

A comparative analysis on the bisecting K-means and the PDDP clustering algorithms

Sergio M. Savaresi^{a,*} and Daniel L. Boley^b

^a*Dipartimento di Elettronica e Informazione, Politecnico di Milano, Piazza L. da Vinci, 32, 20133, Milan, Italy*

Tel.: +39 02 2399 3545; Fax: +39 02 2399 3412; E-mail: savaresi@elet.polimi.it

^b*Department of Computer Science and Engineering, University of Minnesota, 4-192 EE/CSci, 200 Union St SE, Minneapolis, MN 55455, USA*

E-mail: boley@cs.umn.edu

Received 19 March 2003

Revised 19 May 2003

Accepted 22 June 2003

Abstract. This paper deals with the problem of clustering a data set. In particular, the bisecting divisive partitioning approach is here considered. We focus on two algorithms: the celebrated K-means algorithm, and the recently proposed Principal Direction Divisive Partitioning (PDDP) algorithm. A comparison of the two algorithms is given, under the assumption that the data set is uniformly distributed within an ellipsoid. In particular, the dynamic behavior of the K-means iterative procedure is studied and discussed; for the 2-dimensional case a closed-form model is given.

Keywords: Unsupervised clustering; K-means; principal direction divisive partitioning

1. Introduction and problem statement

The problem this paper focuses on is the unsupervised clustering of a data set. The data set is given by the matrix $M = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{p \times N}$, where each column of M , $x_i \in \mathbb{R}^p$ is a single data point. This is one of the more basic and common problems in fields like pattern analysis and recognition, data mining, document retrieval, image segmentation, decision making, etc. [10,12].

The specific problem we want to solve herein is the partition of M into *two* sub-matrices (or sub-clusters) $M_L \in \mathbb{R}^{p \times N_L}$ and $M_R \in \mathbb{R}^{p \times N_R}$, $N_L + N_R = N$. This problem is known as *bisecting divisive clustering*.

Note that by recursively using a divisive bisecting clustering procedure, the data set can be partitioned into any given number of clusters. Interestingly enough, the clusters so-obtained are structured as a *hierarchical binary tree* (or a *binary taxonomy*). This is the reason why the bisecting divisive approach is very attractive in many applications (e.g. in content-retrieval/indexing problems – see e.g. [4,18]).

Among the divisive clustering algorithms which have been proposed in the literature in the last two decades (see e.g. [12]), in this paper we will focus on two techniques:

– the *bisecting K-means* algorithm;

– the *Principal Direction Divisive Partitioning (PDDP)* algorithm.

K-means is probably the most celebrated and widely used clustering technique. It is the best representative of the class of iterative centroid-based divisive algorithms. PDDP is a recently proposed technique [3,17]) representative of the non-iterative techniques based upon the Singular Value Decomposition (SVD) of a matrix built from the data set.

The objective of this paper is twofold:

- compare the clustering performance of bisecting K-means and PDDP;
- analyze the dynamic behavior of the K-means iterative algorithm..

In the existing literature, both these issues have been considered only empirically. The performance of PDDP and K-means have been recently studied, and have been reported to be somehow similar, on the basis of a few application examples. The main theoretical results known so far on K-means are [5, 16], where it is shown that the K-means iterative procedure is guaranteed to converge; however, little is said about “where” and “how” it converges.

The main contribution of this work is to provide a simple mathematical explanation of some features of K-means and PDDP. This is done under the restrictive assumption that the data are uniformly distributed within a p -dimensional ellipsoid.

As frequently happens in this type of problems, the theoretical analysis is developed in a very simple and (from the practitioner point of view) ideal setting. This kind of analysis however are of great interest also in real applications, since they provides a deep insight in the behavior of the algorithms, and can give very useful hints and guidelines for the algorithm selection and for the optimal usage of the selected algorithm.

The paper is organized as follows: in Section 2 K-means and PDDP are concisely recalled and discussed; in Section 3 they are analyzed when the number of data points tends to infinity, whereas in Section 4 an empirical analysis in the case of finite data sets is proposed.

2. Bisecting K-means and PDDP

As already stated in the Introduction, this paper focuses on two bisecting divisive partitioning algorithms, which belong to different classes of methods: K-means is the most popular iterative centroid-based divisive algorithm; PDDP is the latest development of SVD-based partitioning techniques. The specific algorithms considered herein are now recalled and briefly commented. In such algorithms the definition of *centroid* will be used extensively; specifically, the centroid of M , say w , is given by

$$w = \frac{1}{N} \sum_{j=1}^N M_j, \quad (1)$$

where M_j is the j -th column of M . Similarly, the centroids of the sub-clusters M_L and M_R , say w_L and w_R , are given by:

$$\begin{cases} w_L = \frac{1}{N_L} \sum_{j=1}^{N_L} M_{L,j} \\ w_R = \frac{1}{N_R} \sum_{j=1}^{N_R} M_{R,j} \end{cases} \quad (2)$$

where $M_{L,j}$ and $M_{R,j}$ are the j -th column of M_L and M_R , respectively.

Bisecting K-means: the algorithm.

Step 1. (Initialization). Randomly select a point, say $c_L \in \mathbb{R}^p$; then compute the centroid w of M (see Eq. (1)), and compute $c_R \in \mathbb{R}^p$ as $c_R = w - (c_L - w)$.

Step 2. Divide $M = [x_1, x_2, \dots, x_n]$ into two sub-clusters M_L and M_R , according to the following rule:

$$\text{rule: } \begin{cases} x_i \in M_L & \text{if } \|x_i - c_L\| \leq \|x_i - c_R\| \\ x_i \in M_R & \text{if } \|x_i - c_L\| > \|x_i - c_R\| \end{cases}$$

Step 3. Compute the centroids of M_L and M_R , w_L and w_R , as in Eq. (2).

Step 4. If $w_L = c_L$ and $w_R = c_R$, stop. Otherwise, let $c_L := w_L, c_R := w_R$, and go back to Step 2.

Bisecting K-means: remarks.

The algorithm above presented is the bisecting version of the general K-means algorithm. This bisecting algorithm has been recently discussed and emphasized in [18,20]. In these works it is claimed to be very effective in document-processing and content-retrieving problems. It is worth noting that the algorithm above recalled is the very classical and basic version of K-means, also known (see [6, 8,10,13,21]) as *Forgy's algorithm*. Many variations of this basic version of the algorithm have been proposed, aiming to reduce the computational demand, at the price of (hopefully little) sub-optimality. Since the goal of this paper is to analyze convergence properties and clustering performance, thanks to its simplicity this original version of the K-means algorithm is the most interesting and meaningful.

PDDP: the algorithm.

Step 1. Compute the centroid w of M as in Eq. (1).

Step 2. Compute the auxiliary matrix \tilde{M} as $\tilde{M} = M - we$, where e is a N -dimensional row vector of ones, namely $e = [1, 1, 1, 1, 1, \dots, 1]$.

Step 3. Compute the Singular Value Decompositions (SVD) of \tilde{M} , $\tilde{M} = U\Sigma V^T$, where Σ is a diagonal $p \times N$ matrix, and U and V are orthonormal unitary square matrices having dimension $p \times p$ and $N \times N$, respectively (see [GV96]).

Step 4. Take the first column vector of U , say $u = U_1$, and divide $M = [x_1, x_2, \dots, x_n]$ into two sub-clusters M_L and M_R , according to the following rule:

$$\begin{cases} x_i \in M_L & \text{if } u^T(x_i - w) \leq 0 \\ x_i \in M_R & \text{if } u^T(x_i - w) > 0 \end{cases}$$

PDDP: remarks.

The PDDP algorithm, recently proposed in [3], belongs to the class of SVD-based data-processing algorithms ([2]); among them, the most popular and widely known are the *Latent Semantic Indexing* algorithm (LSI – see [1,7]), and the LSI-related *Linear Least Square Fit* (LLSF) algorithm ([CY95]). PDDP and LSI mainly differ in the fact that the PDDP splits the matrix with a hyperplane passing through its centroid; LSI through the origin. Another major feature of PDDP is that the SVD of \tilde{M} (Step 3.) can be stopped at the first singular value/vector. This makes PDDP significantly less computationally demanding than LSI, especially if the data-matrix is sparse and the principal singular vector is calculated by resorting to the Lanczos technique ([9,14]).

The main difference between K-means and PDDP is that K-means is based upon an iterative procedure, which, in general, provides different results for different initializations, whereas PDDP is a “one-shot”

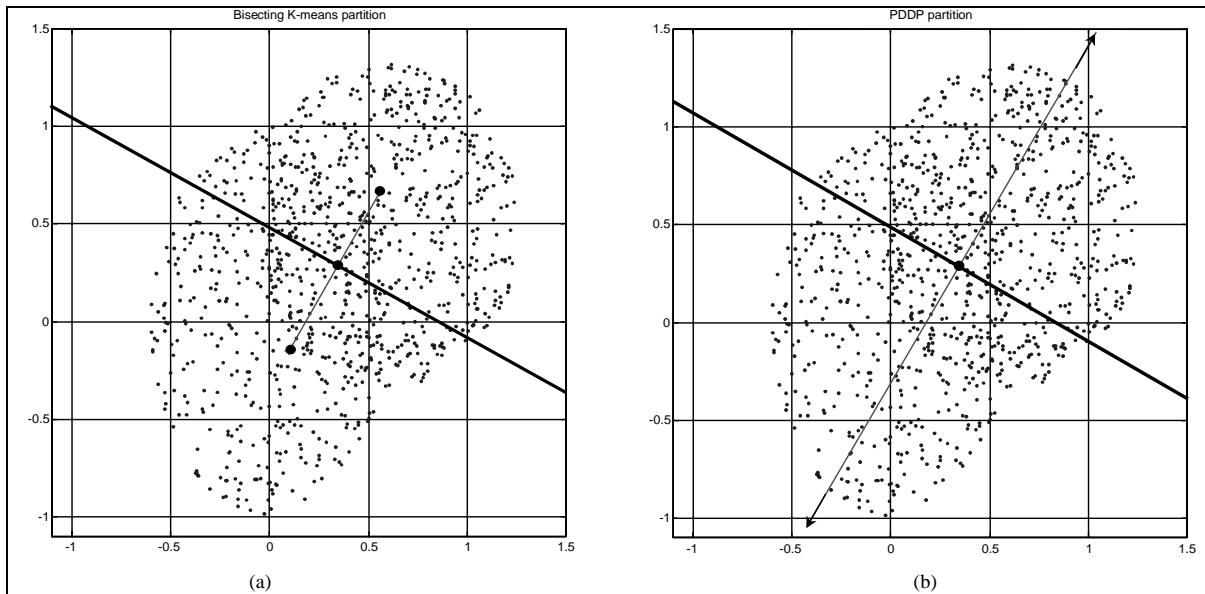


Fig. 1. (a) Partitioning line (bold) of bisecting K-means algorithm. The bullets are the centroids of the data set and of the two sub-clusters; (b) Partitioning line (bold) of PDDP. The bullet is the centroid of the data set. The two arrows show the principal direction of \tilde{M}

algorithm, which provides a unique solution. In order to understand better how K-means and PDDP work, in Fig. 1(a) and Fig. 1(b) the partition of a matrix of dimension provided by K-means and PDDP, respectively, is displayed. From Fig. 1, it is easy to see how K-means and PDDP work:

- the bisecting K-means algorithm splits M with a hyperplane which passes through the centroid w of M , and is perpendicular to the line passing through the centroids w_L and w_R of the sub-clusters M_L and M_R . This is due to the fact that the stopping condition for K-means iterations is that each element of a cluster must be closer to the centroid of that cluster than the centroid of any other cluster.
- PDDP splits M with a hyperplane which passes through the centroid w of M , and is perpendicular to the principal direction of the “unbiased” matrix \tilde{M} (\tilde{M} is the translated version of M , having the origin as centroid). The principal direction of \tilde{M} is its direction of maximum variance (see [GV96]).

At first glance, the two clusters provided by K-means and PDDP look almost indistinguishable. A more careful analysis reveals that the two partitions differ by a few points. Note that this is somewhat unexpected, since the two algorithms differ substantially.

In the rest of the paper we will try to give a rational explanation to the fact that PDDP and bisecting K-means may provide similar results. This will be done by analyzing the dynamic behavior of K-means iteration. Moreover, we will try to clearly outline the *pros* and *cons* of these two seemingly equivalent algorithms.

The analysis presented herein is based upon the restrictive assumption that the points of the data set are uniformly distributed within an ellipsoid. This assumption deserves a few words of comment:

- It is important to point out that an answer to the question “*where does K-means converge?*” can be found only if an assumption of the data distribution is made. Note that this is not mandatory if one only wants an answer to the question “*does the K-means iteration converge?*”. Therefore, the sensible choice of the data distribution becomes the main issue.

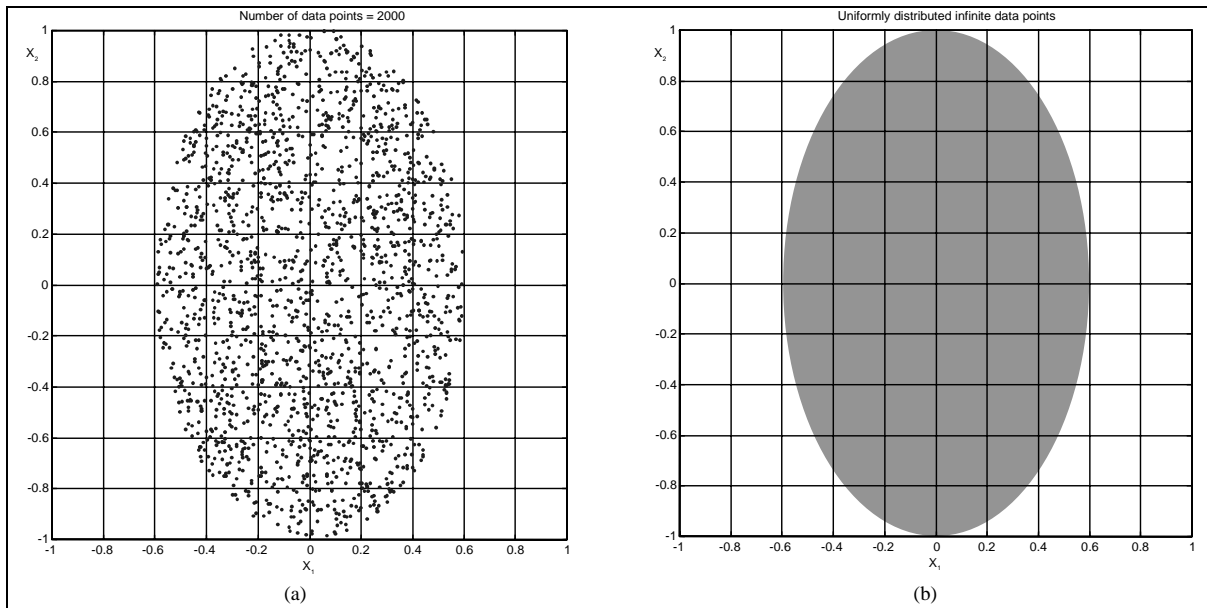


Fig. 2. (a) 2000 data points uniformly distributed within an ellipsoid; (b) Infinite data points uniformly distributed within an ellipsoid.

- Ellipsoid-shaped uniform distribution is the simplest distribution with compact support that, from the clustering point of view, is equivalent to multi-dimensional Gaussian distribution (which is the most typical distribution of experimental data). Henceforth it can be considered the “default” distribution when no a-priori information on the data is available.
- An obvious criticism on the ellipsoid distribution assumption is that the best data clustering is obtained when the data are not distributed like an ellipsoid, but when they are characterized by clearly-separated agglomerations. This is true (even though, in practice, unfortunately this rarely happens), and the analysis presented herein can be extended, in principle, to any given data distribution. However, ellipsoid distribution seems the best compromise between the ambition of considering a realistic data distribution and the need of an easy interpretation of the analysis results.

3. Theoretical results for infinite data sets

In this section the “asymptotic” behavior of bisecting K-means and PDDP will be analyzed. Asymptotic here means that the data set has an infinite number of points, namely $N \rightarrow \infty$. In Fig. 2 the difference between a finite and an infinite set of points is naively depicted.

In the first part of this Section, we will focus on the 2-dimensional case; specifically, it is assumed that each point $x = [x_1, x_2]^T$ of the data set belongs to an ellipsoid centered in the origin and referred to the axes:

$$x = [x_1, x_2]^T \text{ belongs to the data set if: } \frac{x_1^2}{a^2} + x_2^2 \leq 1. \quad (3)$$

The semi-axes lengths of the ellipsoid in Eq. (3) are $a(0 < a \leq 1)$ and 1, respectively.

Given these assumptions, the problem is to find a mathematical description of the dynamic behavior of the bisecting K-means algorithm. This can be done as follows.

- (a) *Parameterization of the splitting line.* First note that the splitting line (the splitting hyperplane in p -dimensional vector spaces) always passes through the origin. This property is preserved even at the first step (see the initialization procedure used in Step.1 – Section 2). Henceforth, the splitting line can be parameterized using one parameter only. The natural choice for this parameter is the angle, say α , between the splitting line and the positive x_1 semi-axis. We shall use the subscript “ t ” to indicate the iteration number. With no loss of generality it is also assumed that $0 \leq \alpha_t \leq \pi/2$.
- (b) *Description of the basic idea.* The basic idea used to compute the mathematical model of the dynamic behavior of bisecting K-means is the following. Given α_t , the next angle α_{t+1} can be calculated by first computing the centroids, say $w_L(\alpha_t)$ and $w_R(\alpha_t)$, of the two semi-clusters induced by the splitting line with angle α_t . The angle α_{t+1} of the next-iteration splitting line then can be easily computed: it is known to be perpendicular to the line connecting $w_L(\alpha_t)$ and $w_R(\alpha_t)$. In this way we obtain a recursive relationship $\alpha_{t+1} = f(\alpha_t)$, which provides a complete description of the dynamic behavior of bisecting K-means.
- (c) *Computation of the centroids.* Due to the infinite number of uniformly distributed points in the data set, the centroids of the two sub-clusters induced by the splitting line with angle α_t must be computed using integral calculus. Using x_2 as integration variable, the computation of the position of w_L (which is the centroid of the “Left” cluster, bordered with a dashed line in Fig. 3) must be split into the computation of the centroids of two sub-pieces of the Left cluster (which are separated by the dashed-dotted line in Fig. 3). The position of w_L hence is given by:

$$w_L = \begin{bmatrix} w_{L1} \\ w_{L2} \end{bmatrix} = \begin{bmatrix} \frac{\int_{-s}^s \frac{1}{2} \left(\frac{\cos(\alpha_t)}{\sin(\alpha_t)} x_2 - a\sqrt{1-x_2^2} \right) \cdot \left(\frac{\cos(\alpha_t)}{\sin(\alpha_t)} x_2 + \sqrt{1-x_2^2} \right) dx_2}{\int_{-s}^s \left(\frac{\cos(\alpha_t)}{\sin(\alpha_t)} x_2 + a\sqrt{1-x_2^2} \right) dx_2} \\ \frac{\int_{-s}^s x_2 \cdot 2a\sqrt{1-x_2^2} dx_2}{\int_{-s}^s 2a\sqrt{1-x_2^2} dx_2} + \frac{\int_{-s}^s x_2 \cdot \left(\frac{\cos(\alpha_t)}{\sin(\alpha_t)} x_2 + a\sqrt{1-x_2^2} \right) dx_2}{\int_{-s}^s \left(\frac{\cos(\alpha_t)}{\sin(\alpha_t)} x_2 + a\sqrt{1-x_2^2} \right) dx_2} \end{bmatrix}, \quad (4)$$

where S is the x_2 -coordinate of the intersection between the splitting line and the ellipsoid in the first quadrant (see Fig. 3); its expression is given by:

$$S = \frac{a \cdot \sin(\alpha_t)}{\sqrt{\cos^2(\alpha_t) + a^2 \sin^2(\alpha_t)}}. \quad (5)$$

Both Eqs (4) and (5) hold for $0 < a \leq 1$ and $0 \leq \alpha_t \leq \pi/2$. The integration of Eq. (4) is complicated, and calls for a symbolic manipulation tool. Fortunately, Eq. (4) can be explicitly computed and significantly simplified. The simplified expression of w_L is given by:

$$w_L = \begin{bmatrix} w_{L1} \\ w_{L2} \end{bmatrix} = \begin{bmatrix} \frac{4}{3} \frac{a^2 \sin(\alpha_t)}{\pi \sqrt{\cos^2(\alpha_t) + a^2 - a^2 \cos^2(\alpha_t)}} \\ \frac{4}{3} \frac{\cos(\alpha_t)}{\pi \sqrt{\cos^2(\alpha_t) + a^2 - a^2 \cos^2(\alpha_t)}} \end{bmatrix}, \quad 0 < a \leq 1, \quad 0 \leq \alpha_t \leq \pi/2.$$

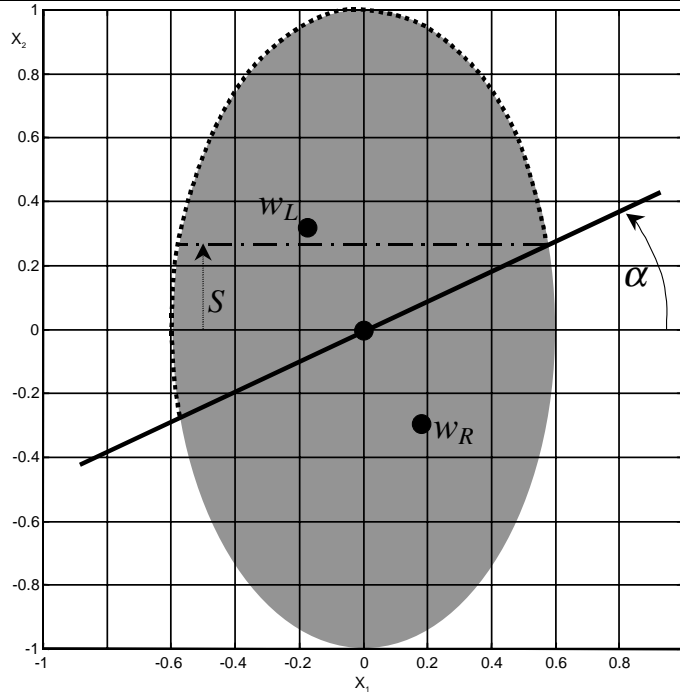


Fig. 3. Parameterization of the splitting-line in K-means.

It is trivial to see that w_R is given by $w_R = -w_L$.

- (d) *The dynamic model of bisecting K-means.* Once $w_L(\alpha_t)$ and $w_R(\alpha_t)$ have been found, it is easy to compute the recursive function $\alpha_{t+1} = f(\alpha_t)$ which models the transition from α_t to the angle α_{t+1} of the next-iteration splitting line. Indeed, this line must be perpendicular to the line passing through w_L and w_R . Henceforth, α_{t+1} can be obtained as:

$$\alpha_{t+1} = \text{atan} \left[-\frac{w_{L1}}{w_{L2}} \right] = \left[\begin{array}{c} -\frac{4}{3} \frac{a^2 \sin(\alpha_t)}{\pi \sqrt{\cos^2(\alpha_t) + a^2} - a^2 \cos^2(\alpha_t)} \\ -\frac{4}{3} \frac{\cos(\alpha_t)}{\pi \sqrt{\cos^2(\alpha_t) + a^2} - a^2 \cos^2(\alpha_t)} \end{array} \right].$$

The following simple expression is finally obtained:

$$\alpha_{t+1} = \text{atan} [a^2 \tan(\alpha_t)], 0 < a \leq 1, 0 \leq \alpha_t \leq \pi/2. \quad (6)$$

Equation (6) is one of the major results of this work, since it provides a rigorous closed-form explicit expression of the dynamic behavior of bisecting K-means. Note that Eq. (6) represents a first order autonomous (i.e. without forcing inputs) non-linear dynamic discrete-time system. As such, it can be analyzed using non-linear systems theory (see e.g. [15,19]). The analysis of Eq. (6) reveals that:

- By solving the steady-state equation

$$\bar{\alpha} = \text{atan} [a^2 \tan(\bar{\alpha})],$$

it is easy to see that the iterative K-means procedure can only have two stationary points, at $\bar{\alpha} = 0$

and $\bar{\alpha} = \pi/2$. In correspondence to these points the ellipsoid is divided by its shorter axis ($\bar{\alpha} = 0$), and by its longer axis ($\bar{\alpha} = \pi/2$), respectively.

- By locally linearizing the dynamic system Eq. (6) about the admissible equilibrium points (namely by computing the tangent model $\delta\alpha_{t+1} = ((\partial f(\alpha_t)/\partial\alpha_t)|_{\alpha_t=\bar{\alpha}})\delta\alpha_t$, where $\delta\alpha_t := \alpha_t - \bar{\alpha}$), we obtain the following two linear dynamic discrete-time systems:

Local dynamic behavior about $\bar{\alpha} = 0$: $\delta\alpha_{t+1} = (a^2)\delta\alpha_t$, $\delta\alpha_t := \alpha_t - 0$;

Local dynamic behavior about $\bar{\alpha} = \pi/2$: $\delta\alpha_{t+1} = (1/a^2)\delta\alpha_t$, $\delta\alpha_t := \alpha_t - \pi/2$.

From linear discrete-time dynamic system theory we know that, if $0 < a < 1$, the linear system about $\bar{\alpha} = 0$ is *asymptotically stable*, and the linear system about $\bar{\alpha} = \pi/2$ is *unstable* (indeed they have poles in a^2 and in $1/a^2$, respectively). This means that bisecting K-means *always converges towards* $\bar{\alpha} = 0$, unless the algorithm is exactly initialized with $\alpha_0 = \pi/2$. In Fig. 4 the function Eq. (6) is displayed, when $a = 0.6$, and a simulated movement of system Eq. (6) is illustrated. Note that, whatever α_0 is (except in the case $\alpha_0 = \pi/2$) the dynamic system $\alpha_{t+1} = f(\alpha_t)$ converges in $\alpha = 0$.

- The value of a strongly affects the number of iterations taken by the algorithm to converge. Thanks to Eq. (6) this number can be given an estimate using dynamic systems theory. First note that the linear system described by the recursive equation $\delta\alpha_{t+1} = (a^2)\delta\alpha_t$ only asymptotically converges at its equilibrium point. A measure of the “speed” at which the system converges towards the equilibrium is given by the so-called *time-constant* τ . τ is defined as the number of steps that $\delta\alpha_t$ takes to decrease its distance from 0 by a factor $1/e$, and it is related to a by the following relationship:

$$\tau = \left(-\frac{1}{\log(a^2)} \right).$$

Due to the discrete nature of the distribution, the bisecting K-means algorithm converges in a finite number of steps, say T . T is a function of the number of the data points N (namely it depends on how densely the data are distributed), which is expected to be *proportional to* τ , namely:

$$T = \gamma(N) \cdot \left(-\frac{1}{\log(a^2)} \right). \quad (7)$$

The value of $\gamma(N)$ is difficult to predicted exactly. A rule of thumb often used by control systems practitioner is that, when $\delta\alpha_t$ has reached the 98% of the distance between the initial value and the equilibrium, the system can be considered, in practice, at steady-state. It is easy to see that this corresponds to $\gamma(N) \approx 4$. In Section 4 a numerical validation of this formula will be provided.

Finally note that T may take very different values. For instance (if $\gamma(N) \approx 4$), K-means is expected to take only 10–15 iterations to converge if $a = 0.7$, about 40 iterations are needed if $a = 0.9$, whereas if $a = 0.95$ the algorithm might need 80 iterations to converge.

The analysis above presented is the main contribution of this Section. It can be concisely summarized with the following two propositions, generalized to p dimensions.

- Proposition 1.** If the data points of a data set are uniformly distributed in a 2-dimensional ellipsoid, the semi-axes of the ellipsoid have lengths equal to 1 and a_1 ($0 < a_1 < 1$), and $N \rightarrow \infty$, then the dynamic discrete-time system which models the K-means iterative algorithm is characterized by 2 equilibrium points; 1 point is locally unstable, and 1 is locally stable. In particular, the dynamic model has the form: $\alpha_{t+1} = a \tan(a^2 \tan(\alpha_t))$, $0 < a < 1$, $0 \leq \alpha_t \leq \pi/2$. The splitting hyperplane corresponding to the equilibrium points pass through the origin and is orthogonal to the main axis of the ellipsoid. The

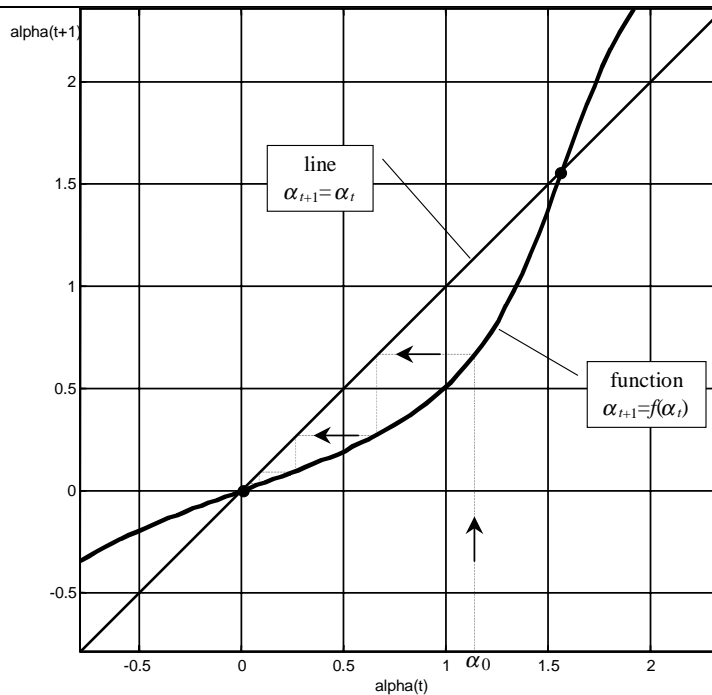


Fig. 4. Function (6) (extended over the range $[-\pi/4; 3\pi/4]$) when $a = 0.6$. The bullets are the equilibria. The thin line is a simulated movement of Eq. (6).

splitting hyperplane corresponding to the stable equilibrium point is orthogonal to the largest axis of the ellipsoid.

Proof. The proof of this result is given in items (a)-(d) above.

Proposition 2. If the data points of a data set are uniformly distributed in a 2-dimensional ellipsoid, the semi-axes of the ellipsoid have lengths equal to 1 and a_1 ($0 < a_1 < 1$), and $N \rightarrow \infty$, then the PDDP algorithm splits the ellipsoid with a hyperplane passing through the origin and orthogonal to the largest axis of the ellipsoid.

Proof. This result is a direct implication of the properties of the SVD. Indeed the 2 singular vectors of a set of points uniformly distributed within an ellipsoid coincide with the direction of the principal axes of the ellipsoid (see [GV96] for details).

Propositions 1 and 2 show that bisecting K-means and PDDP provide the same solution, except in the case when the initialization of K-means exactly corresponds to an unstable equilibrium point of the K-means dynamic model. However, if the initialization is made randomly, this event occurs with probability zero.

These results in principle can be extended to the p -dimensional case, even if the proof is quite cumbersome and lengthy. The p -dimensional case in this paper will be considered in Section 4 with some numerical results.

These asymptotic results are useful to gain a deep insight into the bisecting K-means algorithm, and to explain why, in many cases, K-means and PDDP show a very similar clustering behavior. However, when the data set contains a finite number of data (namely when the number of points is comparatively small), bisecting K-means and PDDP might provide solutions, which, sometimes, are remarkably different. The

finite data set case will be analyzed and discussed in the next section, on the basis of numerical results obtained by simulation.

4. Numerical results for finite data sets

In this section, the bisecting K-means and PDDP will be analyzed when the data set has a finite number of data points. The analysis will be done empirically, using simulated data.

The purpose of this section is twofold:

- validate the theoretical results obtained in the previous section, and see how they change when the data set is finite;
- understand the *pros* and *cons* of K-means and PDDP.

The analysis is structured as follows: in Subsection 4.1 the dynamic model of K-means will be numerically computed for finite data sets, and the problem of local minima will be discussed; in Subsection 4.2 the formula Eq. (7) for the estimation of the number of iterations required by K-means to converge will be validated; in Subsection 4.3 the clustering performance of K-means and PDDP will be compared. Finally, in Subsection 4.4 some conclusions on the *pros* and *cons* of the two algorithms will be drawn.

4.1. The dynamic model of K-means and the problem of local minima

The first problem we consider is the analysis of the K-means dynamic behavior when the data set has a finite number of data. As a first experiment, four sets of data have been considered, characterized by 15, 30, 100 and 2000 data points uniformly distributed within a 2-dimensional ellipsoid with $a = 0.6$. The recursive function $\alpha_{t+1} = f(\alpha_t)$ has been numerically computed for these data sets. The results are displayed in Fig. 5.

- The main difference between the asymptotic function Eq. (6) and the recursive functions corresponding to finite data sets is that the latter are step-wise functions. A major consequence of this function being step-like is that every equilibrium point (namely every point where the function crosses the line $\alpha_{t+1} = \alpha_t$ – see Fig. 5(a)) is locally asymptotically stable, since the local slope of the function about the equilibrium is smaller than 1. *This explains why, in the case of finite data sets, the K-means algorithm is affected by bad “local minima” problems.*
- When the number of data points grows, the finite data set function converges towards the asymptotic function (see Fig. 5(d)). This validates, for the two-dimensional case, the theoretical model developed in the previous section. Moreover, notice that when the number of data points gets large, the number of equilibrium points decreases, and each step gets narrower (see e.g. Fig. 5(c)). *This explains why, when the number of data is sufficiently large, it is the common experience that the problem of local minima tends to vanish.*

As a second experiment, the recursive function $\alpha_{t+1} = f(\alpha_t)$ has been computed for four sets of 15, 30, 100 and 2000 data points uniformly distributed within a 2-dimensional ellipsoid with $a = 0.9$. The results are displayed in Fig. 6. The main difference in the results between the case $a = 0.6$ and $a = 0.9$ is that in the latter the problem of multiple equilibrium points is more severe.

From the inspection of Fig. 5, the following remarks can be done:

These experiments suggest that the problem of local minima for bisecting K-means is expected to:

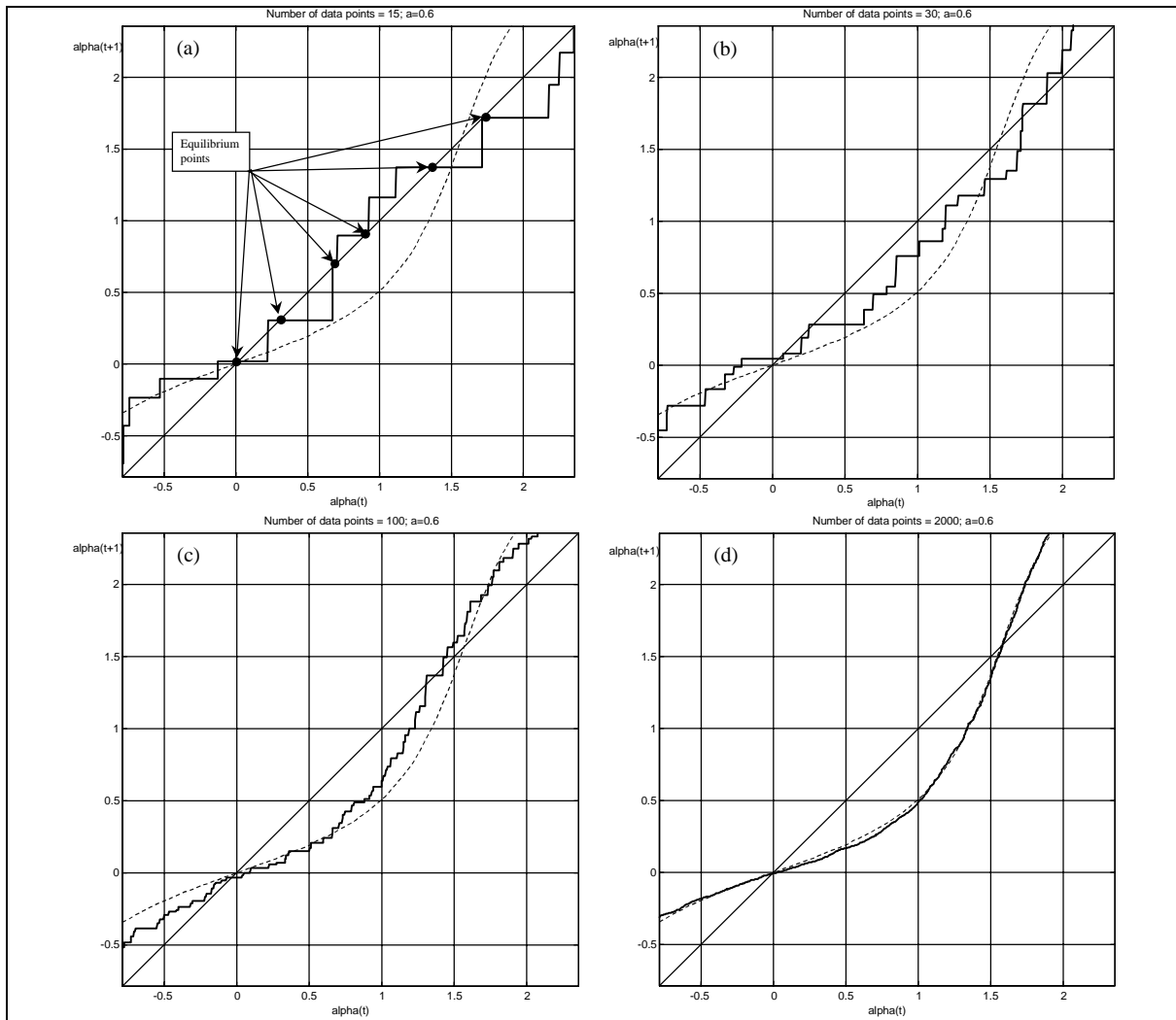


Fig. 5. Recursive function $\alpha_{t+1} = f(\alpha_t)$ estimated from data, when $a = 0.6$. The dashed line is the asymptotic function Eq. (6) computed in Section 4. (a): $N = 15$; (b): $N = 30$; (c): $N = 100$; (d): $N = 2000$.

- *decrease* when the number of data grows;
- *increase* when the size of the “short” semi-axes (a_1, \dots, a_{p-1}) approaches the largest axis.

In order to validate these conjectures, the bisecting K-means algorithm has been extensively tested for different values of a ($a = 0.6, 0.7, 0.8, 0.9$) and for different sizes of the data set (N ranging from 10 to 5000). The average dispersion of the centroids we have obtained (which is directly related to the problem of local minima) is displayed in Fig. 7. In particular, for each value of N , 20 different data sets have been randomly generated; for each data set, 100 different runs of K-means have been done (starting from different initial conditions), so obtaining 100 “dispersed” centroids. The dispersion of these 100 centroids has been computed for each of the 20 data sets, and averaged. Note that the conjectures above outlined are fully confirmed by the data: the centroids dispersion increases with a , and decreases with N .

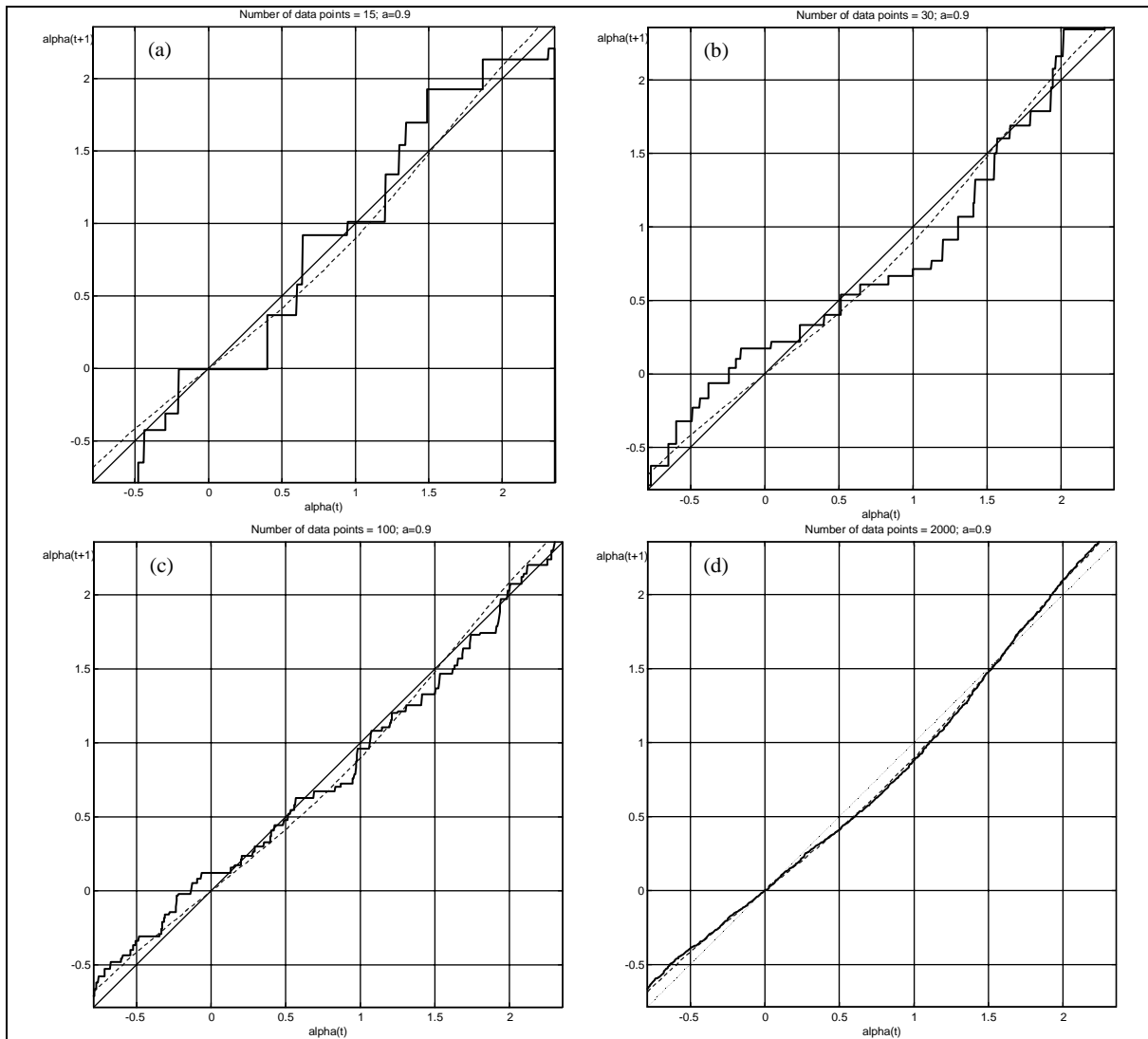


Fig. 6. Recursive function $\alpha_{t+1} = f(\alpha_t)$ estimated from data, when $a = 0.9$. The dashed line is the asymptotic function Eq. (6) computed in Section 4. (a): $N = 15$; (b): $N = 30$; (c): $N = 100$; (d): $N = 2000$.

4.2. The time of convergence of K-means iterations

An interesting result proposed in Section 3, which must be validated, is the prediction of the number of iterations which bisecting K-means needs to converge. Recall that expression Eq. (7) is expected to hold approximately if the data set is large. For small data sets the convergence is expected to be faster.

To this end, the number of iterations required by K-means to converge has been experimentally measured for different values of a in the range $[0.7, 0.95]$, using data sets of size $N = 20000$. The results are in Fig. 8. Notice the very good fit between the predicted and the estimated results (used in Fig. 8 to predict the number of iterations of K-means is $\gamma(N) = 4$, which is the “rule-of-thumb value” suggested in Section 3).

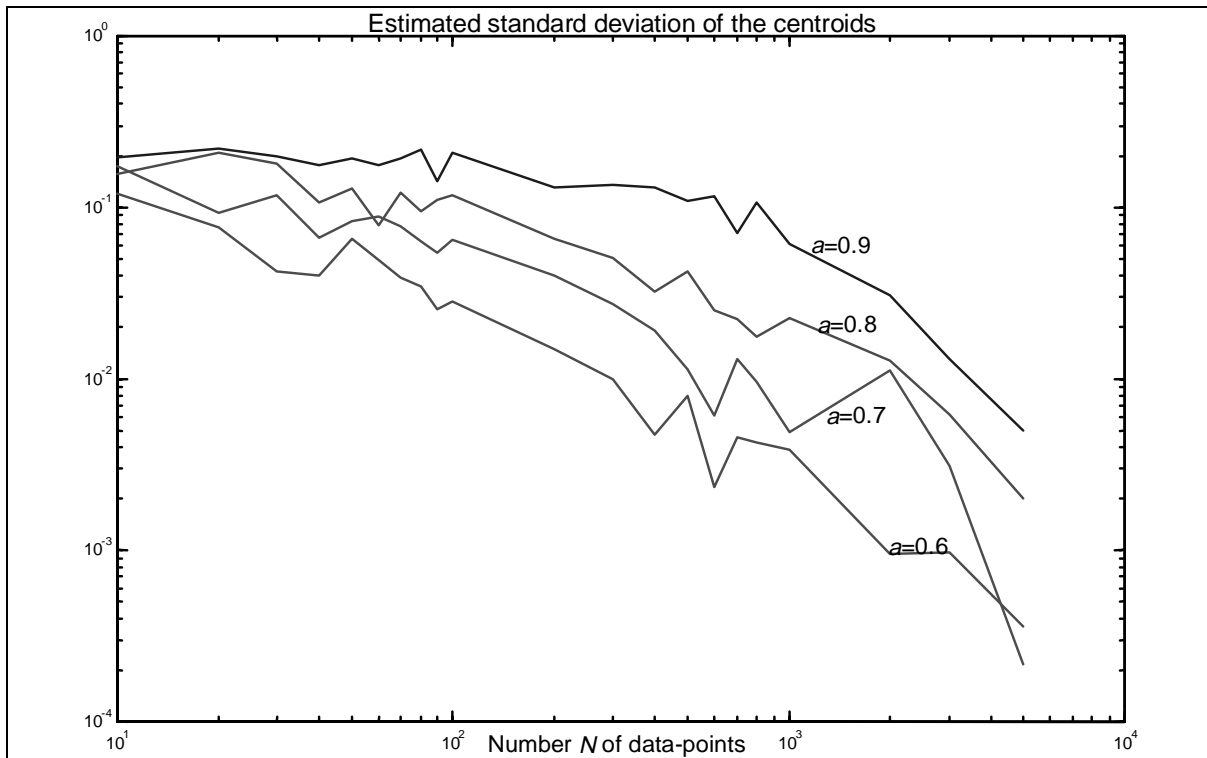


Fig. 7. Average dispersion of the centroids M_L and M_R computed via K-means, as a function of the number of data points. The four lines correspond to different values of a .

4.3. Comparing the clustering performance of bisecting K-means and PDDP

The last crucial issue we consider is the analysis of the clustering performance of K-means and PDDP. This issue is immaterial when the data set is very large, since both methods provide the same results, but is very important when the size of the data set is comparatively small.

To make a comparison between the performance of different clustering algorithms, a performance index must be used. Given two sub-matrices $M_L \in \mathfrak{R}^{N_L}$ and $M_R \in \mathfrak{R}^{N_R}$ of the data set $M = [x_1, x_2, \dots, x_N] \in \mathfrak{R}^{p \times N}$, a widely-accepted way of measuring the internal quality of the partition is given by the following penalty index (see e.g. [11,12,16–18]):

$$J(M_L, M_R) = \sum_{x_i \in M_L} \|x_i - w_L\|^2 + \sum_{x_i \in M_R} \|x_i - w_R\|^2, \quad (8)$$

where w_L and w_R are the centroids of M_L and M_R , given by Eq. (2). Note that Eq. (8) is a measure of cohesiveness of each cluster to its centroid: the smaller $J(M_L, M_R)$ is, the better is the partition.

It is worth pointing out that clustering M by direct minimization of $J(M_L, M_R)$ would be, conceptually, the best clustering method. Unfortunately, the minimization of Eq. (8) is known to require exhaustive search which is exponential in time with respect to the number of data points. Note that the clustering algorithms which have been proposed in the literature (including K-means and PDDP) can be interpreted as *alternate ways of tackling the problem of minimizing* Eq. (8). All of them provide a solution with a reasonable computational effort, at the price of some sub-optimality.

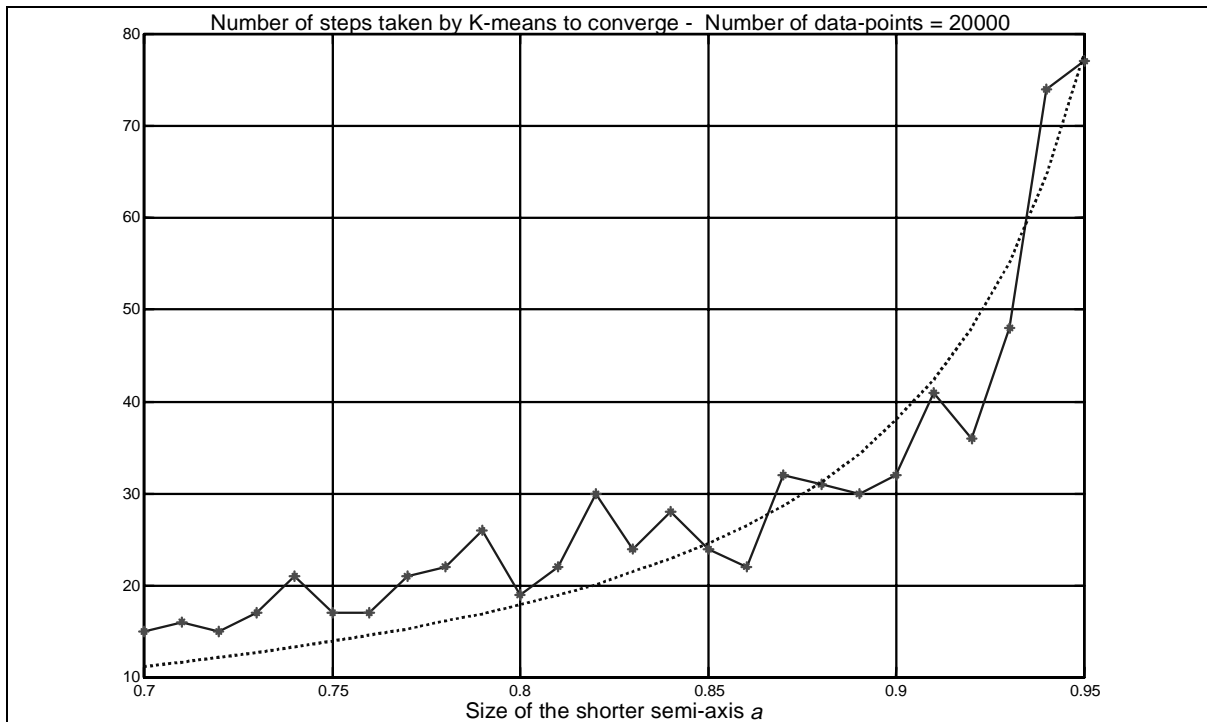


Fig. 8. Estimated number of iterations required by K-means to converge, as a function of a . The dashed line is the number of iterations predicted by Eq. (7), with $\gamma(N) = 4$.

In order to compare the performance of K-means and PDDP, we have considered two sets of data, uniformly distributed in a 100-dimensional ellipsoid. The main semi-axis of the ellipsoid is 1. The size of the remaining 99 semi-axes is in the range $[0.05, 0.95]$. The first data set has 1000 points. The second data set 5000. Note that this number of data points is comparatively small for a 100-dimensional vector space.

For each data set, the following clustering techniques have been used:

- (a) Bisecting K-means, initialized randomly. Specifically, 1000 different initializations have been tested for each data set.
- (b) PDDP.
- (c) Bisecting K-means, initialized with the result provided by PDDP.

At the end of each clustering experiment, the so-obtained partition has been evaluated using Eq. (8). The results are displayed in Fig. 9 ($N = 1000$) and in Fig. 10 ($N = 5000$). The measure of quality in Figs. 9–10 is a normalized version of Eq. (8). Specifically, 0 corresponds to the best clustering performance we have found; 1 corresponds to the “worst-case” situation of non-partitioned cluster (namely $M_L = M$ and $M_R = \emptyset$). The 1000 dots show the clustering performance of K-means randomly initialized; the two horizontal lines show the performance of PDDP, and the performance of K-means initialized via PDDP.

From the inspection of Figs. 9–10, the following remarks can be done:

- The results obtained by random initialization of K-means suffer a remarkably large variation: the corresponding performance index is spread within the $[0, 0.2]$ range. Moreover, notice that K-means may converge towards very “bad” (in terms of clustering performance) solutions (which are

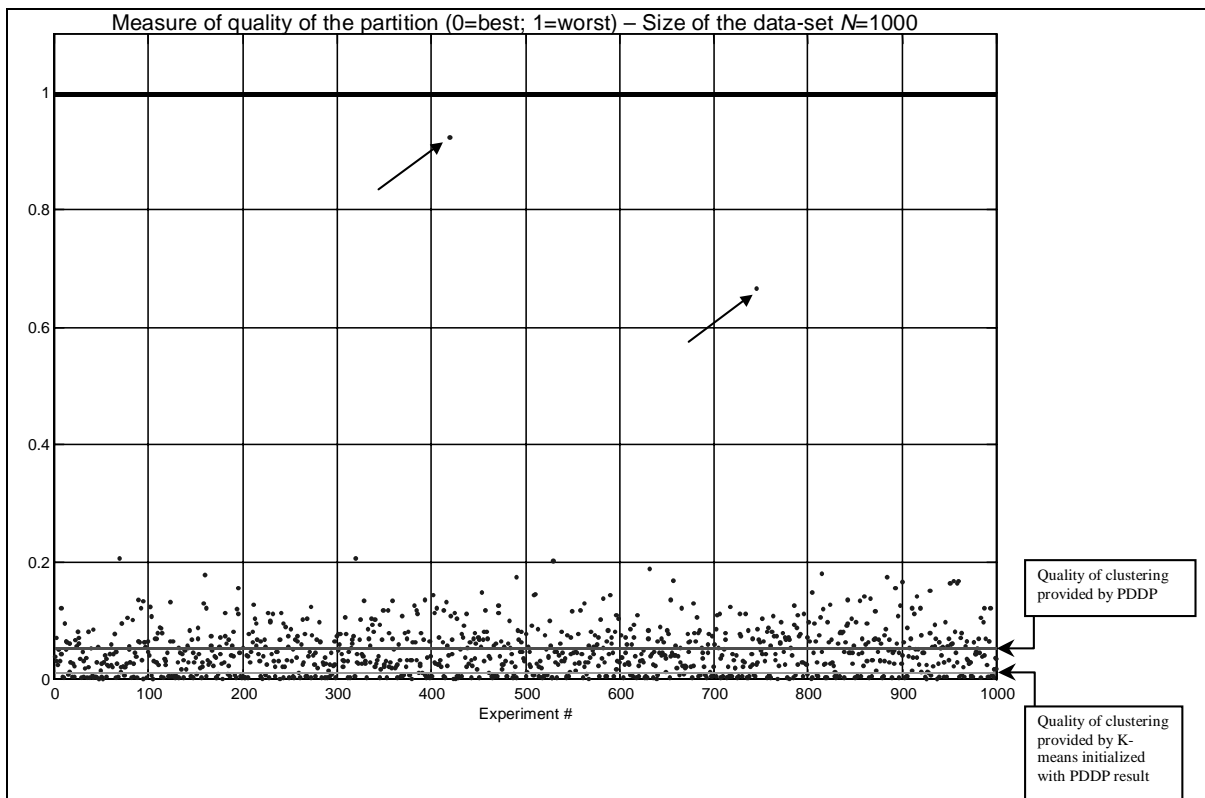


Fig. 9. Measure of quality of bisecting partition of a data set of $N = 1000$ points.

indicated with an arrow in Figs. 9–10). In Fig. 9 there are two “bad” solutions, characterized by a 70% and a 90% (!) performance loss, whereas in Fig.10 the worst solution is characterized by a 40% performance loss.

- In both cases, PDDP slightly under-performs (of about 5%) the best result obtained by K-means. However, it outperforms the worst and the average results of K-means.
- The combination of PDDP and K-means provides very good results. In particular, in the case $N = 1000$ the final performance loss with respect to the best K-means solution is about 1%; in the case $N = 5000$ there is no loss of performance.

To complete this analysis, a few words on the computational power required by the clustering experiments (a)-(c) must be said. To this end, the number of floating point operations (*flops*) spent to cluster the 100×1000 and the 100×5000 data matrices are displayed in Fig. 11.

- The PDDP requires about the number of *flops* required in *average* by a run of K-means. This means that PDDP must be compared with the result of a *single* run of K-means. In light of the performance results displayed in Figs. 9–10, at equal computational power PDDP is expected to provide better performance than K-means.
- As expected, the computational effort required by a run of K-means varies a lot: the minimum and the maximum values of *flops* may differ by an order of magnitude.
- The refinement of the PDDP solution with a run of K-means requires little additional computational power (which – as one intuitively expects – is approximately equal to the minimum of *flops* required

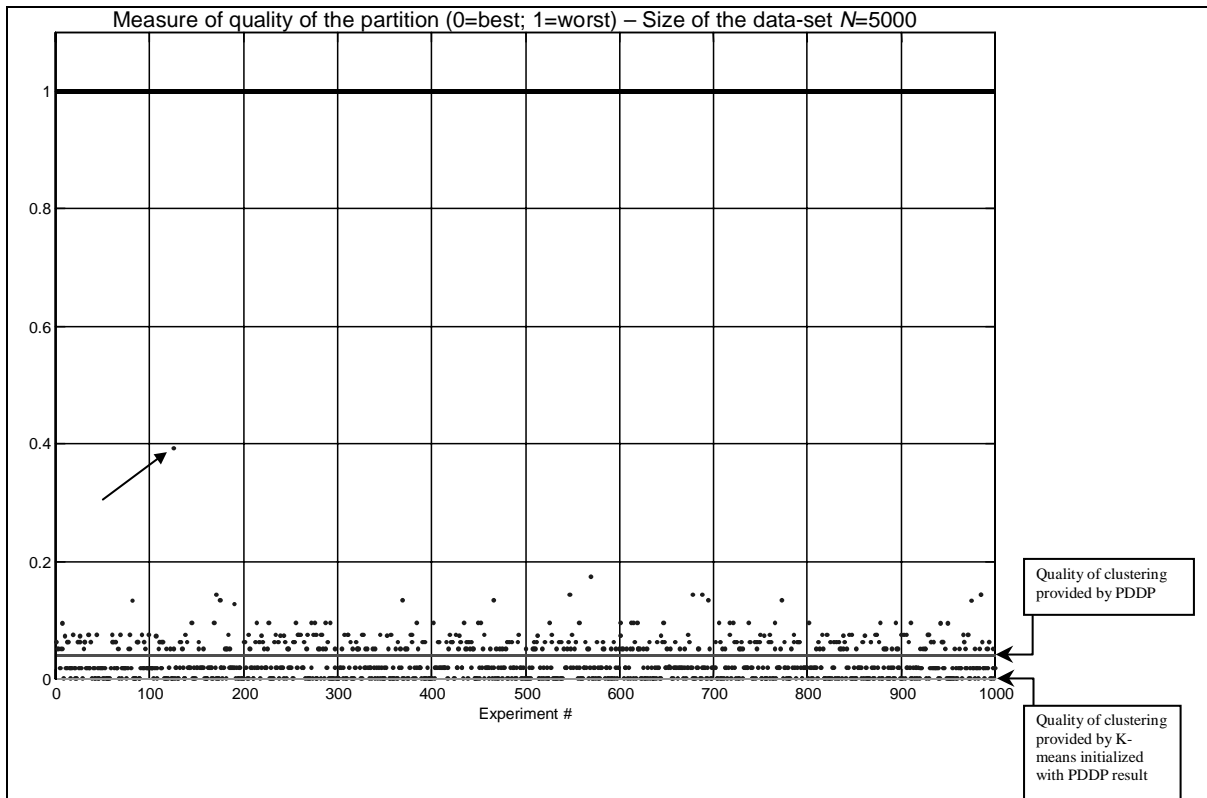


Fig. 10. Measure of quality of bisecting partition of a data set of $N = 5000$ points.

by K-means randomly initialized). This is due to the fact that PDDP provides a solution which is close to an equilibrium point for K-means.

Obviously, the above quick comparison of K-means and PDDP computational demand is far from being exhaustive: the computational power depends on many variables, and is significantly *implementation-dependent* and *data-dependent* (a complete analysis of this issue goes beyond the scope of the present work). For instance, it can be shown that if M tend to be “square” (namely $p \approx N$), PDDP is significantly more demanding than a single run of K-means, whereas, if $p \ll N$, PDDP outperforms K-means (this trend can be clearly observed by comparing the two graphs in Fig. 11). However, it is interesting to see that the above results are very consistent with the results one intuitively expects.

4.4. K-means versus PDDP: concluding remarks

On the basis of the numerical analysis proposed in this Section, we can briefly summarize the *pros* and *cons* of bisecting K-means and PDDP, when the size of the data set is comparatively small:

- K-means is very simple to implement, and tends to give slightly better results in terms of partition quality. However, it is not deterministic (its results strongly depend on the initialization), and it might take a large number of iterations to converge. Hence, if the “best” result is searched for, it is significantly more demanding than PDDP, in terms of computational power.

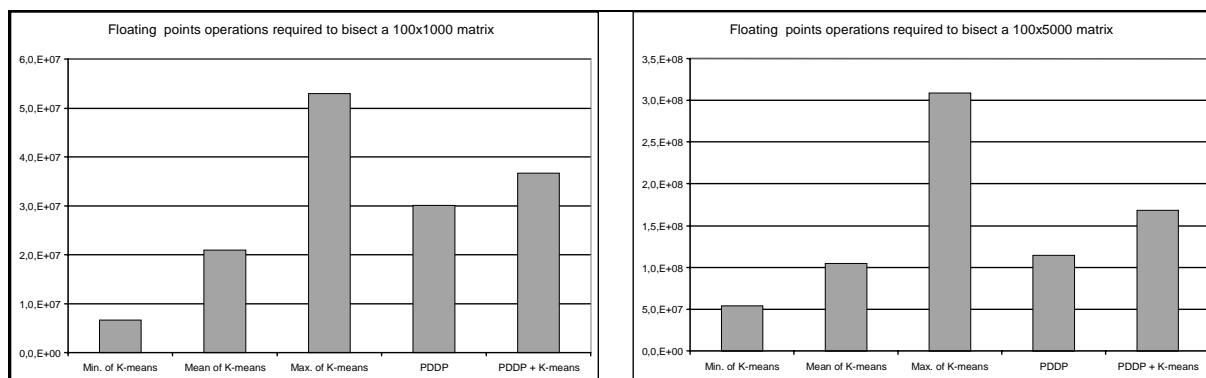


Fig. 11. Comparison of the computational effort required by methods (a)–(c).

- PDDP gives a deterministic result, and it is not affected by the problem of local minima. However it tends to provide results which are slightly worse than the best K-means results.

The peculiar features of K-means and PDDP above outlined can be summarized in the following “rule-of-thumb” for the practitioner:

- The best performance in terms of quality of clustering are obtained by running K-means a large number of times, with different initializations, and picking the best result. However this requires a large computational effort.
- The quickest and safest way of obtaining a “reasonably good” solution is using PDDP.
- The best compromise between computational effort and cluster quality is to use K-means initialized with the PDDP result. This procedure has the additional advantage of providing a deterministic result.

Acknowledgements

First author supported by *Consiglio Nazionale delle Ricerche (CNR) short-term-mobility program*. Second author supported by *NSF grant IIS-9811229*. Thanks are due to Prof. Sergio Bittanti of *Politecnico di Milano*, and to Prof. Giovanna Gazzaniga of *Pavia CNR Institute of Numerical Analysis*.

References

- [1] T. Anderson, On estimation of parameters in latent structure analysis, *Psychometrika* **19** (1954), 1–10.
- [2] M.W. Berry, Z. Drmac, E.R. Jessup, Matrices, Vector spaces, and Information Retrieval, *SIAM Review* **41** (1999), 335–362.
- [3] D.L. Boley, Principal Direction Divisive Partitioning, *Data Mining and Knowledge Discovery* **2**(4) (1998), 325–344.
- [4] D.L. Boley, M. Gini, R. Gross, S. Han, K. Hastings, G. Karypis, V. Kumar, B. Mobasher and J. Moore, Document Categorization and Query Generation on the World Wide Web Using WebACE, *AI Review* **11** (2000), 365–391.
- [5] L. Bottou, Y. Bengio, *Convergence properties of the k-means algorithm*, Advances in Neural Information Processing Systems, 1995.
- [6] C. Chute, Y. Yang An overview of statistical methods for the classification and retrieval of patient events, *Meth. Inform. Med.* **34** (1995), 104–110.
- [7] S. Deerwester, S. Dumais, G. Furnas and R. Harshman, Indexing by latent semantic analysis, *J. Amer. Soc. Inform. Sci.* **41** (1990), 41–50.

- [8] E. Forgy, Cluster Analysis of Multivariate Data: Efficiency versus Interpretability of Classification, *Biometrics* (1965), 768–780.
- [9] G.H. Golub and C.F. van Loan, *Matrix Computations*, (3rd Ed.), The Johns Hopkins University Press, 1996.
- [10] E. Gose, R. Johnsonbaugh and S. Jost, *Pattern Recognition & Image Analysis*, Prentice-Hall, 1996.
- [11] D. Hand, H. Mannila and P. Smyh, *Principles of Data Mining*, MIT Press, 2001.
- [12] A.K. Jain, M.N. Murty and P.J. Flynn, Data Clustering: a Review, *ACM Computing Surveys* **31**(3) (1999), 264–323.
- [13] G.N. Lance and W.T. Williams, A general theory of classificatory sorting strategies, Hierarchical systems, *The Computer Journal* **9** (1967), 373–380.
- [14] C. Lanczos, An iteration method for the solution of the eigenvalue problem of linear differential and integral operators, *J. Res. Nat. Bur. Stand* **45** (1950), 255–282.
- [15] J.P. LaSalle, *The Stability and Control of Discrete Processes*, Springer-Verlag, 1986.
- [16] S.Z. Selim and M.A. Ismail, K-means-type algorithms: a generalized convergence theorem and characterization of local optimality, *IEEE Trans. on Pattern Analysis and Machine Intelligence* **6**(1) (1984), 81–86.
- [17] S.M. Savaresi, D.L. Boley, S. Bittanti and G. Gazzaniga, *Cluster selection in divisive clustering algorithms*, 2nd SIAM International Conference on Data Mining, Arlington, VI, USA, 2002, pp. 299–314.
- [18] M. Steinbach, G. Karipis and V. Kumar, *A comparison of Document Clustering Techniques*, Proceedings of World Text Mining Conference, KDD2000, Boston, 2000.
- [19] M. Vidyasagar, *Nonlinear Systems Analysis*, Prentice-Hall, 1993.
- [20] J.Z. Wang, G. Wiederhold, O. Firschein and S.X. Wei, Content-based image indexing and searching using Daubechies' wavelets, *Int. J. Digit. Library* **1** (1997), 311–328.
- [21] J.H. Ward Jr., Hierarchical groupings to optimise an objective function, *Journal of the American Statistical Association* **58**(301) (1963), 236–244.