

k-means clustering

In data mining, **k-means clustering** is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. This results into a partitioning of the data space into Voronoi cells.

The problem is computationally difficult (NP-hard), however there are efficient heuristic algorithms that are commonly employed and converge fast to a local optimum. These are usually similar to the expectation-maximization algorithm for mixtures of Gaussian distributions via an iterative refinement approach employed by both algorithms. Additionally, they both use cluster centers to model the data, however k -means clustering tends to find clusters of comparable spatial extent, while the expectation-maximization mechanism allows clusters to have different shapes.

Description

Given a set of observations ($\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$), where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k sets ($k \leq n$) $\mathbf{S} = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares (WCSS):

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

where $\boldsymbol{\mu}_i$ is the mean of points in S_i .

History

The term " k -means" was first used by James MacQueen in 1967,^[1] though the idea goes back to Hugo Steinhaus in 1957.^[2] The standard algorithm was first proposed by Stuart Lloyd in 1957 as a technique for pulse-code modulation, though it wasn't published until 1982.^[3]

Algorithms

Standard algorithm

The most common algorithm uses an iterative refinement technique. Due to its ubiquity it is often called the **k-means algorithm**; it is also referred to as **Lloyd's algorithm**, particularly in the computer science community.

Given an initial set of k means $\mathbf{m}_1^{(1)}, \dots, \mathbf{m}_k^{(1)}$ (see below), the algorithm proceeds by alternating between two steps:^[4]

Assignment step: Assign each observation to the cluster with the closest mean (i.e. partition the observations according to the Voronoi diagram generated by the means).

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\| \leq \|x_p - m_j^{(t)}\| \forall 1 \leq j \leq k\}$$

Where each x_p goes into exactly one $S_i^{(t)}$, even if it could go in two of them.

Update step: Calculate the new means to be the centroid of the observations in the cluster.

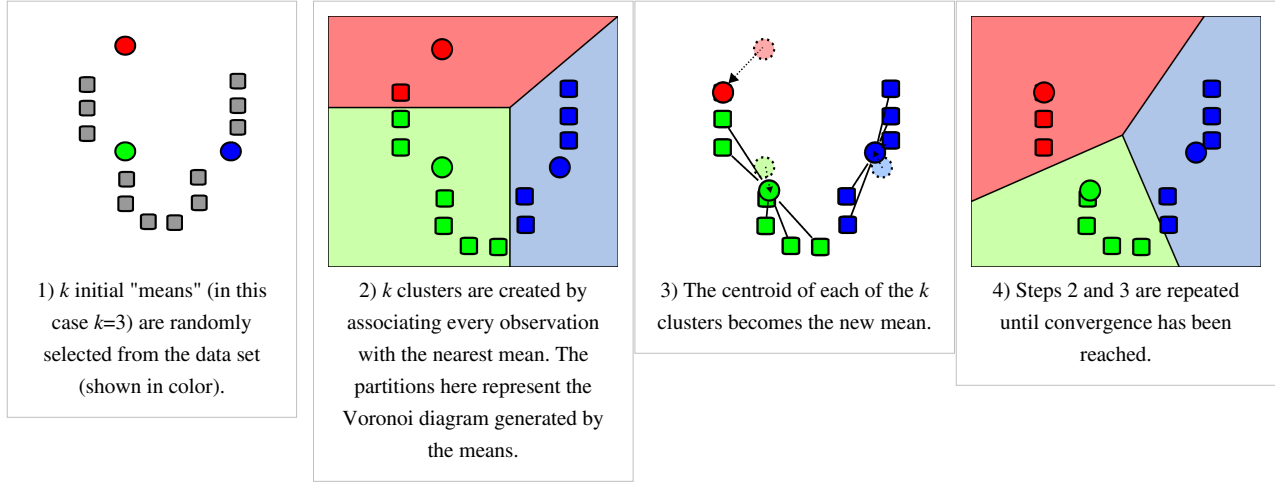
$$\mathbf{m}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j$$

The algorithm is deemed to have converged when the assignments no longer change.

Commonly used initialization methods are Forgy and Random Partition.^[5] The Forgy method randomly chooses k observations from the data set and uses these as the initial means. The Random Partition method first randomly

assigns a cluster to each observation and then proceeds to the Update step, thus computing the initial means to be the centroid of the cluster's randomly assigned points. The Forgy method tends to spread the initial means out, while Random Partition places all of them close to the center of the data set. According to Hamerly et al.,^[5] the Random Partition method is generally preferable for algorithms such as the k -harmonic means and fuzzy k -means. For expectation maximization and standard k -means algorithms, the Forgy method of initialization is preferable.

Demonstration of the standard algorithm



As it is a heuristic algorithm, there is no guarantee that it will converge to the global optimum, and the result may depend on the initial clusters. As the algorithm is usually very fast, it is common to run it multiple times with different starting conditions. However, in the worst case, k -means can be very slow to converge: in particular it has been shown that there exist certain point sets, even in 2 dimensions, on which k -means takes exponential time, that is $2^{\Omega(n)}$, to converge.^[6] These point sets do not seem to arise in practice: this is corroborated by the fact that the smoothed running time of k -means is polynomial.^[7]

The "assignment" step is also referred to as **expectation step**, the "update step" as **maximization step**, making this algorithm a variant of the *generalized* expectation-maximization algorithm.

Complexity

Regarding computational complexity, the k -means clustering problem for observations in d dimensions is:

- NP-hard in general Euclidean space d even for 2 clusters^{[8][9]}
- NP-hard for a general number of clusters k even in the plane^[10]
- If k and d are fixed, the problem can be exactly solved in time $O(n^{dk+1} \log n)$, where n is the number of entities to be clustered^[11]

Thus, a variety of heuristic algorithms are generally used.

- The k -means algorithm discussed below has polynomial smoothed running time. It is shown that^[7] for arbitrary set of n points in $[0, 1]^d$, if each point is independently perturbed by a normal distribution with mean 0 and variance σ^2 , then the expected running time of k -means algorithm is bounded by $O(n^{34} k^{34} d^8 \log^4(n) / \sigma^6)$, which is a polynomial in n , k , d and $1/\sigma$.
- Better bounds are proved for simple cases. For example,^[12] showed that the running time of k -means algorithm is bounded by $O(dn^4 M^2)$ for n points in an integer lattice $\{1, \dots, M\}^d$.

Variations

- Fuzzy C-Means Clustering is a soft version of K-means, where each data point has a fuzzy degree of belonging to each cluster.
- Gaussian mixture models trained with expectation-maximization algorithm (EM algorithm) maintains probabilistic assignments to clusters, instead of deterministic assignments, and multivariate Gaussian distributions instead of means.
- Several methods have been proposed to choose better starting clusters.^[13] One recent proposal is k-means++.
- The filtering algorithm uses kd-trees to speed up each k-means step.^[14]
- Some methods attempt to speed up each k-means step using coresets^[15] or the triangle inequality.^[16]
- Escape local optima by swapping points between clusters.^[17]
- The Spherical k-means clustering algorithm is suitable for directional data.^[18]
- The Minkowski metric weighted k-means deals with the problem of noise features by assigning weights to each feature per cluster^[19]

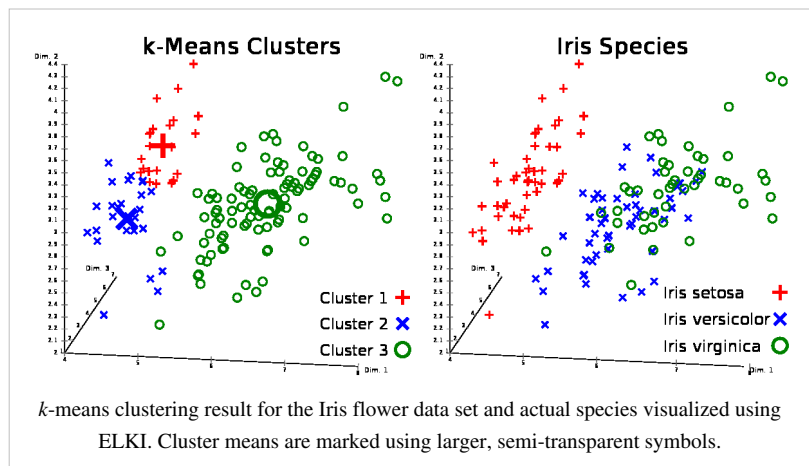
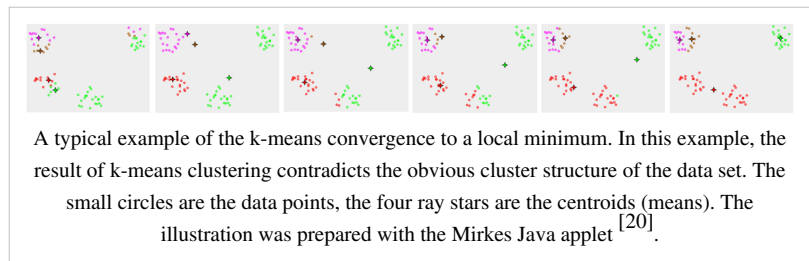
Discussion

The two key features of *k*-means which make it efficient are often regarded as its biggest drawbacks:

- Euclidean distance is used as a metric and variance is used as a measure of cluster scatter.
- The number of clusters *k* is an input parameter: an inappropriate choice of *k* may yield poor results. That is why, when performing *k*-means, it is important to run diagnostic checks for determining the number of clusters in the data set.
- Convergence to a local minimum may produce counterintuitive ("wrong") results (see example in Fig.).

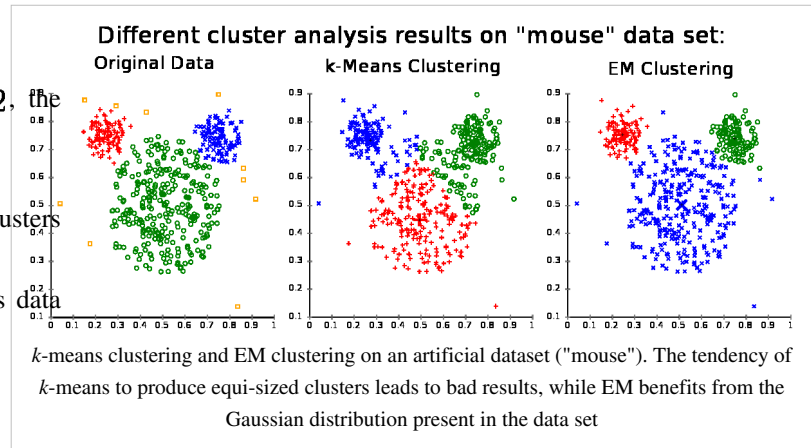
A key limitation of *k*-means is its cluster model. The concept is based on spherical clusters that are separable in

a way so that the mean value converges towards the cluster center. The clusters are expected to be of similar size, so that the assignment to the nearest cluster center is the correct assignment. When for example applying *k*-means with a value of $k = 3$ onto the well-known Iris



Iris

flower data set, the result often fails to separate the three Iris species contained in the data set. With $k = 2$, the two visible clusters (one containing two species) will be discovered, whereas with $k = 3$ one of the two clusters will be split into two even parts. In fact, $k = 2$ is more appropriate for this data set, despite the data set containing 3 classes. As with any other clustering algorithm, the k -means result relies on the data set to satisfy the assumptions made by the clustering algorithms. It works well on some data sets, while failing on others.



The result of k -means can also be seen as the Voronoi cells of the cluster means. Since data is split halfway between cluster means, this can lead to suboptimal splits as can be seen in the "mouse" example. The Gaussian models used by the Expectation-maximization algorithm (which can be seen as a generalization of k -means) are more flexible here by having both variances and covariances. The EM result is thus able to accommodate clusters of variable size much better than k -means as well as correlated clusters (not in this example).

Applications of the algorithm

k -means clustering in particular when using heuristics such as Lloyd's algorithm is rather easy to implement and apply even on large data sets. As such, it has been successfully used in various topics, ranging from market segmentation, computer vision, geostatistics,^[21] and astronomy to agriculture. It often is used as a preprocessing step for other algorithms, for example to find a starting configuration.

Relation to other statistical machine learning algorithms

k -means clustering, and its associated expectation-maximization algorithm, is a special case of a Gaussian mixture model, specifically, the limit of taking all covariances as diagonal, equal, and small. It is often easy to generalize a k -means problem into a Gaussian mixture model.^[22]

Mean shift clustering

Basic mean shift clustering algorithms maintain a set of data points the same size as the input data set. Initially, this set is copied from the input set. Then this set is iteratively replaced by the mean of those points in the set that are within a given distance of that point. By contrast, k -means restricts this updated set to k points usually much less than the number of points in the input data set, and replaces each point in this set by the mean of all points in the *input set* that are closer to that point than any other (e.g. within the Voronoi partition of each updating point). A mean shift algorithm that is similar then to k -means, called *likelihood mean shift*, replaces the set of points undergoing replacement by the mean of all points in the input set that are within a given distance of the changing set.^[23] One of the advantages of mean shift over k -means is that there is no need to choose the number of clusters, because mean shift is likely to find only a few clusters if indeed only a small number exist. However, mean shift can be much slower than k -means. Mean shift has soft variants much as k -means does.

Principal components analysis (PCA)

It has been shown^{[24][25]} that the relaxed solution of k-means clustering, specified by the cluster indicators, is given by the PCA (principal component analysis) principal components, and the PCA subspace spanned by the principal directions is identical to the cluster centroid subspace.

Bilateral filtering

k-means implicitly assumes that the ordering of the input data set does not matter. The bilateral filter is similar to K-means and mean shift in that it maintains a set of data points that are iteratively replaced by means. However, the bilateral filter restricts the calculation of the (kernel weighted) mean to include only points that are close in the ordering of the input data.^[23] This makes it applicable to problems such as image denoising, where the spatial arrangement of pixels in an image is of critical importance.

Similar problems

The set of squared error minimizing cluster functions also includes the k-medoids algorithm, an approach which forces the center point of each cluster to be one of the actual points, i.e., it uses medoids in place of centroids.

Software

Free

- Apache Mahout k-Means^[26]
- CrimeStat implements two spatial K-means algorithms, one of which allows the user to define the starting locations.
- ELKI contains k-means (with Lloyd and MacQueen iteration, along with different initializations such as k-means++ initialization) and various more advanced clustering algorithms
- MLPACK^[27] contains a K-Means implementation
- R kmeans^[28] implements a variety of algorithms^{[1][3][17]}
- SciPy vector-quantization^[29]
- Silverlight widget demonstrating k-means algorithm^[30]
- PostgreSQL extension for k-means^[31]
- CMU's GraphLab Clustering library^[32] Efficient multicore implementation for large scale data.
- Weka contains k-means and a few variants of it, including k-means++ and x-means.

Commercial

- IDL Cluster, Clust_Wts
 - *Mathematica* ClusteringComponents function^[33]
 - MATLAB kmeans^[34]
 - SAS FASTCLUS^[35]
 - VisuMap kMeans Clustering^[36]
-

Source code

- ELKI and Weka are written in Java and include k-means and variations
- K-means application in PHP,^[37] using VB,^[38] using Perl,^[39] using C++,^[40] using Matlab,^[41] using Ruby,^{[42][43]} using Python with scipy,^[44] using X10^[45]
- A parallel out-of-core implementation in C^[46]
- An open-source collection of clustering algorithms, including k-means, implemented in Javascript.^[47] Online demo.^[48]

Visualization, animation and examples

- ELKI can visualize k-means using Voronoi cells and Delaunay triangulation for 2D data. In higher dimensionality, only cluster assignments and cluster centers are visualized
- Demos of the K-means-algorithm^{[49][50][51][52][53]}
- K-means and K-medoids (Applet), University of Leicester^[20]
- Clustergram - cluster diagnostic plot - for visual diagnostics of choosing the number of (k) clusters (R code)^[54]

References

- [1] MacQueen, J. B. (1967). "Some Methods for classification and Analysis of Multivariate Observations" (<http://projecteuclid.org/euclid.bsmmsp/1200512992>). **1**. Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability. University of California Press. pp. 281–297. MR0214227. Zbl 0214.46201. . Retrieved 2009-04-07.
- [2] Steinhaus, H. (1957). "Sur la division des corps matériels en parties" (in French). *Bull. Acad. Polon. Sci.* **4** (12): 801–804. MR0090073. Zbl 0079.16403.
- [3] Lloyd, S. P. (1957). "Least square quantization in PCM". *Bell Telephone Laboratories Paper*. Published in journal much later: Lloyd, S. P. (1982). "Least squares quantization in PCM" (<http://www.cs.toronto.edu/~roweis/csc2515-2006/readings/lloyd57.pdf>). *IEEE Transactions on Information Theory* **28** (2): 129–137. doi:10.1109/TIT.1982.1056489. . Retrieved 2009-04-15.
- [4] MacKay, David (2003). "Chapter 20. An Example Inference Task: Clustering" (<http://www.inference.phy.cam.ac.uk/mackay/itprnn/ps/284.292.pdf>). *Information Theory, Inference and Learning Algorithms* (<http://www.inference.phy.cam.ac.uk/mackay/itila/book.html>). Cambridge University Press. pp. 284–292. ISBN 0-521-64298-1. MR2012999. .
- [5] Hamerly, G. and Elkan, C. (2002). "Alternatives to the k-means algorithm that find better clusterings" (<http://charlotte.ucsd.edu/users/elkan/cikm02.pdf>). *Proceedings of the eleventh international conference on Information and knowledge management (CIKM)*. .
- [6] Vattani, A. (2011). "k-means requires exponentially many iterations even in the plane" (<http://cseweb.ucsd.edu/users/avattani/papers/kmeans-journal.pdf>). *Discrete and Computational Geometry* **45** (4): 596–616. doi:10.1007/s00454-011-9340-1. .
- [7] Arthur, D.; Manthey, B.; Roeglin, H. (2009). "k-means has polynomial smoothed complexity". *Proceedings of the 50th Symposium on Foundations of Computer Science (FOCS)*.
- [8] Aloise, D.; Deshpande, A.; Hansen, P.; Popat, P. (2009). "NP-hardness of Euclidean sum-of-squares clustering". *Machine Learning* **75**: 245–249. doi:10.1007/s10994-009-5103-0.
- [9] Dasgupta, S. and Freund, Y. (July 2009). "Random Projection Trees for Vector Quantization". *Information Theory, IEEE Transactions on* **55**: 3229–3242. arXiv:0805.1390. doi:10.1109/TIT.2009.2021326.
- [10] Mahajan, M.; Nimbhorkar, P.; Varadarajan, K. (2009). "The Planar k-Means Problem is NP-Hard". *Lecture Notes in Computer Science* **5431**: 274–285. doi:10.1007/978-3-642-00202-1_24.
- [11] Inaba, M.; Katoh, N.; Imai, H. (1994). "Applications of weighted Voronoi diagrams and randomization to variance-based k-clustering". *Proceedings of 10th ACM Symposium on Computational Geometry*. pp. 332–339. doi:10.1145/177424.178042.
- [12] Arthur; Abhishek Bhowmick (2009). "A theoretical analysis of Lloyd's algorithm for k-means clustering".
- [13] <http://www.public.iastate.edu/~apghosh/files/IEEEclust2.pdf> A systematic evaluation of different methods for initializing the K-means clustering algorithm. Anna D. Peterson, Arka P. Ghosh and Ranjan Maitra. 2010.
- [14] Kanungo, T.; Mount, D. M.; Netanyahu, N. S.; Piatko, C. D.; Silverman, R.; Wu, A. Y. (2002). "An efficient k-means clustering algorithm: Analysis and implementation" (<http://www.cs.umd.edu/~mount/Papers/pami02.pdf>). *IEEE Trans. Pattern Analysis and Machine Intelligence* **24**: 881–892. doi:10.1109/TPAMI.2002.1017616. . Retrieved 2009-04-24.
- [15] Frahling, G.; Sohler, C. (2006). "A fast k-means implementation using coresets" ([http://www.frahling.de/Gereon_Frahling/Publications_files/A fast k-means implementation using Coresets \(Frahling, Sohler\).pdf](http://www.frahling.de/Gereon_Frahling/Publications_files/A%20fast%20k-means%20implementation%20using%20Coresets%20(Frahling,%20Sohler).pdf)). *Proceedings of the twenty-second annual symposium on Computational geometry (SoCG)*. .
- [16] Elkan, C. (2003). "Using the triangle inequality to accelerate k-means" (<http://www-cse.ucsd.edu/~elkan/kmeansicml03.pdf>). *Proceedings of the Twentieth International Conference on Machine Learning (ICML)*. .
- [17] Hartigan, J. A.; Wong, M. A. (1979). "Algorithm AS 136: A K-Means Clustering Algorithm". *Journal of the Royal Statistical Society, Series C (Applied Statistics)* **28** (1): 100–108. JSTOR 2346830.

- [18] Dhillon, I. S.; Modha, D. M. (2001). "Concept decompositions for large sparse text data using clustering". *Machine Learning* **42** (1): 143–175.
- [19] Amorim, R. C.; Mirkin, B (2012). "Minkowski metric, feature weighting and anomalous cluster initializing in K-Means clustering". *Pattern Recognition* **45** (3): 1061–1075. doi:10.1016/j.patcog.2011.08.012.
- [20] E.M. Mirkes, K-means and K-medoids applet (http://www.math.le.ac.uk/people/ag153/homepage/KmeansKmedoids/Kmeans_Kmedoids.html). University of Leicester, 2011.
- [21] Honarkhah, M and Caers, J, 2010, *Stochastic Simulation of Patterns Using Distance-Based Pattern Modeling* (<http://dx.doi.org/10.1007/s11004-010-9276-7>), Mathematical Geosciences, 42: 487 - 517
- [22] Press, WH; Teukolsky, SA; Vetterling, WT; Flannery, BP (2007). "Section 16.1. Gaussian Mixture Models and k-Means Clustering" (<http://apps.nrbook.com/empanel/index.html#pg=842>). *Numerical Recipes: The Art of Scientific Computing* (3rd ed.). New York: Cambridge University Press. ISBN 978-0-521-88068-8. .
- [23] Little, M.A.; Jones, N.S. (2011). "Generalized Methods and Solvers for Piecewise Constant Signals: Part I" (http://www.maxlittle.net/publications/pwc_filtering_arxiv.pdf). *Proc. Roy. Soc. A*. .
- [24] H. Zha, C. Ding, M. Gu, X. He and H.D. Simon (Dec. 2001). "Spectral Relaxation for K-means Clustering" (<http://ranger.uta.edu/~chqding/papers/Zha-Kmeans.pdf>). *Neural Information Processing Systems vol.14 (NIPS 2001)* (Vancouver, Canada): 1057–1064. .
- [25] Chris Ding and Xiaofeng He (July 2004). "K-means Clustering via Principal Component Analysis" (<http://ranger.uta.edu/~chqding/papers/KmeansPCA1.pdf>). *Proc. of Int'l Conf. Machine Learning (ICML 2004)*: 225–232. .
- [26] <http://cwiki.apache.org/MAHOUT/k-means-clustering.html>
- [27] <http://www.mlpack.org>
- [28] <http://stat.ethz.ch/R-manual/R-patched/library/stats/html/kmeans.html>
- [29] <http://docs.scipy.org/doc/scipy/reference/cluster.vq.html>
- [30] <http://www.codeding.com/?article=14>
- [31] <http://pgxn.org/dist/kmeans/>
- [32] <http://graphlab.org/clustering.html>
- [33] <http://reference.wolfram.com/mathematica/ref/ClusteringComponents.html>
- [34] <http://www.mathworks.com/access/helpdesk/help/toolbox/stats/kmeans.html>
- [35] http://support.sas.com/documentation/cdl/en/statug/63033/HTML/default/fastclus_toc.htm
- [36] <http://www.visumap.com/index.aspx?p=Products>
- [37] <http://www25.brinkster.com/densshade/kmeans.php.htm>
- [38] <http://people.revoledu.com/kardi/tutorial/kMean/download.htm>
- [39] http://www.lwebzem.com/cgi-bin/k_means/test3.cgi
- [40] <http://codingplayground.blogspot.com/2009/03/k-means-in-c.html>
- [41] http://people.revoledu.com/kardi/tutorial/kMean/matlab_kMeans.htm
- [42] <http://ai4r.rubyforge.org/index.html>
- [43] <http://github.com/reddavis/K-Means/tree/master>
- [44] <http://docs.scipy.org/doc/scipy/reference/cluster.vq.html>
- [45] <http://dist.codehaus.org/x10/applications/samples/KMeansDist.x10>
- [46] <http://www.cs.princeton.edu/~wdong/kmeans/>
- [47] <http://code.google.com/p/figure/FIGUE>
- [48] <http://web.science.mq.edu.au/~jydelort/figure/demo.html>
- [49] http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletKM.html
- [50] <http://siebn.de/other/yakmeans/>
- [51] (<http://informationandvisualization.de/blog/kmeans-and-voronoi-tesselation-built-processing>)
- [52] <http://www.javaworld.com/javaworld/jw-11-2006/jw-1121-thread.html>
- [53] <http://www.leet.it/home/lale/clustering/>
- [54] <http://www.r-statistics.com/2010/06/clustergram-visualization-and-diagnostics-for-cluster-analysis-r-code/>

Article Sources and Contributors

k-means clustering *Source:* <http://en.wikipedia.org/w/index.php?oldid=491976380> *Contributors:* Osm0sm0, 3mta3, Agor153, Alai, Andkaha, Annabel, Ashwin, BenFrantzDale, BlueScreenD, CBM, Charibdis, Charles Matthews, Chire, Chrike, ChrisDing, Chrisahn, Cincoutprabu, Corvus cornix, Cronholm144, David Epstein, Den fjättrade ankan, Denshade, Duncharris, Erniepan, FedeLebron, Fnielsen, Foma84, Foobarhoge, Gazpacho, Gfxguy, Giftlite, Golddan Gin, Greenleaf, Hakkinen, Headbomb, Helwr, Hgkamath, Homncruse, Honkkis, Illuminated, Ixfd64, Jnothman, John of Reading, JohnBlackburne, Jonsafari, Jsanchezalmeida, June8th, Killerandy, Kzafer, Leishi, Lessbread, LilHelpa, Mahlon, Manyu aditya, Mark Arsten, MarkPundurs, Materialsscientist, Mathias126, Mati22081979, Mauls, Maxlittle2007, Mcld, Melcombe, Memming, Michael Hardy, Miserlou, NedLevine, Nick Number, Ntvuok, Ostrouchov, Phoolimin, Pot, Qwertyus, Qwfp, Railwaycat, Ranumao, Ratiocinate, Rcalhoun, Rich Farmbrough, Ricky81682, Robert K S, SamuelRiv, Sanchom, SciCompTeacher, SeparateWays, Simeon87, Smartcat, Soren.harward, SpuriousQ, Sundirac, Talgalili, Tbmurphy, Toninowiki, Turketwh, UkPaolo, Utacsecd, Weston.pace, Wonderful597, Woolleynick, Zanetu, 178 anonymous edits

Image Sources, Licenses and Contributors

Image:K Means Example Step 1.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:K_Means_Example_Step_1.svg *License:* Creative Commons Attribution-ShareAlike 3.0 Unported *Contributors:* Weston.pace

Image:K Means Example Step 2.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:K_Means_Example_Step_2.svg *License:* Creative Commons Attribution-ShareAlike 3.0 Unported *Contributors:* Weston.pace

Image:K Means Example Step 3.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:K_Means_Example_Step_3.svg *License:* Creative Commons Attribution-ShareAlike 3.0 Unported *Contributors:* Weston.pace

Image:K Means Example Step 4.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:K_Means_Example_Step_4.svg *License:* Creative Commons Attribution-ShareAlike 3.0 Unported *Contributors:* Weston.pace

File:K-means convergence to a local minimum.png *Source:* http://en.wikipedia.org/w/index.php?title=File:K-means_convergence_to_a_local_minimum.png *License:* Creative Commons Attribution-Sharealike 3.0 *Contributors:* User:Agor153

File:Iris Flowers Clustering kMeans.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:Iris_Flowers_Clustering_kMeans.svg *License:* Public Domain *Contributors:* Chire

File:ClusterAnalysis Mouse.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:ClusterAnalysis_Mouse.svg *License:* Public Domain *Contributors:* Chire

License

Creative Commons Attribution-Share Alike 3.0 Unported
[//creativecommons.org/licenses/by-sa/3.0/](http://creativecommons.org/licenses/by-sa/3.0/)