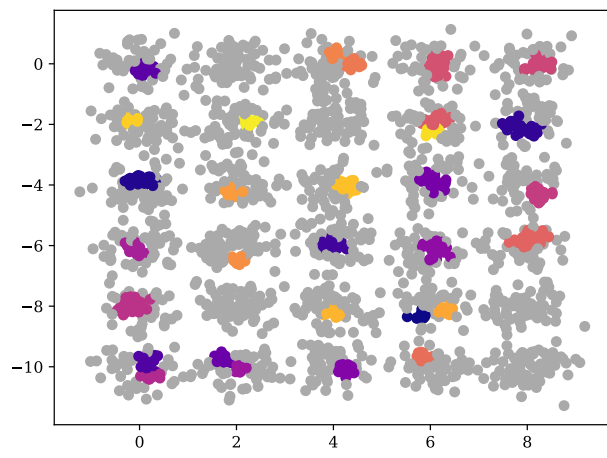
(c) EM $k = 3$



(d) DBSCAN

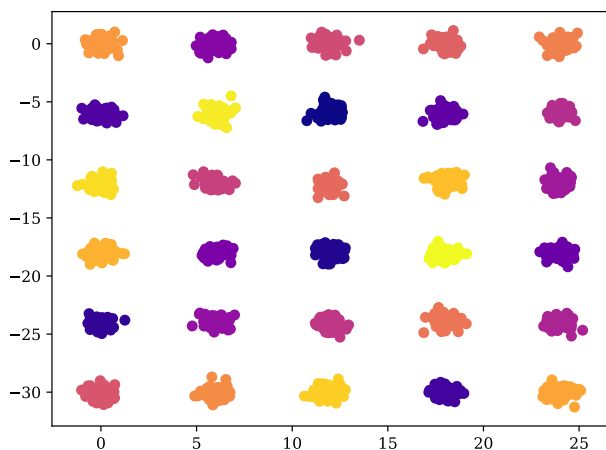


(e) EM $k = 30$

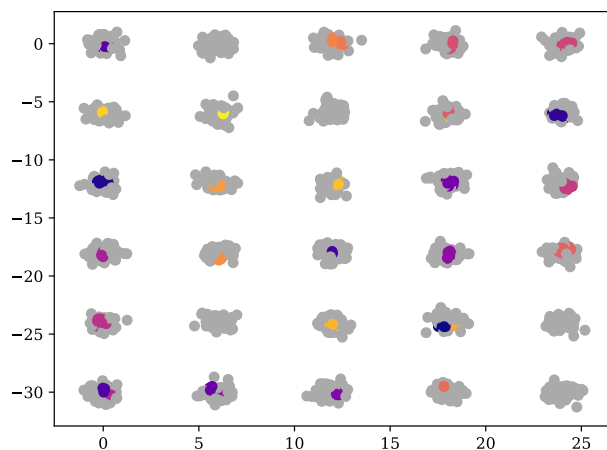


(f) DBSCAN

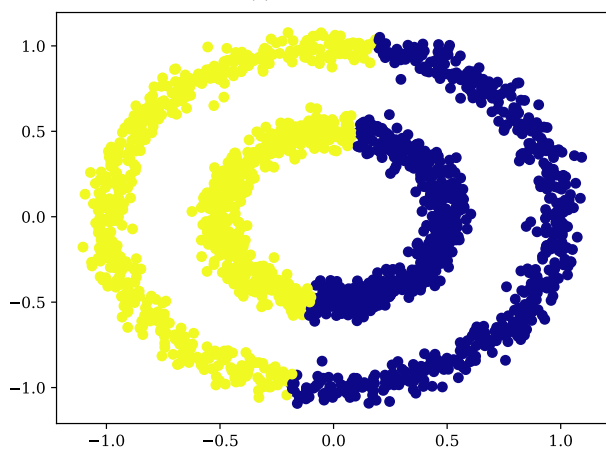
Figure 2: Predictions



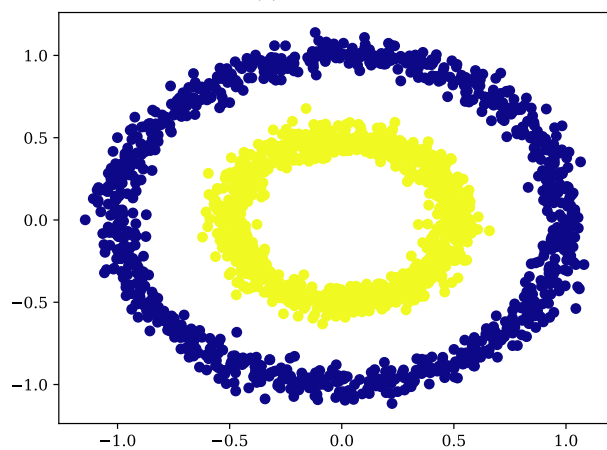
(a) EM $k = 30$



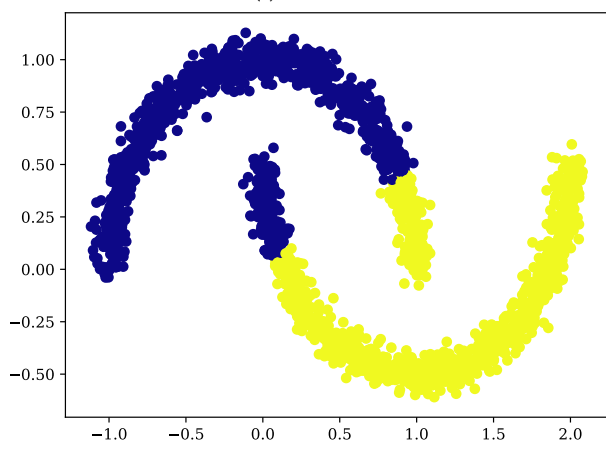
(b) DBSCAN



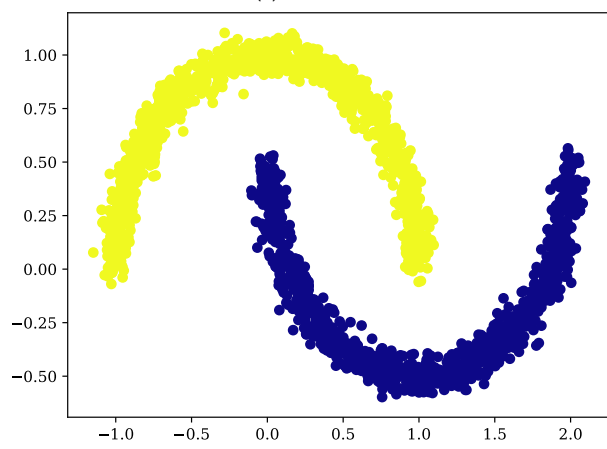
(c) EM $k = 2$



(d) DBSCAN



(e) EM $k = 2$



(f) DBSCAN

Figure 3: Predictions

These could be expected because of the design of the algorithm, explained in section 2.

On the other hand, DBSCAN, with the appropriate parameters, can also predict with high accuracy the data with blob shape, but it fails with a huge number of clusters. Unlike EM, it is also really accurate with the predictions of moons and circles, since it uses the previously mentioned ‘connection’ of the points in order to decide the clusters. Furthermore, the prediction of noise can be really useful with some data.

In both cases, the prediction of the real dataset was horrible. There are problems that the clustering algorithms are just not made for them, and the Heart disease dataset is one of those. Better use the classification algorithms in these cases.

REFERENCES

- [1] DEMPSTER, A. P., LAIRD, N. M., AND RUBIN, D. B. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)* 39, 1 (1977), 1–38.
- [2] ESTER, M., KRIEGLER, H.-P., SANDER, J., AND XU, X. A density-based algorithm for discovering clusters in large spatial databases with noise. AAAI Press, pp. 226–231.
- [3] HUNTER, J. D. Matplotlib: A 2d graphics environment. *Computing In Science & Engineering* 9, 3 (2007), 90–95.
- [4] JANOSI, A., STEINBRUNN, W., PEISTERER, M., AND DETRANO, R. Heart disease. <https://archive.ics.uci.edu/ml/datasets/heart+Disease>, 2019.
- [5] PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B., GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COUNAPEAU, D., BRUCHER, M., PERROT, M., AND DUCHESNAY, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [6] RAPIDMINER. Expectation maximization clustering. https://docs.rapidminer.com/latest/studio/operators/modeling/segmentation/expectation_maximization_clustering.html, 2019.
- [7] WIKIPEDIA. Density-based spatial clustering of applications with noise. <https://en.wikipedia.org/wiki/DBSCAN>, 2019.
- [8] WIKIPEDIA. Expectation–maximization algorithm. https://en.wikipedia.org/wiki/Expectation%E2%80%93maximization_algorithm, 2019.