



Photo: Vladimir Batagelj, *UNI-LJ*

Metoda voditeljev

Vladimir Batagelj

FMF, Univerza Ljubljani

1143. Sredin seminar, Ljubljana, 21. marec 2007

Outline

1	Leaders method	1
2	The dynamic clusters method	2

Leaders method

In order to support our intuition in further development we shall briefly describe a simple version of dynamic clusters method – the *leaders* or k -means method, which is the basis of the ISODATA program (one among the most popular clustering programs) and several recent 'data-mining' methods. In the leaders method the criterion function has the form SR.

The basic scheme of leaders method is simple:

determine C_0 ; $C := C_0$;

repeat

 determine for each $C \in C$ its leader \bar{C} ;

 the new clustering C is obtained by assigning each unit
 to its nearest leader

until leaders stabilize

To obtain a 'good' solution and an impression of its quality we can repeat this procedure with different (random) C_0 .

The dynamic clusters method

The dynamic clusters method is a generalization of the above scheme. Let us denote:

- Λ – set of *representatives*
- $L \subseteq \Lambda$ – *representation*
- Ψ – set of *feasible representations*
- $W : \Phi \times \Psi \rightarrow \mathbb{R}_0^+$ – *extended criterion function*
- $G : \Phi \times \Psi \rightarrow \Psi$ – *representation function*
- $F : \Phi \times \Psi \rightarrow \Phi$ – *clustering function*

and

Basic scheme of the dynamic clusters method

the following conditions have to be satisfied:

$$W0. \quad P(\mathbf{C}) = \min_{L \in \Psi} W(\mathbf{C}, L)$$

the functions G and F tend to improve (diminish) the value of the extended criterion function W :

$$W1. \quad W(\mathbf{C}, G(\mathbf{C}, L)) \leq W(\mathbf{C}, L)$$

$$W2. \quad W(F(\mathbf{C}, L), L) \leq W(\mathbf{C}, L)$$

then the *dynamic clusters method* can be described by the scheme:

$\mathbf{C} := \mathbf{C}_0; L := L_0;$

repeat

$L := G(\mathbf{C}, L);$

$\mathbf{C} := F(\mathbf{C}, L)$

until the clustering stabilizes

Properties of DCM

To this scheme corresponds the sequence $v_n = (\mathbf{C}_n, \mathbf{L}_n), n \in \mathbb{N}$ determined by relations

$$\mathbf{L}_{n+1} = G(\mathbf{C}_n, \mathbf{L}_n) \quad \text{and} \quad \mathbf{C}_{n+1} = F(\mathbf{C}_n, \mathbf{L}_{n+1})$$

and the sequence of values of the extended criterion function $u_n = W(\mathbf{C}_n, \mathbf{L}_n)$. Let us also denote $u^* = P(\mathbf{C}^*)$. Then it holds:

Theorem 1.1 *For every $n \in \mathbb{N}$, $u_{n+1} \leq u_n$, $u^* \leq u_n$, and if for $k > m$, $v_k = v_m$ then $\forall n \geq m : u_n = u_m$.*

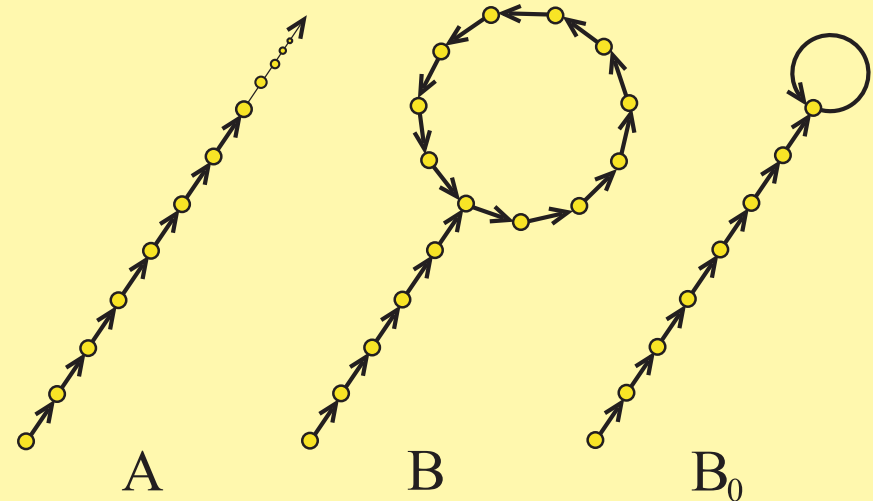
The Theorem 2.6 states that the sequence u_n is monotonically decreasing and bounded, therefore it is convergent. Note that the limit of u_n is not necessarily u^* – the dynamic clusters method is a local optimization method.

... types of DCM sequences

Type A: $\neg \exists k, m \in \mathbb{N}, k > m : v_k = v_m$

Type B: $\exists k, m \in \mathbb{N}, k > m : v_k = v_m$

Type B₀: Type B with $k = m + 1$



The DCM sequence (v_n) is of type B if

- sets Φ and Ψ are both finite.

For example, when we select a representative of C among its members.

- $\exists \delta > 0 : \forall n \in \mathbb{N} : (v_{n+1} \neq v_n \Rightarrow u_n - u_{n+1} > \delta)$

Because the sets U and consequently Φ are finite we expect from a good dynamic clusters procedure to stabilize in finite number of steps – is of type B.

Additional requirement

The conditions W0, W1 and W2 are not strong enough to ensure this. We shall try to compensate the possibility that the set of representations Ψ is infinite by the additional requirement:

$$W3. \quad W(\mathbf{C}, G(\mathbf{C}, L)) = W(\mathbf{C}, L) \Rightarrow L = G(\mathbf{C}, L)$$

With this requirement the 'symmetry' between Φ and Ψ is destroyed. We could reestablish it by the requirement:

$$W4. \quad W(F(\mathbf{C}, L, L)) = W(\mathbf{C}, L) \Rightarrow \mathbf{C} = F(\mathbf{C}, L)$$

but it turns out that W4 often fails. For this reason we shall avoid it.

Theorem 1.2 *If W3 holds and if there exists $m \in \mathbb{N}$ such that $u_{m+1} = u_m$, then also $L_{m+1} = L_m$.*

Simple clustering and representation functions

Usually, in the applications of the DCM, the clustering function takes the form $F : \Psi \rightarrow \Phi$. In this case the condition W2 simplifies to: $W(F(L), L) \leq W(\mathbf{C}, L)$ which can be expressed also as $F(L) \in \text{Min}_{\mathbf{C} \in \Phi} W(\mathbf{C}, L)$. For such, *simple* clustering functions it holds:

Theorem 1.3 *If the clustering function F is simple and if there exists $m \in \mathbb{N}$ such that $L_{m+1} = L_m$, then for every $n \geq m : v_n = v_m$.*

What can be said about the case when G is *simple* – has the form $G : \Phi \rightarrow \Psi$?

Theorem 1.4 *If W3 holds and the representation function G is simple then:*

- a. $G(\mathbf{C}) = \arg \min_{L \in \Psi} W(\mathbf{C}, L)$
- b. $\exists k, m \in \mathbb{N}, k > m \forall i \in \mathbb{N} : v_{k+i} = v_{m+i}$
- c. $\exists m \in \mathbb{N} \forall n \geq m : u_n = u_m$
- d. *if also F is simple then $\exists m \in \mathbb{N} \forall n \geq m : v_n = v_m$*

Original DCM

In the original dynamic clusters method (Diday, 1979) both functions F and G are simple – $F : \Psi \rightarrow \Phi$ and $G : \Phi \rightarrow \Psi$.

We proved, if also W3 holds and the functions F and G are simple, then:

$$G0. \quad G(\mathbf{C}) = \operatorname{argmin}_{L \in \Psi} W(\mathbf{C}, L)$$

and

$$F0. \quad F(L) \in \operatorname{Min}_{\mathbf{C} \in \Phi} W(\mathbf{C}, L)$$

In other words, given an extended criterion function W , the relations G0 and F0 define an appropriate pair of functions G and F such that the DCM stabilizes in finite number of steps.