

DETERMINISTIC INITIALIZATION OF THE K-MEANS ALGORITHM USING HIERARCHICAL CLUSTERING

M. EMRE CELEBI

*Department of Computer Science
Louisiana State University, Shreveport, LA, USA
ecelebi@lsus.edu*

HASSAN A. KINGRAVI

*School of Electrical and Computer Engineering
Georgia Institute of Technology, Atlanta, GA, USA
kingravi@gatech.edu*

Received 17 July 2012

Accepted 29 October 2012

Published 21 December 2012

K-means is undoubtedly the most widely used partitional clustering algorithm. Unfortunately, due to its gradient descent nature, this algorithm is highly sensitive to the initial placement of the cluster centers. Numerous initialization methods have been proposed to address this problem. Many of these methods, however, have superlinear complexity in the number of data points, making them impractical for large data sets. On the other hand, linear methods are often random and/or order-sensitive, which renders their results unrepeatable. Recently, Su and Dy proposed two highly successful hierarchical initialization methods named Var-Part and PCA-Part that are not only linear, but also deterministic (nonrandom) and order-invariant. In this paper, we propose a discriminant analysis based approach that addresses a common deficiency of these two methods. Experiments on a large and diverse collection of data sets from the UCI machine learning repository demonstrate that Var-Part and PCA-Part are highly competitive with one of the best random initialization methods to date, i.e. k-means++, and that the proposed approach significantly improves the performance of both hierarchical methods.

Keywords: Partitional clustering; sum of squared error criterion; k-means; cluster center initialization; thresholding.

1. Introduction

Clustering, the unsupervised classification of patterns into groups, is one of the most important tasks in exploratory data analysis.¹ Primary goals of clustering include gaining insight into data (detecting anomalies, identifying salient features, etc.), classifying data, and compressing data. Clustering has a long and rich history

in a variety of scientific disciplines including anthropology, biology, medicine, psychology, statistics, mathematics, engineering, and computer science. As a result, numerous clustering algorithms have been proposed since the early 1950s.²

Clustering algorithms can be broadly classified into two groups: hierarchical and partitional.² Hierarchical algorithms recursively find nested clusters either in a top-down (divisive) or bottom-up (agglomerative) fashion. In contrast, partitional algorithms find all the clusters simultaneously as a partition of the data and do not impose a hierarchical structure. Most hierarchical algorithms have quadratic or higher complexity in the number of data points¹ and therefore are not suitable for large data sets, whereas partitional algorithms often have lower complexity.

Given a data set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ in \mathbb{R}^D , i.e. N points (vectors) each with D attributes (components), hard partitional algorithms divide \mathcal{X} into K exhaustive and mutually exclusive clusters $\mathcal{P} = \{P_1, P_2, \dots, P_K\}$, $\bigcup_{i=1}^K P_i = \mathcal{X}$, $P_i \cap P_j = \emptyset$ for $1 \leq i \neq j \leq K$. These algorithms usually generate clusters by optimizing a criterion function. The most intuitive and frequently used criterion function is the sum of squared error (SSE) given by:

$$\text{SSE} = \sum_{i=1}^K \sum_{\mathbf{x}_j \in P_i} \|\mathbf{x}_j - \mathbf{c}_i\|_2^2, \tag{1}$$

where $\|\cdot\|_2$ denotes the Euclidean (\mathcal{L}_2) norm and $\mathbf{c}_i = 1/|P_i| \sum_{\mathbf{x}_j \in P_i} \mathbf{x}_j$ is the centroid of cluster P_i whose cardinality is $|P_i|$. The optimization of (1) is often referred to as the minimum SSE clustering (MSSC) problem.

The number of ways in which a set of N objects can be partitioned into K nonempty groups is given by Stirling numbers of the second kind:

$$\mathcal{S}(N, K) = \frac{1}{K!} \sum_{i=0}^K (-1)^{K-i} \binom{K}{i} i^N, \tag{2}$$

which can be approximated by $K^N/K!$. It can be seen that a complete enumeration of all possible clusterings to determine the global minimum of (1) is clearly computationally prohibitive except for very small data sets.³ In fact, this nonconvex optimization problem is proven to be NP-hard even for $K = 2$ (Ref. 4) or $D = 2$ (Ref. 5). Consequently, various heuristics have been developed to provide approximate solutions to this problem.⁶ Among these heuristics, Lloyd's algorithm,⁷ often referred to as the (batch) k-means algorithm, is the simplest and most commonly used one. This algorithm starts with K arbitrary centers, typically chosen uniformly at random from the data points. Each point is assigned to the nearest center and then each center is recalculated as the mean of all points assigned to it. These two steps are repeated until a predefined termination criterion is met.

The k-means algorithm is undoubtedly the most widely used partitional clustering algorithm.² Its popularity can be attributed to several reasons. First, it is

conceptually simple and easy to implement. Virtually every data mining software includes an implementation of it. Second, it is versatile, i.e. almost every aspect of the algorithm (initialization, distance function, termination criterion, etc.) can be modified. This is evidenced by hundreds of publications over the last fifty years that extend k-means in various ways. Third, it has a time complexity that is linear in N , D , and K (in general, $D \ll N$ and $K \ll N$). For this reason, it can be used to initialize more expensive clustering algorithms such as expectation maximization,⁸ DBSCAN,⁹ and spectral clustering.¹⁰ Furthermore, numerous sequential^{11,12} and parallel¹³ acceleration techniques are available in the literature. Fourth, it has a storage complexity that is linear in N , D , and K . In addition, there exist disk-based variants that do not require all points to be stored in memory.¹⁴ Fifth, it is guaranteed to converge¹⁵ at a quadratic rate.¹⁶ Finally, it is invariant to data ordering, i.e. random shufflings of the data points.

On the other hand, k-means has several significant disadvantages. First, it requires the number of clusters, K , to be specified in advance. The value of this parameter can be determined automatically by means of various internal/relative cluster validity measures.¹⁷ Second, it can only detect compact, hyperspherical clusters that are well separated. This can be alleviated by using a more general distance function such as the Mahalanobis distance, which permits the detection of hyperellipsoidal clusters.¹⁸ Third, due its utilization of the squared Euclidean distance, it is sensitive to noise and outlier points since even a few such points can significantly influence the means of their respective clusters. This can be addressed by outlier pruning¹⁹ or using a more robust distance function such as City-block (\mathcal{L}_1) distance. Fourth, due to its gradient descent nature, it often converges to a local minimum of the criterion function.¹⁵ For the same reason, it is highly sensitive to the selection of the initial centers.²⁰ Adverse effects of improper initialization include empty clusters, slower convergence, and a higher chance of getting stuck in bad local minima.²¹ Fortunately, except for the first two, these drawbacks can be remedied by using an adaptive initialization method (IM).

A large number of IMs have been proposed in the literature.^{22,23,21,20} Unfortunately, many of these have superlinear complexity in N ,^{24–26,3,27–32} which makes them impractical for large data sets (note that k-means itself has linear complexity). In contrast, linear IMs are often random and/or order-sensitive,^{33–38,8,39} which renders their results unrepeatable. Su and Dy proposed two divisive hierarchical initialization methods named Var-Part and PCA-Part that are not only linear, but also deterministic and order-invariant.⁴⁰ In this study, we propose a simple modification to these methods that improves their performance significantly.

The rest of the paper is organized as follows. Section 2 presents a brief overview of some of the most popular linear, order-invariant k-means IMs and the proposed modification to Var-Part and PCA-Part. Section 3 presents the experimental results, while Sec. 4 analyzes these results. Finally, Sec. 5 gives the conclusions.

2. Linear, Order-Invariant Initialization Methods for K-Means

2.1. Overview of the existing methods

Forgy’s method³³ assigns each point to one of the K clusters uniformly at random. The centers are then given by the centroids of these initial clusters. This method has no theoretical basis, as such random clusters have no internal homogeneity.⁴¹

MacQueen³⁵ proposed two different methods. The first one, which is the default option in the quick cluster procedure of IBM SPSS Statistics,⁴² takes the first K points in \mathcal{X} as the centers. An obvious drawback of this method is its sensitivity to data ordering. The second method chooses the centers randomly from the data points. The rationale behind this method is that random selection is likely to pick points from dense regions, i.e. points that are good candidates to be centers. However, there is no mechanism to avoid choosing outliers or points that are too close to each other.⁴¹ Multiple runs of this method is the standard way of initializing k-means.⁸ It should be noted that this second method is often mistakenly attributed to Forgy.³³

The maximin method⁴³ chooses the first center \mathbf{c}_1 arbitrarily and the i th ($i \in \{2, 3, \dots, K\}$) center \mathbf{c}_i is chosen to be the point that has the greatest minimum-distance to the previously selected centers, i.e. $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_{i-1}$. This method was originally developed as a two-approximation to the K -center clustering problem.^a

The k-means++ method³⁹ interpolates between MacQueen’s second method and the maximin method. It chooses the first center randomly and the i th ($i \in \{2, 3, \dots, K\}$) center is chosen to be $\mathbf{x}' \in \mathcal{X}$ with a probability of $\frac{md(\mathbf{x}')^2}{\sum_{j=1}^N md(\mathbf{x}_j)^2}$, where $md(\mathbf{x})$ denotes the minimum-distance from a point \mathbf{x} to the previously selected centers. This method yields an $\Theta(\log K)$ approximation to the MSSC problem.

The PCA-Part method⁴⁰ uses a divisive hierarchical approach based on PCA (Principal Component Analysis).⁴⁴ Starting from an initial cluster that contains the entire data set, the method iteratively selects the cluster with the greatest SSE and divides it into two subclusters using a hyperplane that passes through the cluster centroid and is orthogonal to the principal eigenvector of the cluster covariance matrix. This procedure is repeated until K clusters are obtained. The centers are then given by the centroids of these clusters. The Var-Part method⁴⁰ is an approximation to PCA-Part, where the covariance matrix of the cluster to be split is assumed to be diagonal. In this case, the splitting hyperplane is orthogonal to the coordinate axis with the greatest variance.

Figure 1 illustrates the Var-Part procedure on a toy data set with four natural clusters.⁴⁵ In iteration 1, the initial cluster that contains the entire data set is split into two subclusters along the y -axis using a line (one-dimensional hyperplane) that

^a Given a set of N points in a metric space, the goal of K -center clustering is to find K representative points (centers) such that the maximum distance of a point to a center is minimized. Given a minimization problem, a two-approximation algorithm is one that finds a solution whose cost is at most twice the cost of the optimal solution.

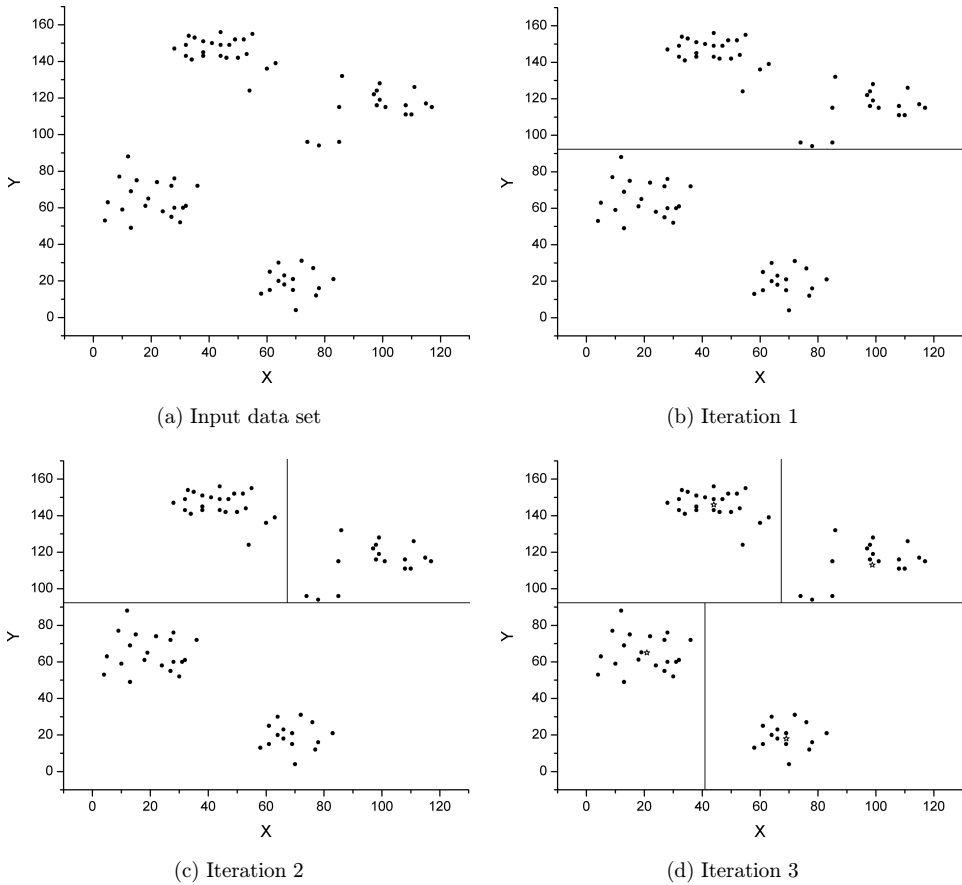


Fig. 1. Illustration of Var-Part on the Ruspini data set.

passes through the mean point (92.026667). Between the resulting two clusters, the one above the line has a greater SSE. In iteration 2, this cluster is therefore split along the x -axis at the mean point (66.975000). In the final iteration, the cluster with the greatest SSE, i.e. the bottom cluster, is split along the x -axis at the mean point (41.057143). In Fig. 1(d), the centroids of the final four clusters are denoted by stars.

2.2. Proposed modification to Var-part and PCA-part

Su and Dy⁴⁰ demonstrated that, besides being computationally efficient, Var-Part and PCA-Part perform very well on a variety of data sets. Recall that in each iteration these methods select the cluster with the greatest SSE and then project the D -dimensional points in this cluster on a partitioning axis. The difference between the two methods is the choice of this axis. In Var-Part, the partitioning axis is the coordinate axis with the greatest variance, whereas in PCA-Part it is the major axis. After the projection operation, both methods use the *mean* point on the partitioning

axis as a “threshold” to divide the points between two clusters. In other words, each point is assigned to one of the two subclusters depending on which side of the *mean* point its projection falls to. It should be noted that the choice of this threshold is primarily motivated by computational convenience. Here, we propose a better alternative based on discriminant analysis.

The projections of the points on the partitioning axis can be viewed as a discrete probability distribution, which can be conveniently represented by a histogram. The problem of dividing a histogram into two partitions is a well studied one in the field of image processing. A plethora of histogram partitioning, a.k.a. thresholding, methods has been proposed in the literature with the early ones dating back to the 1960s.⁴⁶ Among these, Otsu’s method⁴⁷ has become the method of choice as confirmed by numerous comparative studies.^{48–50,46,51,52}

Given an image represented by L gray levels $\{0, 1, \dots, L - 1\}$, a thresholding method partitions the image pixels into two classes $\mathcal{C}_0 = \{0, 1, \dots, t\}$ and $\mathcal{C}_1 = \{t + 1, t + 2, \dots, L - 1\}$ (object and background, or vice versa) at gray level t . In other words, pixels with gray levels less than or equal to the threshold t are assigned to \mathcal{C}_0 , whereas the remaining pixels are assigned to \mathcal{C}_1 .

Let n_i be the number of pixels with gray level i . The total number of pixels in the image is then given by $n = \sum_{i=0}^{L-1} n_i$. The normalized gray level histogram of the image can be regarded as a probability mass function:

$$p_i = \frac{n_i}{n}, \quad p_i \geq 0, \quad \sum_{i=0}^{L-1} p_i = 1.$$

Let $p_0(t) = \sum_{i=0}^t p_i$ and $p_1(t) = 1 - p_0(t)$ denote the probabilities of \mathcal{C}_0 and \mathcal{C}_1 , respectively. The means of the respective classes are then given by:

$$\begin{aligned} \mu_0(t) &= \mu(t)/p_0(t) \\ \mu_1(t) &= (\mu_T - \mu(t))/p_1(t) \end{aligned}$$

where $\mu(t) = \sum_{i=0}^t ip_i$ and $\mu_T = \mu(L - 1)$ denote the first moment of the histogram up to gray level t and mean gray level of the image, respectively.

Otsu’s method adopts between-class variance, i.e. $\sigma_B^2(t) = p_0(t)p_1(t)[\mu_0(t) - \mu_1(t)]^2$, from the discriminant analysis literature as its objective function and determines the optimal threshold t^* as the gray level that maximizes $\sigma_B^2(t)$, i.e. $t^* = \operatorname{argmax}_{t \in \{0, 1, \dots, L-1\}} \sigma_B^2(t)$. Between-class variance can be viewed as a measure of class separability or histogram bimodality. It can be efficiently calculated using:

$$\sigma_B^2(t) = \frac{[\mu_T p_0(t) - \mu(t)]^2}{p_0(t)p_1(t)}.$$

It should be noted that the efficiency of Otsu’s method can be attributed to the fact that it operates on histogrammed pixel gray values, which are non-negative integers. Var-Part and PCA-Part, on the other hand, operate on the projections of the points on the partitioning axis, which are often fractional. This problem can be

circumvented by linearly scaling the projection values to the limits of the histogram, i.e. 0 and $L - 1$. Let y_i be the projection of a point \mathbf{x}_i on the partitioning axis. y_i can be mapped to histogram bin b given by:

$$b = \left\lfloor \frac{L(y_i - \min_j y_j)}{\max_j y_j - \min_j y_j} \right\rfloor,$$

where $\lfloor z \rfloor$ is the floor function which returns the largest integer less than or equal to z .

The computational complexities of histogram construction and Otsu’s method are $\mathcal{O}(N_i)$ (N_i : number of points in the cluster) and $\mathcal{O}(L)$, respectively. L is constant in our experiments and therefore the proposed modification does not alter the linear time complexity of Var-Part and PCA-Part.

Figure 2 shows a histogram where using the mean point as a threshold leads to poor results. This histogram is constructed during the first iteration of PCA-Part from the projections of the points in the Shuttle data set (see Table 1). As marked on the figure, the mean point of this histogram is 61, whereas Otsu’s method gives a threshold of 105. The SSE of the initial cluster is 1836. When the mean point of the histogram is used a threshold, the resulting two subclusters have SSE’s of 408 and 809. This means that splitting the initial cluster with a hyperplane orthogonal to the principal eigenvector of the cluster covariance matrix at the mean point results in approximately 34% reduction in the SSE. On the other hand, when Otsu’s threshold is used, the subclusters have SSE’s of 943 and 101, which translates to about 43%

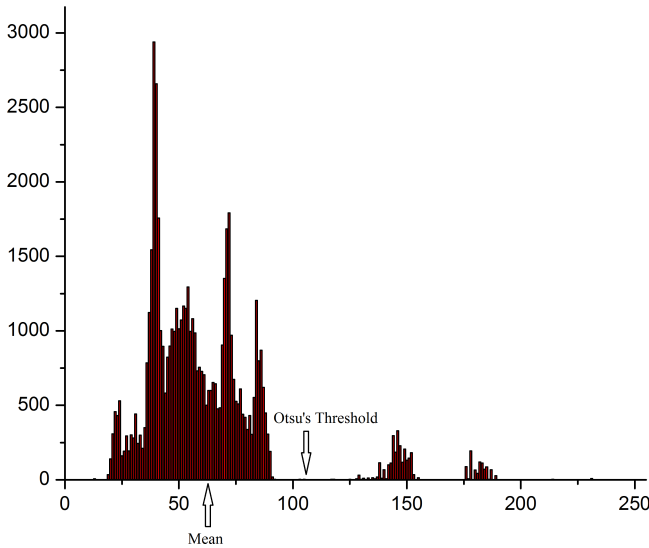


Fig. 2. Comparison of mean point and Otsu’s thresholds.

reduction in the SSE. In the next section, we will demonstrate that using Otsu's threshold instead of the mean point often leads to significantly better initial clusterings on a variety of data sets.

3. Experimental Results

The experiments were performed on 24 commonly used data sets from the UCI machine learning repository.⁵³ Table 1 gives the data set descriptions. For each data set, the number of clusters (K) was set equal to the number of classes (K'), as commonly seen in the related literature.^{23,39,40,29–32,54}

In clustering applications, normalization is a common preprocessing step that is necessary to prevent attributes with large ranges from dominating the distance calculations and also to avoid numerical instabilities in the computations. Two commonly used normalization schemes are linear scaling to unit range (min–max normalization) and linear scaling to unit variance (z-score normalization). Several studies revealed that the former scheme is preferable to the latter since the latter is likely to eliminate valuable between-cluster variation.^{55,40} As a result, we used min–max normalization to map the attributes of each data set to the $[0,1]$ interval.

Table 1. Data Set Descriptions (N : # points, D : # attributes, K' : # classes).

ID	Data Set	N	D	K'
1	Abalone	4177	7	28
2	Breast Cancer Wisconsin (Original)	683	9	2
3	Breast Tissue	106	9	6
4	Ecoli	336	7	8
5	Glass Identification	214	9	6
6	Heart Disease	297	13	5
7	Ionosphere	351	34	2
8	Iris (Bezdek)	150	4	3
9	ISOLET	7797	617	26
10	Landsat Satellite (Statlog)	6435	36	6
11	Letter Recognition	20,000	16	26
12	MAGIC Gamma Telescope	19,020	10	2
13	Multiple Features (Fourier)	2000	76	10
14	Musk (Clean2)	6598	166	2
15	Optical Digits	5620	64	10
16	Page Blocks Classification	5473	10	5
17	Pima Indians Diabetes	768	8	2
18	Shuttle (Statlog)	58,000	9	7
19	Spambase	4601	57	2
20	SPECTF Heart	267	44	2
21	Wall-Following Robot Navigation	5456	24	4
22	Wine Quality	6497	11	7
23	Wine	178	13	3
24	Yeast	1484	8	10

The performance of the IMs was quantified using two effectiveness (quality) and two efficiency (speed) criteria:

- **Initial SSE:** This is the SSE value calculated after the initialization phase, before the clustering phase. It gives us a measure of the effectiveness of an IM by itself.
- **Final SSE:** This is the SSE value calculated after the clustering phase. It gives us a measure of the effectiveness of an IM when its output is refined by k-means. Note that this is the objective function of the k-means algorithm, i.e. (1).
- **Number of Iterations:** This is the number of iterations that k-means requires until reaching convergence when initialized by a particular IM. It is an efficiency measure independent of programming language, implementation style, compiler, and CPU architecture.
- **CPU Time:** This is the total CPU time in milliseconds taken by the initialization and clustering phases.

All of the methods were implemented in the C language, compiled with the gcc v4.4.3 compiler, and executed on an Intel Xeon E5520 2.26 GHz machine. Time measurements were performed using the `getrusage` function, which is capable of measuring CPU time to an accuracy of a microsecond. The MT19937 variant of the Mersenne Twister algorithm was used to generate high quality pseudorandom numbers.⁵⁶

The convergence of k-means was controlled by the disjunction of two criteria: the number of iterations reaches a maximum of 100 or the relative improvement in SSE between two consecutive iterations drops below a threshold, i.e. $(SSE_{i-1} - SSE_i)/SSE_i \leq \varepsilon$, where SSE_i denotes the SSE value at the end of the i th ($i \in \{1, 2, \dots, 100\}$) iteration. The convergence threshold was set to $\varepsilon = 10^{-6}$.

In this study, we focus on IMs that have time complexity linear in N . This is because k-means itself has linear complexity, which is perhaps the most important reason for its popularity. Therefore, an IM for k-means should not diminish this advantage of the algorithm. The proposed methods, named Otsu Var-Part (OV) and Otsu PCA-Part (OP), were compared to six popular, linear, order-invariant IMs: Forgy's method (F), MacQueen's second method (M), maximin (X), k-means++ (K), Var-Part (V), and PCA-Part (P). It should be noted that among these methods F, M, and K are random, whereas X,^b V, P, OV, and OP are deterministic.

We first examine the influence of L (number of histogram bins) on the performance of OV and OP. Tables 2 and 3 show the initial and final SSE values obtained by respectively OV and OP for $L = 64, 128, 256, 512, 1024$ on four of the largest data sets (the best values are underlined). It can be seen that the performances of both methods are relatively insensitive to the value of L . Therefore, in the subsequent experiments we report the results for $L = 256$.

In the remaining experiments, each random method was executed a 100 times and statistics such as *minimum*, *mean*, and *standard deviation* were collected for

^bThe first center is chosen as the centroid of the data set.

Table 2. Influence of L on the effectiveness of Otsu Var-Part (L : # histogram bins).

ID	Criterion	$L = 64$	$L = 128$	$L = 256$	$L = 512$	$L = 1024$
9	Initial SSE	<u>143859</u>	144651	144658	144637	144638
	Final SSE	<u>118267</u>	119127	<u>118033</u>	118033	118034
10	Initial SSE	1987	1920	<u>1919</u>	1920	1920
	Final SSE	<u>1742</u>	<u>1742</u>	<u>1742</u>	<u>1742</u>	<u>1742</u>
11	Initial SSE	3242	<u>3192</u>	3231	3202	3202
	Final SSE	2742	<u>2734</u>	<u>2734</u>	<u>2734</u>	<u>2734</u>
15	Initial SSE	<u>17448</u>	17504	17504	17504	17504
	Final SSE	<u>14581</u>	<u>14581</u>	<u>14581</u>	<u>14581</u>	<u>14581</u>

Table 3. Influence of L on the effectiveness of Otsu PCA-Part (L : # histogram bins).

ID	Criterion	$L = 64$	$L = 128$	$L = 256$	$L = 512$	$L = 1024$
9	Initial SSE	123527	123095	<u>122528</u>	123129	123342
	Final SSE	118575	118577	<u>119326</u>	<u>118298</u>	118616
10	Initial SSE	1855	<u>1807</u>	1835	1849	1848
	Final SSE	<u>1742</u>	<u>1742</u>	<u>1742</u>	<u>1742</u>	<u>1742</u>
11	Initial SSE	2994	2997	2995	2995	<u>2991</u>
	Final SSE	<u>2747</u>	<u>2747</u>	<u>2747</u>	<u>2747</u>	<u>2747</u>
15	Initial SSE	15136	15117	15118	<u>15116</u>	15117
	Final SSE	<u>14650</u>	<u>14650</u>	<u>14650</u>	<u>14650</u>	<u>14650</u>

each performance criteria. The *minimum* and *mean* statistics represent the *best* and *average* case performance, respectively, while *standard deviation* quantifies the *variability* of performance across different runs. Note that for a deterministic method, the *minimum* and *mean* values are always identical and the *standard deviation* is always 0.

Tables 4–7 give the performance measurements for each method with respect to initial SSE, final SSE, number of iterations, and CPU time, respectively. For each of the initial SSE, final SSE, and number of iterations criteria, we calculated the ratio of the *minimum/mean/standard deviation* value obtained by each method to the best (least) *minimum/mean/standard deviation* value on each data set. These ratios will be henceforth referred to as the “normalized” performance criteria. For example, on the Abalone data set (1), the *minimum* initial SSE of F is 424.92 and the best *minimum* initial SSE is 21.57 and thus the normalized initial SSE is about 20. This simply means that on this data set F obtained approximately 20 times worse *minimum* initial SSE than the best method. We then averaged^c these normalized values over the 24 data sets to quantify the *overall* performance of each method with respect

^cDue to outliers, the “median” statistic rather than the “mean” was used to summarize the normalized standard deviation values.

Table 4. Initial SSE comparison of the initialization methods.

	F	M	K	X	V	P	OV	OP
1	min	425	33	29	95	24	23	22
	mean	483 ± 20	46 ± 10	34 ± 2	95 ± 0	24 ± 0	23 ± 0	22 ± 0
2	min	534	318	304	498	247	240	239
	mean	575 ± 15	706 ± 354	560 ± 349	498 ± 0	247 ± 0	240 ± 0	239 ± 0
3	min	20	11	9	19	8	8	7
	mean	27 ± 3	20 ± 8	13 ± 2	19 ± 0	8 ± 0	8 ± 0	7 ± 0
4	min	54	26	26	48	20	19	20
	mean	61 ± 2	40 ± 7	33 ± 5	48 ± 0	20 ± 0	19 ± 0	20 ± 0
5	min	42	24	25	45	21	20	18
	mean	48 ± 2	40 ± 9	32 ± 5	45 ± 0	21 ± 0	20 ± 0	18 ± 0
6	min	372	361	341	409	249	249	244
	mean	396 ± 8	463 ± 58	450 ± 49	409 ± 0	249 ± 0	249 ± 0	244 ± 0
7	min	771	749	720	827	632	629	629
	mean	814 ± 12	1246 ± 463	1237 ± 468	827 ± 0	632 ± 0	629 ± 0	629 ± 0
8	min	26	9	9	18	8	8	7
	mean	34 ± 4	28 ± 23	16 ± 6	18 ± 0	8 ± 0	8 ± 0	7 ± 0
9	min	218965	212238	210387	221163	145444	124958	122528
	mean	223003 ± 1406	224579 ± 5416	223177 ± 4953	221163 ± 0	145444 ± 0	124958 ± 0	122528 ± 0
10	min	7763	2637	2458	4816	2050	2116	1835
	mean	8057 ± 98	4825 ± 1432	3561 ± 747	4816 ± 0	2050 ± 0	2116 ± 0	1835 ± 0
11	min	7100	4203	4158	5632	3456	3101	2995
	mean	7225 ± 30	4532 ± 165	4501 ± 176	5632 ± 0	3456 ± 0	3101 ± 0	2995 ± 0
12	min	4343	3348	3296	4361	3056	2927	2923
	mean	4392 ± 13	5525 ± 1816	5346 ± 1672	4361 ± 0	3056 ± 0	2927 ± 0	2923 ± 0
13	min	4416	5205	5247	4485	3354	3266	3180
	mean	4475 ± 25	5693 ± 315	5758 ± 283	4485 ± 0	3354 ± 0	3266 ± 0	3180 ± 0
14	min	53508	56841	56822	54629	37334	37142	36375
	mean	54312 ± 244	82411 ± 14943	75532 ± 12276	54629 ± 0	37334 ± 0	37142 ± 0	36375 ± 0

Table 4. (Continued)

	F	M	K	X	V	P	OV	OP
15	min	25466	25492	24404	25291	17476	15714	15118
	mean	25811 ± 99	28596 ± 1550	27614 ± 1499	25291 ± 0	17476 ± 0	15714 ± 0	15118 ± 0
16	min	633	275	250	635	300	230	222
	mean	648 ± 6	423 ± 74	372 ± 72	635 ± 0	300 ± 0	230 ± 0	222 ± 0
17	min	152	144	141	156	124	122	121
	mean	156 ± 1	216 ± 44	219 ± 61	156 ± 0	124 ± 0	122 ± 0	121 ± 0
18	min	1788	438	328	1818	316	309	268
	mean	1806 ± 6	946 ± 290	494 ± 115	1818 ± 0	316 ± 0	309 ± 0	268 ± 0
19	min	834	873	881	772	782	783	765
	mean	838 ± 1	1186 ± 386	1124 ± 244	772 ± 0	782 ± 0	783 ± 0	765 ± 0
20	min	269	295	297	277	232	222	214
	mean	281 ± 4	384 ± 88	413 ± 159	277 ± 0	232 ± 0	222 ± 0	214 ± 0
21	min	10976	11834	11829	11004	8517	7805	7802
	mean	11082 ± 34	14814 ± 1496	14435 ± 1276	11004 ± 0	8517 ± 0	7805 ± 0	7802 ± 0
22	min	719	473	449	733	386	361	351
	mean	729 ± 4	601 ± 59	567 ± 64	733 ± 0	386 ± 0	361 ± 0	351 ± 0
23	min	78	76	70	87	51	53	51
	mean	87 ± 3	113 ± 22	101 ± 20	87 ± 0	51 ± 0	53 ± 0	51 ± 0
24	min	144	89	83	115	77	63	63
	mean	149 ± 2	110 ± 8	101 ± 9	115 ± 0	77 ± 0	63 ± 0	63 ± 0

Table 5. Final SSE comparison of the initialization methods.

	F	M	K	X	V	P	OV	OP
1	min 21 23 ± 1	22 22 ± 1	21 22 ± 0	25 25 ± 0	21 21 ± 0	21 21 ± 0	21 21 ± 0	21 21 ± 0
2	min 239 239 ± 0	239 239 ± 0	239 239 ± 0	239 239 ± 0	239 239 ± 0	239 239 ± 0	239 239 ± 0	239 239 ± 0
3	min 7 8 ± 1	7 9 ± 1	7 8 ± 1	7 7 ± 0	7 7 ± 0	7 7 ± 0	8 8 ± 0	7 7 ± 0
4	min 17 19 ± 1	17 19 ± 2	17 19 ± 1	19 19 ± 0	17 17 ± 0	18 18 ± 0	18 18 ± 0	18 18 ± 0
5	min 18 20 ± 1	18 21 ± 2	18 20 ± 2	23 23 ± 0	19 19 ± 0	19 19 ± 0	20 20 ± 0	18 18 ± 0
6	min 243 252 ± 8	243 252 ± 8	243 252 ± 8	249 249 ± 0	248 248 ± 0	243 243 ± 0	248 248 ± 0	243 243 ± 0
7	min 629 629 ± 0	629 643 ± 50	629 641 ± 47	826 826 ± 0	629 629 ± 0	629 629 ± 0	629 629 ± 0	629 629 ± 0
8	min 7 8 ± 1	7 8 ± 2	7 7 ± 1	7 7 ± 0	7 7 ± 0	7 7 ± 0	7 7 ± 0	7 7 ± 0
9	min 117872 119650 ± 945	117764 119625 ± 947	117710 119536 ± 934	135818 135818 ± 0	118495 118495 ± 0	118386 118386 ± 0	118033 118033 ± 0	119326 119326 ± 0
10	min 1742 1742 ± 0	1742 1742 ± 0	1742 1744 ± 28	1742 1742 ± 0	1742 1742 ± 0	1742 1742 ± 0	1742 1742 ± 0	1742 1742 ± 0
11	min 2723 2772 ± 29	2718 2757 ± 19	2716 2751 ± 19	2749 2749 ± 0	2735 2735 ± 0	2745 2745 ± 0	2734 2734 ± 0	2747 2747 ± 0
12	min 2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0	2923 2923 ± 0
13	min 3127 3166 ± 31	3128 3172 ± 29	3128 3173 ± 35	3316 3316 ± 0	3137 3137 ± 0	3214 3214 ± 0	3143 3143 ± 0	3153 3153 ± 0
14	min 36373 37296 ± 1902	36373 37163 ± 1338	36373 37058 ± 1626	36373 36373 ± 0	36373 36373 ± 0	36373 36373 ± 0	36373 36373 ± 0	36373 36373 ± 0

Table 5. (Continued)

	F	M	K	X	V	P	OV	OP
15	min	14559	14559	14679	14581	14807	14581	14650
	mean	14687 ± 216	14752 ± 236	14679 ± 0	14581 ± 0	14807 ± 0	14581 ± 0	14650 ± 0
16	min	215	215	230	227	215	229	216
	mean	217 ± 4	216 ± 2	230 ± 0	227 ± 0	215 ± 0	229 ± 0	216 ± 0
17	min	121	121	121	121	121	121	121
	mean	121 ± 0	122 ± 5	121 ± 0	121 ± 0	121 ± 0	121 ± 0	121 ± 0
18	min	235	235	726	235	274	274	235
	mean	317 ± 46	272 ± 23	726 ± 0	235 ± 0	274 ± 0	274 ± 0	235 ± 0
19	min	765	765	765	778	778	778	765
	mean	778 ± 3	779 ± 14	765 ± 0	778 ± 0	778 ± 0	778 ± 0	765 ± 0
20	min	214	214	214	214	214	214	214
	mean	214 ± 0	215 ± 5	214 ± 0	214 ± 0	214 ± 0	214 ± 0	214 ± 0
21	min	7772	7772	7772	7774	7774	7774	7772
	mean	7799 ± 93	7821 ± 124	7772 ± 0	7774 ± 0	7774 ± 0	7774 ± 0	7772 ± 0
22	min	334	334	399	335	334	335	335
	mean	335 ± 1	336 ± 3	399 ± 0	335 ± 0	334 ± 0	335 ± 0	335 ± 0
23	min	49	49	63	49	49	49	49
	mean	49 ± 0	49 ± 2	63 ± 0	49 ± 0	49 ± 0	49 ± 0	49 ± 0
24	min	58	58	61	69	59	69	59
	mean	64 ± 6	69 ± 6	61 ± 0	69 ± 0	59 ± 0	69 ± 0	59 ± 0

Table 6. Number of iterations comparison of the initialization methods.

		F	M	K	X	V	P	OV	OP
1	min	59	29	<u>22</u>	100	50	43	31	38
	mean	90 ± 11	68 ± 19	48 ± <u>17</u>	100 ± 0	50 ± 0	43 ± 0	<u>31</u> ± 0	38 ± 0
2	min	4	4	4	8	4	4	5	<u>3</u>
	mean	5 ± 0	6 ± 1	6 ± 1	8 ± 0	4 ± 0	4 ± 0	5 ± 0	<u>3</u> ± 0
3	min	5	5	<u>3</u>	7	6	7	5	<u>3</u>
	mean	10 ± 2	9 ± 3	7 ± <u>2</u>	7 ± 0	6 ± 0	7 ± 0	5 ± 0	<u>3</u> ± 0
4	min	8	<u>6</u>	7	14	17	7	12	<u>6</u>
	mean	15 ± 6	15 ± <u>5</u>	14 ± 5	14 ± 0	17 ± 0	7 ± 0	12 ± 0	<u>6</u> ± 0
5	min	6	5	<u>4</u>	6	6	5	9	<u>4</u>
	mean	10 ± 3	11 ± 4	9 ± <u>3</u>	6 ± 0	6 ± 0	5 ± 0	9 ± 0	<u>4</u> ± 0
6	min	5	5	5	12	<u>3</u>	4	<u>3</u>	4
	mean	11 ± 3	10 ± 3	9 ± 3	12 ± 0	<u>3</u> ± 0	4 ± 0	<u>3</u> ± 0	4 ± 0
7	min	4	3	3	3	3	3	4	<u>2</u>
	mean	5 ± 1	7 ± 2	8 ± 2	3 ± 0	3 ± 0	3 ± 0	4 ± 0	<u>2</u> ± 0
8	min	4	4	<u>3</u>	6	4	4	6	<u>3</u>
	mean	9 ± 3	8 ± 2	7 ± <u>3</u>	6 ± 0	4 ± 0	4 ± 0	6 ± 0	<u>3</u> ± 0
9	min	18	19	<u>14</u>	32	82	45	59	39
	mean	43 ± 15	40 ± 14	36 ± <u>13</u>	<u>32</u> ± 0	82 ± 0	45 ± 0	59 ± 0	39 ± 0
10	min	12	12	11	53	28	27	<u>10</u>	23
	mean	28 ± 8	33 ± 10	29 ± 9	53 ± 0	28 ± 0	27 ± 0	<u>10</u> ± 0	23 ± 0
11	min	39	37	<u>31</u>	73	100	83	67	85
	mean	75 ± 19	72 ± 18	76 ± <u>18</u>	73 ± 0	100 ± 0	83 ± 0	<u>67</u> ± 0	85 ± 0
12	min	<u>9</u>	10	10	35	25	10	26	<u>9</u>
	mean	18 ± <u>5</u>	18 ± 5	20 ± 6	35 ± 0	25 ± 0	10 ± 0	26 ± 0	<u>9</u> ± 0
13	min	<u>13</u>	14	<u>13</u>	37	14	25	17	<u>13</u>
	mean	29 ± <u>10</u>	30 ± 10	30 ± <u>11</u>	37 ± 0	14 ± 0	25 ± 0	17 ± 0	<u>13</u> ± 0
14	min	4	4	4	8	5	5	5	<u>3</u>
	mean	6 ± 1	6 ± 1	6 ± 1	8 ± 0	5 ± 0	5 ± 0	5 ± 0	<u>3</u> ± 0
15	min	<u>12</u>	<u>12</u>	14	36	16	22	15	59
	mean	31 ± <u>13</u>	33 ± <u>14</u>	30 ± 10	36 ± 0	16 ± 0	22 ± 0	<u>15</u> ± 0	59 ± 0
16	min	14	12	<u>9</u>	27	25	15	19	16
	mean	27 ± 9	31 ± 14	24 ± <u>11</u>	27 ± 0	25 ± 0	15 ± 0	19 ± 0	<u>16</u> ± 0
17	min	8	<u>4</u>	<u>4</u>	19	11	10	8	5
	mean	13 ± 2	12 ± <u>5</u>	11 ± <u>4</u>	19 ± 0	11 ± 0	10 ± 0	8 ± 0	<u>5</u> ± 0
18	min	10	8	9	22	30	16	<u>7</u>	27
	mean	25 ± 9	25 ± 11	23 ± 9	22 ± 0	30 ± 0	16 ± 0	<u>7</u> ± 0	27 ± 0
19	min	6	<u>3</u>	<u>3</u>	5	9	10	12	<u>3</u>
	mean	12 ± 5	14 ± <u>6</u>	12 ± <u>7</u>	5 ± 0	9 ± 0	10 ± 0	12 ± 0	<u>3</u> ± 0
20	min	6	5	4	7	7	7	6	<u>2</u>
	mean	8 ± 1	8 ± 2	7 ± 2	7 ± 0	7 ± 0	7 ± 0	6 ± 0	<u>2</u> ± 0
21	min	9	11	11	24	20	<u>8</u>	21	19
	mean	20 ± 8	22 ± 8	20 ± 8	24 ± 0	20 ± 0	<u>8</u> ± 0	21 ± 0	19 ± 0
22	min	<u>15</u>	17	18	20	62	50	33	49
	mean	41 ± <u>20</u>	42 ± 18	40 ± 18	<u>20</u> ± 0	62 ± 0	50 ± 0	33 ± 0	49 ± 0
23	min	<u>4</u>	<u>4</u>	<u>4</u>	9	5	7	5	5
	mean	7 ± <u>2</u>	8 ± <u>3</u>	7 ± <u>3</u>	9 ± 0	<u>5</u> ± 0	7 ± 0	<u>5</u> ± 0	<u>5</u> ± 0
24	min	<u>13</u>	<u>13</u>	15	73	33	21	28	32
	mean	29 ± <u>10</u>	31 ± <u>11</u>	29 ± 10	73 ± 0	33 ± 0	<u>21</u> ± 0	28 ± 0	32 ± 0

Table 7. CPU time comparison of the initialization methods.

	F	M	K	X	V	P	OV	OP
1	min 43 ± 7	10 33 ± 10	10 24 ± 9	40 40 ± 0	20 20 ± 0	20 20 ± 0	20 20 ± 0	10 10 ± 0
2	min 0 ± 1	0 0 ± 1	0 0 ± 2	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
3	min 0 ± 1	0 0 ± 1	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
4	min 0 ± 2	0 0 ± 1	0 0 ± 2	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
5	min 0 ± 1	0 0 ± 1	0 0 ± 1	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
6	min 0 ± 1	0 0 ± 2	0 1 ± 2	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
7	min 0 ± 2	0 0 ± 2	0 0 ± 2	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
8	min 0 ± 0	0 0 ± 1	0 0 ± 1	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0	0 0 ± 0
9	min 3691 ± 1229	1630 3370 ± 1178	1580 3397 ± 1055	2570 2570 ± 0	6920 6920 ± 0	12160 12160 ± 0	5040 5040 ± 0	12460 12460 ± 0
10	min 30 ± 9	10 32 ± 11	10 32 ± 10	50 50 ± 0	30 30 ± 0	50 50 ± 0	10 10 ± 0	50 50 ± 0
11	min 710 ± 174	380 673 ± 168	320 724 ± 166	670 670 ± 0	960 960 ± 0	790 790 ± 0	620 620 ± 0	810 810 ± 0
12	min 19 ± 7	10 19 ± 7	10 21 ± 6	40 40 ± 0	20 20 ± 0	10 10 ± 0	30 30 ± 0	20 20 ± 0
13	min 52 ± 17	20 52 ± 17	20 57 ± 20	60 60 ± 0	20 20 ± 0	80 80 ± 0	30 30 ± 0	40 40 ± 0
14	min 22 ± 7	10 21 ± 5	10 25 ± 8	30 30 ± 0	30 30 ± 0	210 210 ± 0	20 20 ± 0	200 200 ± 0

Table 7. (Continued)

	F	M	K	X	V	P	OV	OP
15	min	50	50	140	60	140	70	280
	mean	122 ± 50	126 ± 53	140 ± 0	60 ± 0	140 ± 0	70 ± 0	280 ± 0
16	min	0	0	10	10	10	10	10
	mean	9 ± 5	11 ± 6	10 ± 0	10 ± 0	10 ± 0	10 ± 0	10 ± 0
17	min	0	0	0	0	0	0	0
	mean	0 ± 2	1 ± 2	0 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0
18	min	40	30	50	100	70	30	100
	mean	87 ± 27	87 ± 37	50 ± 0	100 ± 0	70 ± 0	30 ± 0	100 ± 0
19	min	0	0	10	10	30	10	20
	mean	11 ± 5	12 ± 6	10 ± 0	10 ± 0	30 ± 0	10 ± 0	20 ± 0
20	min	0	0	0	0	0	0	0
	mean	0 ± 2	0 ± 2	0 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0
21	min	10	10	20	20	20	20	30
	mean	21 ± 9	20 ± 8	20 ± 0	20 ± 0	20 ± 0	20 ± 0	30 ± 0
22	min	10	20	10	60	50	40	40
	mean	40 ± 19	40 ± 18	10 ± 0	60 ± 0	50 ± 0	40 ± 0	40 ± 0
23	min	0	0	0	0	0	0	0
	mean	0 ± 1	0 ± 1	0 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0
24	min	0	0	10	0	0	10	0
	mean	6 ± 5	6 ± 5	10 ± 0	0 ± 0	0 ± 0	10 ± 0	0 ± 0

Table 8. Overall performance comparison of the initialization methods.

Statistic	Criterion	F	M	K	X	V	P	OV	OP
Min	Init. SSE	2.968	1.418	1.348	2.184	1.107	1.043	1.067	<u>1.002</u>
	Final SSE	1.001	1.003	<u>1.000</u>	1.163	1.019	1.018	1.031	<u>1.005</u>
	# Iters.	1.488	1.284	<u>1.183</u>	2.978	2.469	2.013	2.034	1.793
Mean	Init. SSE	3.250	2.171	1.831	2.184	1.107	1.043	1.067	<u>1.002</u>
	Final SSE	1.047	1.049	1.032	1.161	1.017	1.016	1.029	<u>1.003</u>
	# Iters.	2.466	2.528	2.314	2.581	2.057	1.715	1.740	<u>1.481</u>
Stdev	Init. SSE	<u>1.000</u>	13.392	12.093	—	—	—	—	—
	Final SSE	<u>1.013</u>	1.239	1.172	—	—	—	—	—
	# Iters.	<u>1.000</u>	1.166	1.098	—	—	—	—	—

to each statistic (see Table 8). Note that we did not attempt to summarize the CPU time values in the same manner due to the sparsity of the data (see Table 7). For convenient visualization, Fig. 3 shows box plots of the normalized performance criteria. Here, the bottom and top end of the whiskers of a box represent the *minimum*

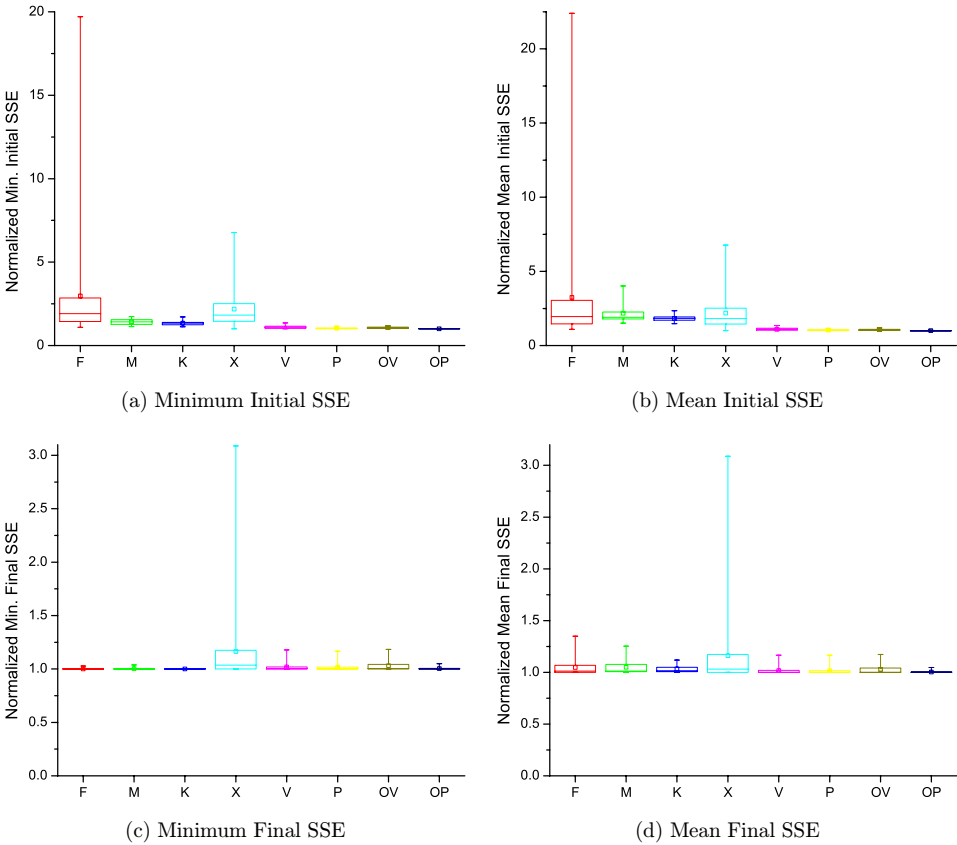


Fig. 3. Box plots of the normalized performance criteria.

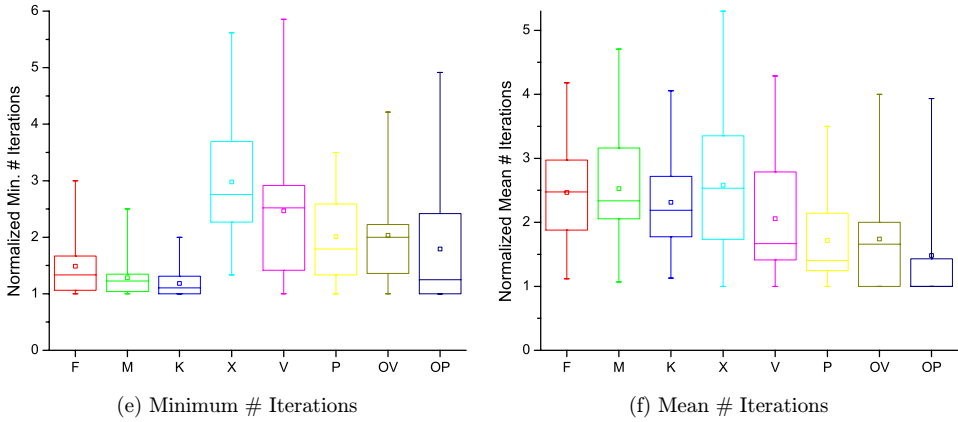


Fig. 3. (Continued)

and *maximum*, respectively, whereas the bottom and top of the box itself are the 25th percentile ($Q1$) and 75th percentile ($Q3$), respectively. The line that passes through the box is the 50th percentile ($Q2$), i.e. the *median*, while the small square inside the box denotes the *mean*.

4. Discussion

4.1. Best case performance analysis

With respect to the *minimum* statistic, the following observations can be made:

- **Initial SSE:** OP is the best method, followed closely by P, OV, and V. On the other hand, F is the worst method, followed by X. These two methods give 2–3 times worse *minimum* initial SSE than the best method. It can be seen that multiple runs of random methods do not produce good initial clusterings. In contrast, only a single run of OP, P, OV, or V often gives very good results. This is because these methods are approximate clustering algorithms by themselves and thus they give reasonable results even without k-means refinement.
- **Final SSE:** X is the worst method, while the remaining methods exhibit very similar performance. This homogeneity in performance is because k-means can take two disparate initial configurations to similar (or even identical) local minima. Given the abundance of local minima even in data sets of moderate size and/or dimensionality and the gradient descent nature of k-means, it is not surprising that the deterministic methods (except X) perform slightly worse than the random methods as the former methods were executed only once, whereas the latter ones were executed a 100 times.
- **Number of Iterations:** K is the best method, followed by M and F. X is the worst method. As in the case of final SSE, random methods outperform deterministic methods due to their multiple-run advantage.

4.2. Average case performance analysis

With respect to the *mean* statistic, the following observations can be made:

- **Initial SSE:** OP is the best method, followed closely by P, OV, and V. The remaining methods give 1.7–3.2 times worse *mean* initial SSE than any of these hierarchical methods. Random methods exhibit significantly worse *average* performance than the deterministic ones because the former methods can produce highly variable results across different runs (see the standard deviation values in Table 4).
- **Final SSE:** This is similar to the case of *minimum* final SSE, with the difference that deterministic methods (except X) are now slightly better than the random ones. Once again, this is because random methods can produce highly variable results due to their stochastic nature.
- **Number of Iterations:** The ranking of the methods is similar to the case of *mean* final SSE.

4.3. Consistency analysis

With respect to the *standard deviation* statistic, F is significantly better than both M and K. If, however, the application requires absolute consistency, i.e. exactly the same clustering in every run, a deterministic IM should be used.

4.4. CPU time analysis

It can be seen from Table 7 that, on about half of the data sets, each of the IMs require less than a few milliseconds of CPU time. On the other hand, on large and/or high-dimensional data sets efficiency differences become more prominent. It should be noted that each of the values reported in this table corresponds to a single k-means “run”. In practice, a random method is typically executed R times, e.g. in this study $R = 100$, and the output of the run that gives the least final SSE is taken as the result. Therefore, the total computational cost of a random method is often significantly higher than that of a deterministic method. For example, on the ISOLET data set, which has the greatest $N \times D \times K$ value among all the data sets, K took on the average 3397 ms, whereas OP took 12460 ms. The latter method, however, required about 27 times less CPU time than the former one since the former was executed a total of 100 times.

4.5. Relative performance analysis

We also determined the number of data sets (out of 24) on which OV and OP respectively performed {worse than/same as/better than}V and P. Tables 9 and 10 present the results for OV and OP, respectively. It can be seen that, with respect to initial SSE and number of iterations criteria, OV outperforms V more often than not. On the other hand, OP frequently outperforms P with respect to both criteria. As for final SSE, OP performs slightly better than P, whereas OV performs slightly

Table 9. Performance of Otsu Var-Part relative to Var-Part (V).

Criterion	Worse than V	Same as V	Better than V
Init. SSE	6	3	15
Final SSE	6	16	2
# Iters.	8	3	13

Table 10. Performance of Otsu PCA-Part relative to PCA-Part (P).

Criterion	Worse than P	Same as P	Better than P
Init. SSE	1	2	21
Final SSE	4	14	6
# Iters.	6	1	17

worse than V. It appears that Otsu’s method benefits P more than it benefits V. This is most likely due to the fact that histograms of projections over the major axis necessarily have a greater dynamic range and variability and thus are more amenable to thresholding compared to histograms of projections over any coordinate axis.

4.6. Recommendations for practitioners

Based on the analyses presented above, the following recommendations can be made:

- In general, X should not be used. As mentioned in Sec. 2, this method was not designed specifically as a k-means initializer.⁴³ It is easy to understand and implement, but is mostly ineffective and unreliable. Furthermore, despite its low overhead, this method does not offer significant time savings since it often results in slow k-means convergence.
- In applications that involve small data sets, e.g. $N < 1000$, K should be used. It is computationally feasible to run this method hundreds of times on such data sets given that one such run takes only a few milliseconds.
- In time-critical applications that involve large data sets or applications that demand determinism, the hierarchical methods should be used. These methods need to be executed only once and they lead to reasonably fast k-means convergence. The efficiency difference between V/OV and P/OP is noticeable only on high-dimensional data sets such as ISOLET ($D = 617$) and Musk ($D = 166$). This is because V/OV calculates the direction of split by determining the coordinate axis with the greatest variance (in $\mathcal{O}(D)$ time), whereas P/OP achieves this by calculating the principal eigenvector of the cluster covariance matrix (in $\mathcal{O}(D^2)$ time using the power method⁴⁴). Note that despite its higher computational complexity, P/OP can, in some cases, be more efficient than V/OV (see Table 7). This is because the former converges significantly faster than the latter (see Table 8). The

main disadvantage of these methods is that they are more complicated to implement due to their hierarchical formulation.

- In applications where an approximate clustering of the data set is desired, the hierarchical methods should be used. These methods produce very good initial clusterings, which makes it possible to use them as standalone clustering algorithms.
- Among the hierarchical methods, the ones based on PCA, i.e. P and OP, are preferable to those based on variance, i.e. V and OV. Furthermore, the proposed OP and OV methods generally outperform their respective counterparts, i.e. P and V, especially with respect to initial SSE and number of iterations.

5. Conclusion

In this paper, we presented a simple modification to Var-Part and PCA-Part, two hierarchical k-means initialization methods that are linear, deterministic, and order-invariant. We compared the original methods and their modified versions to some of the most popular linear initialization methods, namely Forgy's method, Macqueen's second method, maximin, and k-means++, on a large and diverse collection of data sets from the UCI machine learning repository. The results demonstrated that, despite their deterministic nature, Var-Part and PCA-Part are highly competitive with one of the best random initialization methods to date, i.e. k-means++. In addition, the proposed modification significantly improves the performance of both hierarchical methods. The presented Var-Part and PCA-Part variants can be used to initialize k-means effectively, particularly in time-critical applications that involve large data sets. Alternatively, they can be used as approximate clustering algorithms without additional k-means refinement.

Acknowledgments

This publication was made possible by grants from the Louisiana Board of Regents (LEQSF2008-11-RD-A-12) and National Science Foundation (0959583, 1117457).

References

1. A. K. Jain, M. N. Murty and P. J. Flynn, Data clustering: A review, *ACM Comput. Surveys* **31**(3) (1999) 264–323.
2. A. K. Jain, Data clustering: 50 years beyond k-means, *Pattern Recog. Lett.* **31**(8) (2010) 651–666.
3. L. Kaufman and P. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis* (Wiley-Interscience, 1990).
4. D. Aloise, A. Deshpande, P. Hansen and P. Popat, NP-hardness of Euclidean sum-of-squares clustering, *Machine Learn.* **75**(2) (2009) 245–248.
5. M. Mahajan, P. Nimbhorkar and K. Varadarajan, The planar k-means problem is NP-hard, *Theor. Comput. Sci.* **442** (2012) 13–21.

6. A. Tarsitano, A computational study of several relocation methods for k-means algorithms, *Pattern Recog.* **36**(12) (2003) 2955–2966.
7. S. Lloyd, Least squares quantization in PCM, *IEEE Trans. Inf. Theor.* **28**(2) (1982) 129–136.
8. P. S. Bradley and U. Fayyad, Refining initial points for k-means clustering, in *Proc. 15th Int. Conf. Machine Learning*, 1998, pp. 91–99.
9. M. Dash, H. Liu and X. Xu, ‘1 + 1 > 2’: Merging distance and density based clustering, in *Proc. 7th Int. Conf. Database Systems for Advanced Applications*, 2001, pp. 32–39.
10. W. Y. Chen, Y. Song, H. Bai, C. J. Lin and E. Y. Chang, Parallel spectral clustering in distributed systems, *IEEE Trans. Pattern Anal. Machine Intell.* **33**(3) (2011) 568–586.
11. T. Kanungo, D. Mount, N. Netanyahu, C. Piatko, R. Silverman and A. Wu, An efficient k-means clustering algorithm: Analysis and implementation, *IEEE Trans. Pattern Anal. Machine Intell.* **24**(7) (2002) 881–892.
12. G. Hamerly, Making k-means even faster, in *Proc. 2010 SIAM Int. Conf. Data Mining*, 2010, pp. 130–140.
13. T. W. Chen and S. Y. Chien, Bandwidth adaptive hardware architecture of k-means clustering for video analysis, *IEEE Trans. Very Large Scale Integration (VLSI) Syst.* **18**(6) (2010) 957–966.
14. C. Ordonez and E. Omiecinski, Efficient disk-based k-means clustering for relational databases, *IEEE Trans. Knowledge and Data Eng.* **16**(8) (2004) 909–921.
15. S. Z. Selim and M. A. Ismail, K-mans-type algorithms: A generalized convergence theorem and characterization of local optimality, *IEEE Trans. Pattern Anal. Machine Intell.* **6**(1) (1984) 81–87.
16. L. Bottou and Y. Bengio, Convergence properties of the k-means algorithms, in *Advances in Neural Information Processing Systems 7*, eds. G. Tesauro, D. S. Touretzky and T. K. Leen (MIT Press 1995), pp. 585–592.
17. L. Vendramin, R. J. G. B. Campello and E. R. Hruschka, Relative clustering validity criteria: A comparative overview, *Stat. Anal. Data Mining* **3**(4) (2010) 209–235.
18. J. Mao and A. K. Jain, A self-organizing network for hyperellipsoidal clustering (HEC), *IEEE Trans. Neural Netw.* **7**(1) (1996) 16–29.
19. J. S. Zhang and Y.-W. Leung, Robust clustering by pruning outliers, *IEEE Trans. Syst. Man Cybern. B* **33**(6) (2003) 983–999.
20. M. E. Celebi, H. Kingravi and P. A. Vela, A comparative study of efficient initialization methods for the k-means clustering algorithm, *Expert Syst. Appl.* **40**(1) (2013) 200–210.
21. M. E. Celebi, Improving the performance of k-means for color quantization, *Image Vis. Comput.* **29**(4) (2011) 260–271.
22. J. M. Pena, J. A. Lozano and P. Larranaga, An empirical comparison of four initialization methods for the k-means algorithm, *Pattern Recogn. Lett.* **20**(10) (1999) 1027–1040.
23. J. He, M. Lan, C. L. Tan, S. Y. Sung and H. B. Low, Initialization of cluster refinement algorithms: A review and comparative study, in *Proc. 2004 IEEE Int. Joint Conf. Neural Networks*, 2004, pp. 297–302.
24. G. N. Lance and W. T. Williams, A general theory of classificatory sorting strategies — II. Clustering systems, *Computer J.* **10**(3) (1967) 271–277.
25. M. M. Astrahan, Speech analysis by clustering, or the hyperphoneme method, Technical Report AIM-124, Stanford University, 1970.
26. J. A. Hartigan and M. A. Wong, Algorithm AS 136: A k-means clustering algorithm, *J. R. Stat. Soc. C* **28**(1) (1979) 100–108.
27. A. Likas, N. Vlassis and J. Verbeek, The global k-means clustering algorithm, *Pattern Recog.* **36**(2) (2003) 451–461.

28. M. Al-Daoud, A new algorithm for cluster initialization, in *Proc. 2nd World Enformatika Conf.*, 2005, pp. 74–76.
29. S. J. Redmond and C. Heneghan, A method for initialising the k-means clustering algorithm using kd-trees, *Pattern Recog. Lett.* **28**(8) (2007) 965–973.
30. M. Al Hasan, V. Chaoji, S. Salem and M. Zaki, Robust partitional clustering by outlier and density insensitive seeding, *Pattern Recog. Lett.* **30**(11) (2009) 994–1002.
31. F. Cao, J. Liang and G. Jiang, An initialization method for the k-means algorithm using neighborhood model, *Comput. Math. Appl.* **58**(3) (2009) 474–483.
32. P. Kang and S. Cho, K-means clustering seeds initialization based on centrality, sparsity, and isotropy, in *Proc. 10th Int. Conf. Intelligent Data Engineering and Automated Learning*, 2009, pp. 109–117.
33. E. Forgy, Cluster analysis of multivariate data: Efficiency vs. interpretability of classification, *Biometrics* **21** (1965) 768.
34. R. C. Jancey, Multidimensional group analysis, *Australian J. Botany* **14**(1) (1966) 127–130.
35. J. MacQueen, Some methods for classification and analysis of multivariate observations, in *Proc. 5th Berkeley Symp. Mathematical Statistics and Probability*, 1967, pp. 281–297.
36. G. H. Ball and D. J. Hall, A clustering technique for summarizing multivariate data, *Behav. Sci.* **12**(2) (1967) 153–155.
37. J. T. Tou and R. C. Gonzales, *Pattern Recognition Principles* (Addison-Wesley, 1974).
38. H. Späth, Computational experiences with the exchange method: Applied to four commonly used partitioning cluster analysis criteria, *Eur. J. Oper. Res.* **1**(1) (1977) 23–31.
39. D. Arthur and S. Vassilvitskii, k-means++: The advantages of careful seeding, in *Proc. 18th Annual ACM-SIAM Symp. Discrete Algorithms*, 2007, pp. 1027–1035.
40. T. Su and J. G. Dy, In search of deterministic methods for initializing k-means and gaussian mixture clustering, *Intell. Data Anal.* **11**(4) (2007) 319–338.
41. M. R. Anderberg, *Cluster Analysis for Applications* (Academic Press, 1973).
42. M. J. Norušis, *IBM SPSS Statistics 19 Statistical Procedures Companion* (Addison Wesley, 2011).
43. T. Gonzalez, Clustering to minimize the maximum intercluster distance, *Theor. Comput. Sci.* **38**(2–3) (1985) 293–306.
44. H. Hotelling, Simplified calculation of principal components, *Psychometrika* **1**(1) (1936) 27–35.
45. E. H. Ruspini, Numerical methods for fuzzy clustering, *Inf. Sci.* **2**(3) (1970) 319–350.
46. M. Sezgin and B. Sankur, Survey over image thresholding techniques and quantitative performance evaluation, *J. Electron. Imaging* **13**(1) (2004) 146–165.
47. N. Otsu, A threshold selection method from gray level histograms, *IEEE Trans. Syst. Man Cybern.* **9**(1) (1979) 62–66.
48. P. K. Sahoo, S. Soltani and A. K. C. Wong, A survey of thresholding techniques, *Comput. Vis. Graphics Image Process.* **41**(2) (1988) 233–260.
49. S. U. Lee, S. Y. Chung and R. H. Park, A comparative performance study of several global thresholding techniques, *Comput. Vis. Graphics Image Process.* **52**(2) (1990) 171–190.
50. O. D. Trier and T. Taxt, Evaluation of binarization methods for document images, *IEEE Trans. Pattern Anal. Machine Intell.* **17**(3) (1995) 312–315.
51. R. Medina-Carnicer, F. J. Madrid-Cuevas, N. L. Fernandez-Garcia and A. Carmona-Poyato, Evaluation of global thresholding techniques in non-contextual edge detection, *Pattern Recog. Lett.* **26**(10) (2005) 1423–1434.
52. M. E. Celebi, Q. Wen, S. Hwang, H. Iyatomi and G. Schaefer, Lesion border detection in dermoscopy images using ensembles of thresholding methods, to appear in *Skin Res. Technol.* (2013).

53. A. Frank and A. Asuncion, UCI machine learning repository, University of California, Irvine, School of Information and Computer Sciences, <http://archive.ics.uci.edu/ml>, 2012.
 54. T. Onoda, M. Sakai and S. Yamada, Careful seeding method based on independent components analysis for k-means clustering, *J. Emerging Technol. Web Intell.* **4**(1) (2012) 51–59.
 55. G. Milligan and M. C. Cooper, A study of standardization of variables in cluster analysis, *J. Class.* **5**(2) (1988) 181–204.
 56. M. Matsumoto and T. Nishimura, Mersenne twister: A 623-dimensionally equidistributed uniform pseudo-random number generator, *ACM Trans. Model. Comput. Simul.* **8**(1) (1998) 3–30.
-



M. Emre Celebi received his B.Sc. degree in Computer Engineering from the Middle East Technical University (Ankara, Turkey) in 2002. He received his M.Sc. and Ph.D. degrees in Computer Science and Engineering from the University of

Texas at Arlington (Arlington, TX, USA) in 2003 and 2006, respectively. He is currently an Associate Professor and the founding director of the Image Processing and Analysis Laboratory in the Department of Computer Science at the Louisiana State University in Shreveport.

Dr. Celebi has actively pursued research in the field of image processing and analysis with an emphasis on medical image analysis and color image processing. He has worked on several projects funded by the US National Science Foundation (NSF) and National Institutes of Health (NIH) and published over 100 articles in premier journals and conference proceedings. His research contributions are covered in two recent books published by Wiley Interscience: *Image Processing — Principles and Applications* (Acharya and Ray, 2005) and *Clustering* (Xu and Wunsch, 2009). His recent research is funded by grants from the NSF and Louisiana Board of Regents.

Dr. Celebi is an editorial board member of 6 international journals, reviews for over 60 international journals, and served on the program committee of more than 60 international conferences. He has been invited to be a speaker at several colloquia, workshops, and conferences, is the organizer of several workshops, and the editor of several journal special issues and books. He is a senior member of the IEEE and SPIE.



Hassan A. Kingravi received his B.Sc. degree in Computer Science and Engineering from the University of Texas at Arlington in 2006. He is currently pursuing a Ph.D. degree in Electrical and Computer Engineering at Georgia Institute of Technology. His research interests include pattern recognition and machine learning.

include pattern recognition and machine learning.