

UNCERTAIN CENTROID BASED PARTITIONAL CLUSTERING OF UNCERTAIN DATA

Francesco Gullo * Andrea Tagarelli †

* Yahoo! Research
Barcelona, Spain

† Dept. of Electronics, Computer and Systems Science
University of Calabria, Italy

38th International Conference on Very Large Data Bases (VLDB)
August 27-31, 2012
Istanbul, Turkey

Uncertainty

Uncertainty inherently affects data from a wide range of emerging application domains:

- sensor data
- location-based services (e.g., moving objects data)
- biomedical and biometric data (e.g., gene expression data)
- distributed applications
- RFID data

Generally due to noisy factors, such as signal noise, instrumental errors, wireless transmission

Uncertain Objects (UO) (1)

Modeling by *regions (domains) of definition and probability density functions (pdfs)*

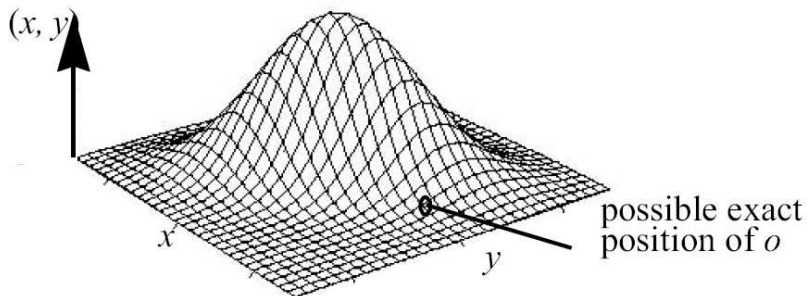


Figure borrowed from [Kriegel and Pfeifle, ICDM 2005]

Uncertain Objects (UO) (2)

- m -dimensional region
- multivariate pdf defined over the region

Definition (uncertain object)

An **uncertain object** o is a pair (\mathcal{R}, f) :

- $\mathcal{R} \subseteq \mathbb{R}^m$ is the m -dimensional domain region in which o is defined
- $f : \mathbb{R}^m \rightarrow \mathbb{R}_0^+$ is the probability density function of o at each point $\vec{x} \in \mathbb{R}^m$ such that:

$$f(\vec{x}) > 0, \forall \vec{x} \in \mathcal{R} \quad \text{and} \quad f(\vec{x}) = 0, \forall \vec{x} \in \mathbb{R}^m \setminus \mathcal{R}$$

Clustering Uncertain Objects

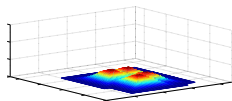
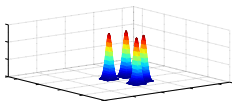
Major approaches:

- partitional approaches:
 - uncertain version of k -Means [Chau et al., PAKDD 2006] and its relative optimizations [Lee et al., ICDM Work. 2007, Kao et al., TKDE 2010, Ngai et al., Information Systems 2011]
 - uncertain version of k -Medoids [Gullo et al., SUM 2008]
- density-based approaches:
 - uncertain version of DBSCAN [Kriegel and Pfeifle, KDD 2005]
 - uncertain version of OPTICS [Kriegel and Pfeifle, ICDM 2005]
- hierarchical approaches [Gullo et al., ICDM 2008]

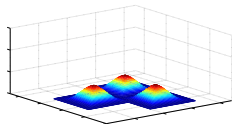
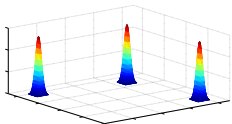
Partitional approaches include the fastest methods so far defined

Intuition

Approaches to partitional clustering of uncertain objects should take into account both **central tendency** and **variance** of the input uncertain objects



Uncertain objects with the same central tendency: lower-variance, more-compact cluster (left) and higher-variance, less-compact cluster (right)



Uncertain objects with different central tendency: lower-variance, less-compact cluster (left) and higher-variance, more-compact cluster (right)

Contributions

- We formally show that existing formulations of partitional clustering of uncertain objects do not comply with the intuition about central tendency and variance
- We propose a novel formulation to the problem of clustering uncertain objects based on the notion of *U-centroid*
- Given that the expression of the U-centroid is not analytically computable, we derive some theoretical properties to be efficiently exploited as *closed-form update rules* for the proposed objective function
- We define an efficient *local-search procedure* based on these rules

UK-means and MMVar

Partitional clustering of uncertain objects relies on two main notions: *cluster centroid* (\bar{C}), and *cluster compactness* (J)

Most prominent existing formulations:

- *UK-means* [Chau et al., PAKDD'06] \rightarrow cluster centroid is a **deterministic** object

$$\bar{C}_{UK} = \frac{1}{|C|} \sum_{o \in C} \vec{\mu}(o)$$

$$J_{UK}(C) = \sum_{o \in C} ED(o, \bar{C}_{UK}), \text{ where } ED(o, \bar{C}_{UK}) = \int_{\vec{x} \in \mathcal{R}} \|\vec{x} - \bar{C}_{UK}\|^2 f(\vec{x}) d\vec{x}$$

- *MMvar* [Gullo et al., ICDM'10] \rightarrow cluster centroid is an **uncertain** object

$$\bar{C}_{MM} = (\bar{\mathcal{R}}_{MM}, \bar{f}_{MM}), \text{ where } \bar{\mathcal{R}}_{MM} = \bigcup_{o \in C} \mathcal{R} \text{ and } \bar{f}_{MM}(\vec{x}) = \frac{1}{|C|} \sum_{o \in C} f(\vec{x})$$

$$J_{MM}(C) = \sigma^2(\bar{C}_{MM})$$

Issues of UK-means and MMVar formulations

- The deterministic centroid representation in UK-means is not able to discriminate among different variances
- The MMvar formulation does not overcome this issue, although its centroid representation involves uncertainty

Proposition

Given a cluster C of m -dimensional uncertain objects, where $o = (\mathcal{R}, f)$, $\forall o \in C$, it holds that $J_{MM}(C) = |C|^{-1} J_{UK}(C)$.

A straightforward (inappropriate) solution

- **Idea:** combine the notions of MMVar centroid with the UK-means cluster compactness criterion

$$\hat{J}(C) = \sum_{o \in C} \widehat{ED}(o, \bar{C}_{MM}),$$

$$\text{where } \widehat{ED}(o, \bar{C}_{MM}) = \int_{\vec{x} \in \mathcal{R}} \int_{\vec{y} \in \bar{\mathcal{R}}_{MM}} \|\vec{x} - \vec{y}\|^2 f(\vec{x}) \bar{f}_{MM}(\vec{y}) \, d\vec{x} \, d\vec{y}$$

- Unfortunately, such an objective function \hat{J} is not appropriate as it is equivalent to functions J_{UK} and J_{MM}

Proposition

Given a cluster C of m -dimensional uncertain objects, where $o = (\mathcal{R}, f)$, $\forall o \in C$, it holds that

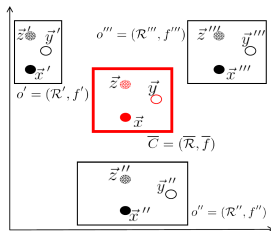
$$\hat{J}(C) = 2 |C| J_{MM}(C) = 2 J_{UK}(C).$$

Our proposal

- 1 Introducing a novel notion of cluster centroid
- 2 Defining a cluster compactness criterion based on this novel cluster centroid definition which meets the requirements about central tendency and variance

U-centroid

Cluster centroid as *random variable* summarizing all possible deterministic representations of the objects in the cluster



Two key advantages:

- Shortcomings of a deterministic centroid notion are addressed
- Clear stochastic meaning

U-centroid: analytical expression (1)

Theorem

Given a cluster $C = \{o_1, \dots, o_{|C|}\}$ of m -dimensional uncertain objects, where $o_i = (\mathcal{R}_i, f_i)$ and $\mathcal{R}_i = [\ell_i^{(1)}, u_i^{(1)}] \times \dots \times [\ell_i^{(m)}, u_i^{(m)}]$, $\forall i \in [1..|C|]$, let $\bar{C} = (\bar{\mathcal{R}}, \bar{f})$ be the U-centroid of C defined by employing the squared Euclidean norm as distance to be minimized. It holds that:

$$\bar{f}(\vec{x}) = \int_{\vec{x}_1 \in \mathcal{R}_1} \dots \int_{\vec{x}_{|C|} \in \mathcal{R}_{|C|}} \mathbb{I}[\vec{x} = \frac{1}{|C|} \sum_{i=1}^{|C|} \vec{x}_i] \prod_{i=1}^{|C|} f_i(\vec{x}_i) d\vec{x}_1 \dots d\vec{x}_{|C|}$$

$$\bar{\mathcal{R}} = \left[\frac{1}{|C|} \sum_{i=1}^{|C|} \ell_i^{(1)}, \frac{1}{|C|} \sum_{i=1}^{|C|} u_i^{(1)} \right] \times \dots \times \left[\frac{1}{|C|} \sum_{i=1}^{|C|} \ell_i^{(m)}, \frac{1}{|C|} \sum_{i=1}^{|C|} u_i^{(m)} \right]$$

where $\mathbb{I}[A]$ is the indicator function, which is 1 when the event A occurs, 0 otherwise.

U-centroid based cluster compactness criterion

Two main requirements for the proposed cluster compactness criterion J :

- It should rely on the U-centroid notion so to meet the requirements about central tendency and variance
- The expression of the pdf \bar{f} in the proposed U-centroid is not analytically computable $\Rightarrow J$ should be such that it can be optimized without requiring to explicitly compute \bar{f}

A first solution: minimizing the U-centroid variance

Minimizing the variance of the U-centroid (similarly to MMVar) does not work, as it is equivalent to minimizing the average variance of the individual uncertain objects in the cluster:

Theorem

Given a cluster $C = \{o_1, \dots, o_{|C|}\}$ of m -dimensional uncertain objects, where $o_i = (\mathcal{R}_i, f_i)$, $\forall i \in [1..|C|]$, let $\bar{C} = (\bar{\mathcal{R}}, \bar{f})$ be the U-centroid of C . It holds that $\sigma^2(\bar{C}) = |C|^{-2} \sum_{i=1}^{|C|} \sigma^2(o_i)$.

Minimizing the expected distance between uncertain objects and U-centroid (1)

$$J(C) = \sum_{o \in C} \widehat{ED}(o, \bar{C})$$

Observation 1: J takes into account both central tendency and variance

Theorem

Let $C = \{o_1, \dots, o_{|C|}\}$ be a cluster of uncertain objects, where $o_i = (\mathcal{R}_i, f_i)$, and $\bar{C} = (\bar{\mathcal{R}}, \bar{f})$ be the U-centroid of C . It holds that:

$$J(C) = \sum_{j=1}^m \left(\frac{\Psi_C^{(j)}}{|C|} + \Phi_C^{(j)} - \frac{\Upsilon_C^{(j)}}{|C|} \right) = \frac{1}{|C|} \sum_{i=1}^{|C|} \sigma^2(o_i) + \sum_{o \in C} ED \left(o, \frac{1}{|C|} \sum_{o \in C} \vec{\mu}(o) \right)$$

where

$$\Psi_C^{(j)} = \sum_{i=1}^{|C|} (\sigma^2)_j(o_i) \quad \Phi_C^{(j)} = \sum_{i=1}^{|C|} (\mu_2)_j(o_i) \quad \Upsilon_C^{(j)} = \left(\sum_{i=1}^{|C|} \mu_j(o_i) \right)^2$$

Minimizing the expected distance between uncertain objects and U-centroid (2)

Observation 2: Given a cluster C , the value of J of any other cluster resulting from adding/removing an object to/from C can be computed according to an efficient closed-form expression

Corollary

Let C be a cluster of uncertain objects, and $C^+ = C \cup \{o^+\}$, $C^- = C \setminus \{o^-\}$ be two clusters defined by adding an object $o^+ \notin C$ to C and removing an object $o^- \in C$ from C , respectively. It holds that:

$$J(C^+) = \sum_{j=1}^m \left(\frac{\Psi_{C^+}^{(j)}}{|C^+|+1} + \Phi_{C^+}^{(j)} - \frac{\Upsilon_{C^+}^{(j)}}{|C^+|+1} \right) \quad J(C^-) = \sum_{j=1}^m \left(\frac{\Psi_{C^-}^{(j)}}{|C^-|-1} + \Phi_{C^-}^{(j)} - \frac{\Upsilon_{C^-}^{(j)}}{|C^-|-1} \right)$$

The UCPC local-search algorithm

Input: A set \mathcal{D} of UO; the number k of output clusters

Output: A partition \mathcal{C} of \mathcal{D} , where $|\mathcal{C}| = k$

- 1: compute $\bar{\mu}(o), \bar{\mu}_2(o), \bar{\sigma}^2(o), \forall o \in \mathcal{D}$
- 2: $\mathcal{C} \leftarrow \text{initialPartition}(\mathcal{D}, k)$, compute $\Psi_{\mathcal{C}}^{(j)}, \Phi_{\mathcal{C}}^{(j)}, \Upsilon_{\mathcal{C}}^{(j)}, J(\mathcal{C})$
- 3: **repeat**
- 4: $V \leftarrow \sum_{\mathcal{C} \in \mathcal{C}} J(\mathcal{C})$
- 5: **for all** $o \in \mathcal{D}$ **do**
- 6: $\mathcal{C}^* \leftarrow \text{argmin}_{\mathcal{C} \in \mathcal{C}} V - [J(\mathcal{C}^o) + J(\mathcal{C})] + [J(\mathcal{C}^o \setminus \{o\}) + J(\mathcal{C} \cup \{o\})]$
- 7: **if** $\mathcal{C}^* \neq \mathcal{C}^o$ **then**
- 8: $\mathcal{C} \leftarrow \mathcal{C} \setminus \{\mathcal{C}^*, \mathcal{C}^o\} \cup \{\mathcal{C}^+, \mathcal{C}^-\}$
- 9: replace $\Psi_{\mathcal{C}^*}^{(j)}, \Phi_{\mathcal{C}^*}^{(j)}, \Upsilon_{\mathcal{C}^*}^{(j)}, J(\mathcal{C}^*)$ with $\Psi_{\mathcal{C}^+}^{(j)}, \Phi_{\mathcal{C}^+}^{(j)}, \Upsilon_{\mathcal{C}^+}^{(j)}, J(\mathcal{C}^+), \forall j \in [1..m]$
- 10: replace $\Psi_{\mathcal{C}^o}^{(j)}, \Phi_{\mathcal{C}^o}^{(j)}, \Upsilon_{\mathcal{C}^o}^{(j)}, J(\mathcal{C}^o)$ with $\Psi_{\mathcal{C}^-}^{(j)}, \Phi_{\mathcal{C}^-}^{(j)}, \Upsilon_{\mathcal{C}^-}^{(j)}, J(\mathcal{C}^-), \forall j \in [1..m]$
- 11: **until** no object in \mathcal{D} is relocated

- UCPC converges to a local optimum of function J in a finite number l of iterations
- UCPC works in $\mathcal{O}(l k |\mathcal{D}| m)$

Evaluation methodology (1)

- Benchmark datasets from UCI (Iris, Wine, Glass, Ecoli, Yeast, Image, Abalone, Letter) where uncertainty is generated **synthetically** and modeled according to *Uniform* (U), *Normal* (N), and *Exponential* (E) pdfs
- Real (gene expression) datasets where uncertainty is inherently present

(a) Benchmark datasets

<i>dataset</i>	<i>obj.</i>	<i>attr.</i>	<i>classes</i>
Iris	150	4	3
Wine	178	13	3
Glass	214	10	6
Ecoli	327	7	5
Yeast	1,484	8	10
Image	2,310	19	7
Abalone	4,124	7	17
Letter	7,648	16	10

(b) Real datasets

<i>dataset</i>	<i>obj.</i>	<i>attr.</i>
Neuroblastoma	22,282	14
Leukaemia	22,690	21

Evaluation methodology (2)

- Evaluation in terms of:
 - **accuracy** (external and internal clustering evaluation)
 - **efficiency**
- Competitors: MMVar (MMV), UK-means (UKM), UK-medoids (UKmed), UAHC, \mathcal{F} DBSCAN (\mathcal{F} DB), \mathcal{F} OPTICS (\mathcal{F} OPT)

Accuracy results: benchmark datasets

		<i>F-measure</i> ($\Theta \in [-1, 1]$)							
		<i>pdf</i>	<i>FDB</i>	<i>FOPT</i>	<i>UAHC</i>	<i>UKmed</i>	<i>UKM</i>	<i>MMV</i>	UCPC
<i>avg score</i>	U		-.189	.055	.089	.210	.081	.193	.429
	N		-.081	-.046	.149	-.028	.019	.199	.287
	E		-.317	-.088	-.008	-.011	-.137	.200	.223
<i>overall avg. score</i>			-.196	-.026	.077	.057	-.012	.198	.313
<i>overall avg. gain</i>			+.509	+.339	+.236	+.256	+.324	+.115	—

		<i>Quality</i> ($Q \in [-1, 1]$)							
		<i>pdf</i>	<i>FDB</i>	<i>FOPT</i>	<i>UAHC</i>	<i>UKmed</i>	<i>UKM</i>	<i>MMV</i>	UCPC
<i>avg score</i>	U		.021	.089	.027	.084	.042	.345	.375
	N		.061	.115	.091	.089	.127	.139	.189
	E		-.001	.025	0	.011	.015	.199	.200
<i>overall avg. score</i>			.027	.076	.039	.061	.061	.228	.255
<i>overall avg. gain</i>			+.228	+.179	+.216	+.194	+.194	+.027	—

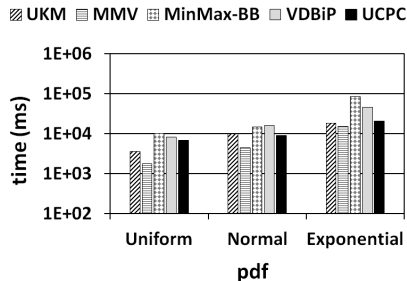
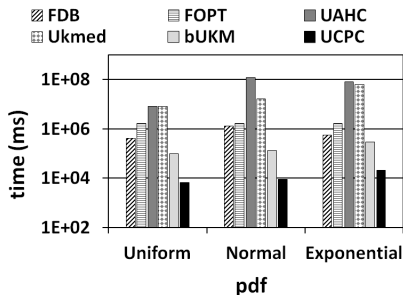
Accuracy results: real datasets

<i>data</i>	<i>#clust.</i>	<i>Quality (Q ∈ [-1, 1])</i>						
		<i>FDB</i>	<i>FOPT</i>	<i>UAHC</i>	<i>UKmed</i>	<i>UKM</i>	<i>MMV</i>	UCPC
<i>Neuro. avg score</i>		-.004	.010	.630	.045	.060	.544	.576
<i>Leuk. avg score</i>		-.018	.190	.192	.231	.430	.433	.471
<i>over. avg score</i>		-.011	.100	.411	.138	.245	.489	.523
<i>over. avg gain</i>		+.534	+.423	+.112	+.385	+.278	+.034	—

Efficiency results: benchmark datasets

- Efficiency evaluation also involves optimized versions of UK-means, i.e., MinMax-BB and VDBiP

Letter

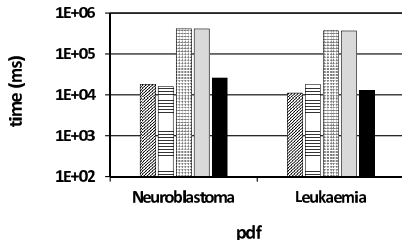
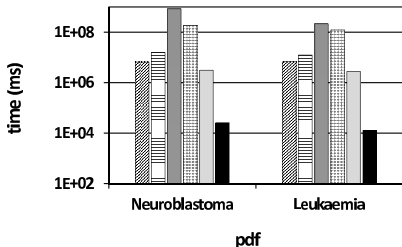


Efficiency results: real datasets

Real datasets

FDB
 FOPT
 UAHC
 Ukmed
 bUKM
 UCPC

UKM
 MMV
 MinMax-BB
 VDBiP
 UCPC



Conclusions

- Existing formulations of partitional clustering of uncertain objects miss some crucial requirements about central tendency and variance of the objects to be clustered
- Novel notion of cluster centroid, called U-centroid
- Effective and efficient U-centroid based cluster compactness criterion
- Efficient local-search heuristic to optimize the proposed objective function

Thanks!