# A feasible algorithm for locating concave and convex zones of interval data and its use in statistics-based clustering

Vladik Kreinovich [a], Eric J. Pauwels [b,c], Scott A. Ferson [d] and Lev Ginzburg [d]

[a] *Department of Computer Science, University of Texas at El Paso, El Paso, TX 79968, USA*
E-mail: vladik@cs.utep.edu
[b] *Centre for Mathematics and Computer Science, CWI, Kruislaan 413, NL-1098, SJ Amsterdam,*
*The Netherlands*
E-mail: eric.pauwels@cwi.nl
[c] *Katolieke University Leuven, P.O. Box 94079, NL-1090, GB Amsterdam, The Netherlands*
[d] *Applied Biomathematics, 100 North Country Road, Setauket, NY 11733, USA*
E-mail: {scott;lev}@ramas.com

Often, we need to divide $n$ objects into clusters based on the value of a certain quantity $x$. For example, we can classify insects in the cotton field into groups based on their size and other geometric characteristics. Within each cluster, we usually have a unimodal distribution of $x$, with a probability density $\rho(x)$ that increases until a certain value $x_0$ and then decreases. It is therefore natural, based on $\rho(x)$, to determine a cluster as the interval between two local minima, i.e., as a union of adjacent increasing and decreasing segments. In this paper, we describe a feasible algorithm for solving this problem.

**Keywords:** clustering, statistical methods, Kolmogorov–Smirnov statistic, convex segments of interval data

**AMS subject classification:** 91C20, 65G20, 52A10, 62D05

## 1. Formulation of the practical problem

In each area of interest, we are studying certain classes of objects: e.g., in astronomy, we study stars, galaxies, etc.; in entomology, we study insects. To study different objects, we perform different measurements on these objects; as a result, each object is characterized by the values $x_1, \ldots, x_d$ of the measured quantities.

Usually, for each of these variables $x_i$, there are physical bounds that bound its possible values, and within these bounds, all values are, in principle, possible. In practice, however, there are a few subintervals that contain the vast majority of the objects, and values in between these subintervals are rare. In other words, most objects belong to *clusters*, with few objects in between.

This fact is very useful in practice: if we want to control objects from a given class, then, instead of designing different control strategies for different values $x = (x_1, \ldots, x_d)$, we can use the fact that objects within each cluster are similar to each other, and design a single control strategy for all the objects from a cluster.

Practical example: cotton contains insects. Some of these insects destroy the cotton crop, some are harmless. If we use insecticides against harmless insects, we pollute the environment; if we do not use insecticides against harmful insects, we may lose the crop. To characterize insects, we can measure different geometric characteristics $x_1, \ldots, x_d$, and then cluster the corresponding points $x = (x_1, \ldots, x_d)$. The resulting clusters crudely correspond to species. Therefore, to distinguish between harmful and harmless insects, it is sufficient to subdivide all the insects into clusters, and then decide, for each cluster, whether its insects are harmful or not (for details, see, e.g., [12]).

It is therefore important to classify objects into clusters. There exist many clustering techniques. In some cases, when we have a large number of objects, we can use statistical techniques to get a statistically validated classification; see, e.g., [10].

In many real-life situations, however, there is not enough data to apply these statistically validated techniques. Instead, practitioners use heuristic clustering techniques that use fuzzy logic, neural networks, etc.; see, e.g., [2,3,5,13]. The problem with these methods is that their results are not completely justified; moreover, most of these methods require that we choose certain parameters, and different choices of these parameters lead to different subdivision into clusters. So, there is a need to design justified clustering methods.

Such methods are described in this paper. In this description, we use several ideas first announced in [9,14,15].

## 2. Main idea: how to reformulate this practical problem in precise mathematical terms

### 2.1. We will only consider 1-D case

In general, we have many parameters $x_i$ that describe different aspects of the objects from our class, so clustering is a multi-dimensional problem. In many practical multi-D cases, it is possible to find a *single* parameter – it can be one of the parameters $x_i$ or a combination of these parameters – that is sufficient to classify objects into clusters. In some cases, one parameter is not sufficient: e.g., when we classify insects by size, we get several classes, but some of these classes actually contain several different clusters. In many such cases, it is sufficient to use repeated 1-D clustering: first, we cluster by one parameter $x_i$, then we subcluster the resulting clusters by another parameter $x_j$, etc.

To serve such situations, in this paper, we consider a 1-D clustering problem. In this problem, we have a single parameter $x$ that characterizes different objects from our class, and we have several ($n$) objects with different values $x^{(1)}, \ldots, x^{(n)}$ of this parameter. We want to divide these objects into clusters.

### 2.2. Traditional statistical approach to clustering

The value of the parameter $x$ depends not only on the cluster to which this object belongs, but also on many other factors. For example, the size of an insect is determined

not only by its species, but also by the weather conditions, by the environment, by the presence or absence of chemicals that are damaging to these insects, etc. As a result, different objects within the same class exhibit random variations from the average value corresponding to this class.

It is therefore reasonable to consider, for each cluster, a probability distribution that describes how frequently different values $x$ occur for objects from this cluster. This is the main idea behind statistical clustering methods: for each class, we measure a lot of objects, determine the corresponding probability distribution; then, for each new object with the value $x$, we get the probability (actually, probability density) $\rho_i(x)$ that $x$ is from cluster $i$, and we assign the object $x$ to the cluster $i_0$ for which this probability is the largest, i.e., for which $\rho_{i_0}(x) \geqslant \rho_i(x)$ for all $i \neq i_0$.

### 2.3. Towards a formal definition of a cluster

In this paper, we consider situations in which we do not have enough observations to determine all the distributions $\rho_i(x)$. Instead, all we observe is a sample $x^{(1)}, \ldots, x^{(n)}$, and we do not know which object corresponds to which cluster. If we knew the probability distributions $\rho_i(x)$ and the frequency $p_i$ of objects from each cluster, then we could say that the observed data are a sample from a mixture distribution, with the density $\rho(x) = \sum p_i \cdot \rho_i(x)$. Since we do not know neither the densities $\rho_i(x)$ nor the frequencies $p_i$, all we can say is that the values $x^{(k)}$ are a sample from *some* probability distribution with an unknown density $\rho(x)$.

Let us relate this observation with the above (informal) description of a cluster: clusters are subintervals that contain the vast majority of the objects, and values in between these subintervals are rarer than inside them. A natural corollary of this description is that immediately outside the cluster subinterval $[\underline{a}, \overline{a}]$, the density is smaller than inside, i.e., that the density function $\rho(x)$ is increasing for $x = \underline{a}$ and decreasing for $x = \overline{a}$. A continuous function that increases at $\underline{a}$ and decreases at $\overline{a} > \underline{a}$ must attain a (local) maximum inside the interval $[\underline{a}, \overline{a}]$; vice versa, if a function has a local maximum, then a sufficiently narrow interval around this maximum is a cluster in this sense.

Therefore, if, within an interval, there are at least two local maxima, this means that we can form at least two clusters. Thus, it is natural to identify clusters with local maxima of the probability density function (pdf) $\rho(x)$. To be more precise, clusters are neighborhoods of local maxima, neighborhoods that go both ways until the corresponding local minimum – the point at which decreasing changes to increasing or vice versa.

In other words, it is natural to define a cluster as the interval between two consequent local minima of the pdf $\rho(x)$, i.e., as a union of adjacent increasing and decreasing segments. Within each cluster, we have a unimodal distribution of $x$, with a probability density $\rho(x)$ that increases until a certain value $x_0$ and then decreases.

### 2.4. A similar problem with known solution and why we cannot use it

If we could determine bounds on $\rho(x)$ based on the empirical data (i.e., on the values $x^{(1)}, \ldots, x^{(n)}$), then we would be able to use the known algorithms for finding local minima and local maxima of interval-valued functions; see, e.g., [18]. The problem is that, based on empirical data, we cannot find bounds on the pdf.

### 2.5. Enter Kolmogorov–Smirnov bounds

What we can find is bounds on the *cumulative density function* (CDF) $F(x)$. Specifically, based on the sample values $x^{(1)}, \ldots, x^{(n)}$, we can determine an *empirical CDF* $F_{\text{emp}}(x)$: for each $x$, $F_{\text{emp}}(x)$ is defined as the ratio $\#\{k \mid x^{(k)} \leqslant x\}/n$. This empirical distribution is the easiest to compute if we first sort the values $x^{(k)}$ in the increasing order $x^{(1)} \leqslant x^{(2)} \leqslant \cdots \leqslant x^{(n)}$; then:

- $F_{\text{emp}}(x) = 0$ for $x < x^{(1)}$;
- $F_{\text{emp}}(x) = k/n$ for $x^{(k)} \leqslant x < x^{(k+1)}$;
- $F_{\text{emp}}(x) = 1$ for $x \geqslant x^{(n)}$.

Kolmogorov–Smirnov theorem (see, e.g., [20]) states that if the actual (unknown) PDF $F(x)$ is located on a known interval, then, for any given confidence level $\alpha$, we can find the value $\varepsilon$ for which, with this confidence, we have $\max_x |F_{\text{emp}}(x) - F(x)| \leqslant \varepsilon$. Thus, with this given confidence level, we know that for every $x$, we have

$$F(x) \in \mathbf{F}(x) = \left[\underline{F}(x), \overline{F}(x)\right] \stackrel{\text{def}}{=} \left[\max\left(F_{\text{emp}}(x) - \varepsilon, 0\right), \min\left(F_{\text{emp}}(x) + \varepsilon, 1\right)\right]. \quad (1)$$

Based on this interval information, we want to find out whether it is possible that $\rho(x)$ is increasing or decreasing on a given interval of values $x$. The function $\rho(x) = \mathrm{d}F(x)/\mathrm{d}x$ is decreasing (increasing) if and only if $F(x)$ is concave (convex). Thus, to solve our problem, we must determine concave and convex zones of the interval-valued function $\mathbf{F}(x)$.

In this paper, we propose a $O(n \cdot \log(n))$ time algorithm for determining such zones.

## 3. Proposed algorithm

Our algorithm uses known $O(n \cdot \log(n))$ time incremental convex hull algorithms [4,7,8,16]; these algorithms, given $n$ points $p_1, \ldots, p_n$ on the 2-D plane, find the convex hull of these points. These algorithms are called *incremental* because in the time $O(n \cdot \log(n))$, they not only find the convex hull of all $n$ points $p_1, \ldots, p_n$, but also, for all $k$ from 1 to $n$, convex hulls of the sets $\{p_1, \ldots, p_k\}$.

The algorithm for checking whether there exists a convex function $F(x)$ within given bounds is as follows:

- First, we sort all the values $x^{(k)}$ in the increasing order; this sorting takes $O(n \cdot \log(n))$ time. Without losing generality, we will therefore assume that the values $x^{(k)}$ are already sorted, i.e., that $x^{(1)} \leqslant x^{(2)} \leqslant \cdots \leqslant x^{(n)}$.

- We compute the values $\underline{F}_k \overset{\text{def}}{=} \max(k/n - \varepsilon, 0)$ and $\overline{F}_k \overset{\text{def}}{=} \min((k-1)/n + \varepsilon, 1)$ ($\overline{F}_1 = 0$); this computation takes O($n$) time.

- Then, we use an incremental convex hull algorithm to compute, for every $k$ from 1 to $n$, the convex hull $C$ of the points $\overline{p}_1, \ldots, \overline{p}_k$, where $\overline{p}_k \overset{\text{def}}{=} (x^{(k)}, \overline{F}_k)$, and check whether the points $\underline{p}_1, \ldots, \underline{p}_k$, where $\underline{p}_k \overset{\text{def}}{=} (x^{(k)}, \underline{F}_k)$, are outside the interior of this convex hull (i.e., on or below the piecewise-linear curve $F_{\text{conv}}(x)$ describing the lower boundary of this convex hull). Since the points are already sorted, an appropriate version of Graham's scan (see, e.g., [16]) will provide all these checks within O($n$) time.

To find zones of convexity and concavity, we do the following:

- First, we run the above algorithm until we find the last value $k^+$ for which the values $\underline{p}_k$ are outside the interior of the convex hull, i.e., for which it is still possible to have a convex function $F(x) \in \mathbf{F}(x)$ for $x \leqslant x^{(k^+)}$.

- Then, we similarly process the values starting with $x^{(n)}$ backwards and find the smallest value $k^-$ for which it is still possible to have a concave function $F(x) \in \mathbf{F}(x)$ for $x \geqslant x^{(k^-)}$.

- If $k^- \leqslant k^+$, this means that it is possible to have a unimodal distribution $F(x)$, and $[x^{(k^-)}, x^{(k^+)}]$ is the interval of possible locations of its mode. In cluster terms, it means that the data is consistent with having only one cluster, with a "center" at some point $x \in [x^{(k^-)}, x^{(k^+)}]$.

- If $k^- > k^+$, this means that there are several clusters; to find these clusters, we apply the same algorithm to data starting with the $(k^+ + 1)$st point.

## 4. Justification of the proposed algorithm

One can easily see that the above algorithm requires O($n \cdot \log(n)$) time. (Application of this algorithm to real-life data shows that this algorithm is not only theoretically feasible, it is also practically is efficient and useful [9,14,15].)

Let us now show that this algorithm works correctly. Our algorithm is based on the ability to check convexity, so it is sufficient to show that the algorithm for checking convexity is correct.

1°. Let us first show that if there is a convex function $F(x) \in \mathbf{F}(x)$, then none of the points $\underline{p}_k$ are inside the interior of the convex hull $C$.

Indeed, let us assume that there is a convex function $F(x) \in \mathbf{F}(x)$. Then, the area $G \overset{\text{def}}{=} \{(x, y) \mid y \geqslant F(x)\}$ above this function is convex. The function $\overline{F}(x)$ is constant on each interval $[x^{(k-1)}, x^{(k)})$; thus, for arbitrary small $\delta > 0$, we have $\overline{F}(x^{(k)} - \delta) = \overline{F}(x^{(k-1)})$. From the definition of $G$, we conclude that $(x^{(k)} - \delta, \overline{F}(x^{(k)} - \delta)) = (x^{(k)} - \delta, \overline{F}(x^{(k-1)})) \in G$. In the limit $\delta \to 0$, we conclude that $\overline{p}_k = (x^{(k)}, \overline{F}(x^{(k-1)})) \in G$. Also,

by definition of the set $G$, the pairs $\underline{p}_k = (x^{(k)}, \underline{F}(x^{(k)}))$ are not inside the interior of this area $G$.

By definition, the convex hull $C$ is the intersection of all the convex sets that contain given points, thus, $C \subseteq G$; hence the points $\underline{p}_k$ cannot be inside the interior of $C$ either.

$2°$. Vice versa, let us show that if all the points $\underline{p}_k$ are not inside the interior of the convex hull $C$, then there exists a convex function $F(x) \in \mathbf{F}(x)$.

We will prove that we can take $F_{\text{conv}}(x)$ as the desired function $F(x)$. Indeed, as a lower bound for a convex set, this function is convex.

Since the points $\underline{p}_k = (x^{(k)}, \underline{F}(x^{(k)}))$ are not inside the interior of the convex hull, we conclude that $\underline{F}(x^{(k)}) \leqslant F_{\text{conv}}(x^{(k)})$. Since the function $\underline{F}(x)$ is constant on the interval $[x^{(k)}, x^{(k+1)})$, and the function $F_{\text{conv}}(x)$ is nondecreasing on this interval, we conclude that $\underline{F}(x) = \underline{F}(x^{(k)}) \leqslant F_{\text{conv}}(x^{(k)}) \leqslant F_{\text{conv}}(x)$ – i.e., $\underline{F}(x) \leqslant F_{\text{conv}}(x)$ for all $x$.

Since all the points $\overline{p}_k = (x^{(k)}, \overline{F}(x^{(k-1)}))$ are inside the convex hull, we conclude that $F_{\text{conv}}(x^{(k)}) \leqslant \overline{F}(x^{(k-1)})$; similarly, since the function $\overline{F}(x)$ is constant on the interval $[x^{(k-1)}, x^{(k)})$, and the function $F_{\text{conv}}(x)$ is nondecreasing on this interval, we conclude that $F_{\text{conv}}(x) \leqslant \overline{F}(x)$ for all $x$. Thus, $F(x) \in \mathbf{F}(x)$ for all $x$.

## 5. What if measurements come with interval uncertainty?

The above algorithm can be used to check whether there is a convex function $F(x)$ within arbitrary piecewise-constant bounds $\mathbf{F}(x) = [\underline{F}(x), \overline{F}(x)]$. An important case of this general situation stems from the fact that the values $x^{(k)}$ come from measurements, and measurements are never 100% accurate. As a result, the measured values $\tilde{x}^{(k)}$ are, in general, different from the actual (unknown) values $x^{(k)}$ of the measured characteristics. Usually, we know the upper bound $\Delta$ for the (absolute value of) the measurement error $\Delta x_k \stackrel{\text{def}}{=} \tilde{x}^{(k)} - x^{(k)}$; thus, instead of the exact value of $x^{(k)}$, we only know the *interval* $\mathbf{x}_k = [\underline{x}_k, \overline{x}_k] = [\tilde{x}^{(k)} - \Delta, \tilde{x}^{(k)} + \Delta]$ of possible values of $x^{(k)}$.

In this case, for every $x$, we have $F(x) \in [\underline{F}(x), \overline{F}(x)]$, where $\underline{F}(x) = \max(\underline{F}_{\text{emp}}(x) - \varepsilon, 0)$, $\overline{F}(x) = \min(\overline{F}_{\text{emp}}(x) + \varepsilon, 1)$, $\underline{F}_{\text{emp}}(x) = \#\{k \mid \overline{x}_k \leqslant x\}/n$, and $\overline{F}_{\text{emp}}(x) = \#\{k \mid \underline{x}_k \leqslant x\}/n$. Similarly to the above case, we have piecewise-constant bounds for $F(x)$ that change values only at $\underline{x}_k$ and $\overline{x}_k$. It is therefore sufficient to require that $F(x^{(i)}) \in [\underline{F}(x^{(i)}), \overline{F}(x^{(i)})]$ for $2n$ values $x^{(i)}$ that coincide with either the lower endpoints $\underline{x}_k$ or with upper endpoints $\overline{x}_k$ of the given intervals $\mathbf{x}_k$.

Thus, we can apply the above algorithm – starting with the sorting of the values $x^{(i)}$, i.e., in this case, of $2n$ values $\underline{x}_k$ and $\overline{x}_k$ – and find convexity zones in $O(n \cdot \log(n))$ time. (A similar algorithm for detecting monotonicity zones is described in [11].)

## 6. How to parallelize our algorithm

Although our algorithm is pretty fast, its running time still grows with the number of points $n$. So, when the number of points is large, it is desirable to speed it up. A natural way to speed up an algorithm is to run it in parallel. There exist algorithms that compute the convex hull of $n$ points in $O(\log(n))$ time on $n$ processors [1] (see also [6]). If we use this algorithm, we can check, for every $k$, in $O(\log(n))$ time, whether the values $x^{(1)}, \ldots, x^{(k)}$ are consistent with the convexity of $F(x)$. We can use this check in two different ways:

- we can run $n$ checks in parallel; thus, by using $O(n^2)$ processors, we detect the desired zones in $O(\log(n))$ time;
- alternatively, we can use bisection on the interval $[1, n]$ to find the last value $k$ that is still consistent with convexity; binary search requires $\log(n)$ checks, so we find the zones in $O(\log^2(n))$ time by using $n$ processors.

It is worth mentioning that a similar drastic speed-up is possible if we use parallel computations in a similar problem of detecting areas of monotonicity [19].

## Acknowledgements

## References

[1]  S.G. Akl and K.A. Lyons, *Parallel Computational Geometry* (Prentice-Hall, Englewood Cliffs, NJ, 1993).

[2]  J.C. Bezdek, J. Keller, R. Krisnapuram and N.R. Pal, *Fuzzy Models and Algorithms for Pattern Recognition and Image Processing* (Kluwer, Boston, 1999).

[3]  J.C. Bezdek and S.K. Pal, eds., *Fuzzy Models for Pattern Recognition* (IEEE Press, New York, 1992).

[4]  J.-D. Boissonant and M. Yvinec, *Algorithmic Geometry* (Cambridge Univ. Press, Cambridge, 1998).

[5]  H. Bunke and A. Kandel, eds., *Neuro-Fuzzy Pattern Recognition* (World Scientific, Singapore, 2000).

[6]  W. Chen, K. Wada and K. Kawaguchi, Robust algorithms for constructing strongly convex hulls in parallel, Theoret. Comput. Sci. 289 (2002) 277–295.

[7]  T.H. Cormen, C.E. Leiserson, R.L. Rivest and C. Stein, *Introduction to Algorithms* (MIT Press, Cambridge, MA, and McGraw-Hill, New York, 2001).

[8]  H. Edelsbrunner, *Algorithms in Computational Geometry* (Springer, Berlin, 1987).

[9]  G. Frederix and E.J. Pauwels, Segmentation based on principled self-supervising clustering, in: *Abstracts of the 1st SIAM Conf. on Imaging Science*, Boston, MA, 4–6 March 2002, p. 25.

[10] K. Fukunaga, *Introduction to Statistical Pattern Recognition* (Academic Press, San Diego, CA, 1990).

[11] J. Lorkowski and V. Kreinovich, If we measure a number, we get an interval. What if we measure a function or an operator?, Reliable Comput. 2(3) (1996) 287–298.

[12] H.T. Nguyen, N.R. Prasad, V. Kreinovich and H. Gassoumi, Some practical applications of soft computing and data mining, in: *Data Mining and Computational Intelligence*, eds. A. Kandel, H. Bunke and M. Last (Springer, Berlin, 2001) pp. 273–307.

[13] S.K. Pal and S. Mitra, *Neuro-Fuzzy Pattern Recognition: Methods in Soft Computing* (Wiley, New York, 1999).

[14] E.J. Pauwels and G. Frederix, Image segmentation by nonparametric clustering based on the Kolmogorov–Smirnov distance, in: *Proc. of the 6th European Conf. on Computer Vision ECCV'2000*, ed. D. Vernon, Dublin, Ireland, 26 June–1 July 2000, Part II, Lecture Notes in Computer Science, Vol. 1843 (Springer, New York, 2000) pp. 85–99.

[15] E.J. Pauwels, G. Frederix and G. Caenaen, Image segmentation based on statistically principled clustering, in: *Proc. of 2001 Internat. IEEE Conf. on Image Processing*, Vol. 3, 2001, pp. 66–69.

[16] F.P. Preparata and M.I. Shamos, *Computational Geometry: An Introduction* (Springer, New York, 1989).

[17] Ramas Website `http://www.ramas.com` (2002).

[18] K. Villaverde and V. Kreinovich, A linear-time algorithm that locates local extrema of a function of one variable from interval measurement results, Interval Comput. 4 (1993) 176–194.

[19] K. Villaverde and V. Kreinovich, Parallel algorithm that locates local extrema of a function of one variable from interval measurement results, in: *Extended Abstracts of APIC'95: Internat. Workshop on Applications of Interval Computations*, El Paso, TX, 23–25 February 1995, Supplement of Reliable Computing (1995) pp. 212–219.

[20] H.M. Wadsworth, Jr., ed., *Handbook of Statistical Methods for Engineers and Scientists* (McGraw-Hill, New York, 1990).