Methods of Network Analysis Clustering and Blockmodeling 2. Approaches to Clustering

Vladimir Batagelj University of Ljubljana, Slovenia

University of Konstanz, Algorithms and Data Structures

June 6, 2002, 14-16h, room F 426

# **Approaches to Clustering**

- local optimization
- dynamic programming
- hierarchical methods; agglomerative methods; Lance-Williams formula; dendrogram; inversions; adding methods
- leaders and the dynamic clusters method
- graph theory (next, 3. lecture);

## **Local optimization**

Often for a given optimization problem  $(\Phi, P)$  there exist rules which relate to each element of the set  $\Phi$  some elements of  $\Phi$ . We call them *local transformations*.

The elements which can be obtained from a given element are called *neighbors* – local transformations determine the *neighborhood relation*  $S \subseteq \Phi \times \Phi$  in the set  $\Phi$ . The *neighborhood* of element  $X \in \Phi$  is called the set  $S(X) = \{Y : XSY\}$ . The element  $X \in \Phi$  is a *local minimum* for the *neighborhood structure*  $(\Phi, S)$  iff

 $\forall \mathbf{Y} \in S(\mathbf{X}) : P(\mathbf{X}) \le P(\mathbf{Y})$ 

In the following we shall assume that S is reflexive,  $\forall X \in \Phi : XSX$ . They are the basis of the local activity proved have

They are the basis of the *local optimization procedure* 

select  $X_0$ ;  $X := X_0$ ; while  $\exists Y \in S(X) : P(Y) < P(X)$  do X := Y;

which starting in an element of  $X_0 \in \Phi$  repeats moving to an element determined by local transformation which has better value of the criterion function until no such element exists.

### **Clustering neigborhoods**

Usually the neighborhood relation in local optimization clustering procedures over  $P_k(\mathbf{U})$  is determined by the following two transformations:

• *transition*: clustering C' is obtained from C by moving a unit from one cluster to another

$$\mathbf{C}' = (\mathbf{C} \setminus \{C_u, C_v\}) \cup \{C_u \setminus \{\mathbf{X}_s\}, C_v \cup \{\mathbf{X}_s\}\}$$

• *transposition*: clustering C' is obtained from C by interchanging two units from different clusters

 $\mathbf{C}' = (\mathbf{C} \setminus \{C_u, C_v\}) \cup \{(C_u \setminus \{\mathbf{X}_p\}) \cup \{\mathbf{X}_q\}, (C_v \setminus \{\mathbf{X}_q\}) \cup \{\mathbf{X}_p\}\}$ 

The transpositions preserve the number of units in clusters.

#### Hints

Two basic implementation approaches are usually used: *stored data* approach and *stored dissimilarity matrix* approach.

If the constraints are not too stringent, the relocation method can be applied directly on  $\Phi$ ; otherwise, we can transform using *penalty function method* the problem to an equivalent nonconstrained problem  $(P_k, Q)$  with  $Q(\mathbf{C}) = P(\mathbf{C}) + \alpha K(\mathbf{C})$  where  $\alpha > 0$  is a large constant and  $K(\mathbf{C}) = 0$ , for  $\mathbf{C} \in \Phi$ , and  $K(\mathbf{C}) > 0$  otherwise.

There exist several improvements of the basic relocation algorithm: simulated annealing, tabu search, ... (Aarts and Lenstra, 1997).

The *initial clustering*  $C_0$  can be given; most often we generate it randomly. Let  $c[s] = u \Leftrightarrow X_s \in C_u$ . Fill the vector c with the desired number of units in each cluster and shuffle it:

for p := n downto 2 do begin q := random(1, p); swap(c[p], c[q]) end;

### **Quick scanning of neighbors**

Testing  $P(\mathbf{C}') < P(\mathbf{C})$  is equivalent to  $P(\mathbf{C}) - P(\mathbf{C}') > 0$ . For the S criterion function

$$\Delta P(\mathbf{C}, \mathbf{C}') = P(\mathbf{C}) - P(\mathbf{C}') = p(C_u) + p(C_v) - p(C'_u) - p(C'_v)$$

Additional simplifications can be done considering relations between  $C_u$  and  $C'_u$ , and between  $C_v$  and  $C'_v$ .

Let us illustrate this on the generalized Ward's method. For this purpose it is useful to introduce the quantity

$$a(C_u, C_v) = \sum_{\mathbf{X} \in C_u, \mathbf{Y} \in C_v} w(\mathbf{X}) \cdot w(\mathbf{Y}) \cdot d(\mathbf{X}, \mathbf{Y})$$

Using the quantity  $a(C_u, C_v)$  we can express p(C) in the form  $p(C) = \frac{a(C,C)}{2w(C)}$  and the equality mentioned in the introduction of the generalized Ward clustering problem: if  $C_u \cap C_v = \emptyset$  then

$$w(C_u \cup C_v) \cdot p(C_u \cup C_v) = w(C_u) \cdot p(C_u) + w(C_v) \cdot p(C_v) + a(C_u, C_v)$$

#### $\Delta$ for the generalized Ward's method

Let us analyze the transition of a unit  $X_s$  from cluster  $C_u$  to cluster  $C_v$ : We have  $C'_u = C_u \setminus \{X_s\}$ ,  $C'_v = C_v \cup \{X_s\}$ ,

 $w(C_u) \cdot p(C_u) = w(C'_u) \cdot p(C'_u) + a(X_s, C'_u) = (w(C_u) - w(X_s)) \cdot p(C'_u) + a(X_s, C'_u)$ 

and

$$w(C'_v) \cdot p(C'_v) = w(C_v) \cdot p(C_v) + a(\mathbf{X}_s, C_v)$$

From  $d(X_s, X_s) = 0$  it follows  $a(X_s, C_u) = a(X_s, C'_u)$ . Therefore

$$p(C'_u) = \frac{w(C_u) \cdot p(C_u) - a(X_s, C_u)}{w(C_u) - w(X_s)} \quad p(C'_v) = \frac{w(C_v) \cdot p(C_v) + a(X_s, C_v)}{w(C_v) + w(X_s)}$$

and finally

$$\Delta P(\mathbf{C}, \mathbf{C}') = p(C_u) + p(C_v) - p(C'_u) - p(C'_v) = = \frac{w(X_s) \cdot p(C_v) - a(X_s, C_v)}{w(C_v) + w(X_s)} - \frac{w(X_s) \cdot p(C_u) - a(X_s, C_u)}{w(C_u) - w(X_s)}$$

In the case when d is the squared Euclidean distance it is possible to derive also expression for corrections of centers (Späth, 1977).

## **Dynamic programming**

Suppose that  $Min(\Phi_k, P) \neq \emptyset$ , k = 1, 2, ... Denoting  $P^*(\mathbf{U}, k) = P(\mathbf{C}_k^*(\mathbf{U}))$  we can derive the generalized *Jensen equality* (Batagelj, Korenjak and Klavžar, 1994):

$$P^{*}(\mathbf{U},k) = \begin{cases} p(\mathbf{U}) & \{\mathbf{U}\} \in \mathbf{\Phi}_{1} \\ & \underset{\substack{\emptyset \subset C \subset \mathbf{U} \\ \exists \mathbf{C} \in \mathbf{\Phi}_{k-1}(\mathbf{U} \setminus C): \mathbf{C} \cup \{C\} \in \mathbf{\Phi}_{k}(\mathbf{U})}{\end{cases}} (P^{*}(\mathbf{U} \setminus C, k-1) \oplus p(C)) & k > 1 \end{cases}$$

This is a *dynamic programming* (Bellman) equation which, for some special constrained problems, that keep the size of  $\Phi_k$  small, allows us to solve the clustering problem by the adapted Fisher's algorithm.

### **Hierarchical methods**

The set of feasible clusterings  $\Phi$  determines the *feasibility predicate*  $\Phi(\mathbf{C}) \equiv \mathbf{C} \in \Phi$  defined on  $\mathcal{P}(\mathcal{P}(\mathbf{U}) \setminus \{\emptyset\})$ ; and conversely  $\Phi \equiv \{\mathbf{C} \in \mathcal{P}(\mathcal{P}(\mathbf{U}) \setminus \{\emptyset\}) : \Phi(\mathbf{C})\}$ .

In the set  $\Phi$  the relation of *clustering inclusion*  $\sqsubseteq$  can be introduced by

$$\mathbf{C}_1 \sqsubseteq \mathbf{C}_2 \equiv \forall C_1 \in \mathbf{C}_1, C_2 \in \mathbf{C}_2 : C_1 \cap C_2 \in \{\emptyset, C_1\}$$

we say also that the clustering  $C_1$  is a *refinement* of the clustering  $C_2$ .

It is well known that  $(P(\mathbf{U}), \sqsubseteq)$  is a partially ordered set (even more, semimodular lattice). Because any subset of partially ordered set is also partially ordered, we have: Let  $\mathbf{\Phi} \subseteq P(\mathbf{U})$  then  $(\mathbf{\Phi}, \sqsubseteq)$  is a partially ordered set.

The clustering inclusion determines two related relations (on  $\Phi$ ):

$$\mathbf{C}_1 \sqsubset \mathbf{C}_2 \equiv \mathbf{C}_1 \sqsubseteq \mathbf{C}_2 \land \mathbf{C}_1 \neq \mathbf{C}_2$$
 – strict inclusion, and

 $\mathbf{C}_1 \sqsubseteq \mathbf{C}_2 \equiv \mathbf{C}_1 \sqsubset \mathbf{C}_2 \land \neg \exists \mathbf{C} \in \mathbf{\Phi} : (\mathbf{C}_1 \sqsubset \mathbf{C} \land \mathbf{C} \sqsubset \mathbf{C}_2) \quad -\text{predecessor.}$ 

#### **Conditions on the structure of the set of feasible clusterings**

We shall assume that the set of feasible clusterings  $\Phi \subseteq P(\mathbf{U})$  satisfies the following conditions:

F1.  $\mathbf{O} \equiv \{\{X\} : X \in \mathbf{U}\} \in \mathbf{\Phi}$ 

**F2.** The feasibility predicate  $\Phi$  is *local* – it has the form  $\Phi(\mathbf{C}) = \bigwedge_{C \in \mathbf{C}} \varphi(C)$ where  $\varphi(C)$  is a predicate defined on  $\mathcal{P}(\mathbf{U}) \setminus \{\emptyset\}$  (clusters).

The intuitive meaning of  $\varphi(C)$  is:  $\varphi(C) \equiv$  the cluster *C* is 'good'. Therefore the locality condition can be read: a 'good' clustering  $\mathbf{C} \in \Phi$  consists of 'good' clusters.

**F3.** The predicate  $\Phi$  has the property of *binary heredity* with respect to the *fusibility* predicate  $\psi(C_1, C_2)$ , i.e.,

$$C_1 \cap C_2 = \emptyset \land \varphi(C_1) \land \varphi(C_2) \land \psi(C_1, C_2) \Rightarrow \varphi(C_1 \cup C_2)$$

This condition means: in a 'good' clustering, a fusion of two 'fusible' clusters produces a 'good' clustering.

#### ... conditions

**F4.** The predicate  $\psi$  is *compatible* with clustering inclusion  $\sqsubseteq$ , i.e.,

 $\forall \mathbf{C}_1, \mathbf{C}_2 \in \mathbf{\Phi} : (\mathbf{C}_1 \sqsubset \mathbf{C}_2 \land \mathbf{C}_1 \setminus \mathbf{C}_2 = \{C_1, C_2\} \Rightarrow \psi(C_1, C_2) \lor \psi(C_2, C_1))$ 

**F5.** The *interpolation* property holds in  $\Phi$ , i.e.,  $\forall C_1, C_2 \in \Phi$ :

 $(\mathbf{C}_1 \sqsubset \mathbf{C}_2 \land \operatorname{card}(\mathbf{C}_1) > \operatorname{card}(\mathbf{C}_2) + 1 \Rightarrow \exists \mathbf{C} \in \mathbf{\Phi} : (\mathbf{C}_1 \sqsubset \mathbf{C} \land \mathbf{C} \sqsubset \mathbf{C}_2))$ 

These conditions provide a framework in which the hierarchical methods can be applied also for constrained clustering problems  $\Phi_k(\mathbf{U}) \subset P_k(\mathbf{U})$ .

In the ordinary problem both predicates  $\varphi(C)$  and  $\psi(C_p, C_q)$  are always true – all conditions F1-F5 are satisfied.

9/2

#### **Criterion functions compatible with a dissimilarity between clusters**

We shall call a *dissimilarity between clusters* a function  $D : (C_1, C_2) \to \mathbb{R}_0^+$  which is symmetric, i.e.,  $D(C_1, C_2) = D(C_2, C_1)$ .

Let  $(\mathbb{R}_0^+, \oplus, 0, \leq)$  be an ordered abelian monoid. Then the criterion function  $P(\mathbb{C}) = \bigoplus_{C \in \mathbb{C}} p(C), \forall X \in \mathbb{U} : p(\{X\}) = 0$  is *compatible* with dissimilarity Dover  $\Phi$  iff for all  $C \subseteq \mathbb{U}$  holds:

$$\varphi(C) \wedge \operatorname{card}(C) > 1 \Rightarrow p(C) = \min_{(C_1, C_2) \in \Psi(C)} (p(C_1) \oplus p(C_2) \oplus D(C_1, C_2))$$

**Theorem 2.1** A S criterion function is compatible with dissimilarity D defined by

$$D(C_p, C_q) = p(C_p \cup C_q) - p(C_p) - p(C_q)$$

In this case, let  $\mathbf{C}' = \mathbf{C} \setminus \{C_p, C_q\} \cup \{C_p \cup C_q\}, C_p, C_q \in \mathbf{C}$ , then

$$P(\mathbf{C}') - P(\mathbf{C}) = D(C_p, C_q)$$

#### **Greedy approximation**

**Theorem 2.2** Let P be compatible with D over  $\Phi$ ,  $\oplus$  distributes over min, and F1 - F5 hold, then

$$P(\mathbf{C}_k^*) = \min_{\mathbf{C}\in\mathbf{\Phi}_k} P(\mathbf{C}) = \min_{\substack{C_1, C_2\in\mathbf{C}\in\mathbf{\Phi}_{k+1}\\\psi(C_1, C_2)}} \left(P(\mathbf{C})\oplus D(C_1, C_2)\right)$$

The equality from theorem 2.1 can also be written in the form

$$P(\mathbf{C}_k^*) = \min_{\mathbf{C}\in\Phi_{k+1}} \left( P(\mathbf{C}) \oplus \min_{\substack{C_1, C_2\in\mathbf{C}\\\psi(C_1, C_2)}} D(C_1, C_2) \right)$$

from where we can see the following 'greedy' approximation:

$$P(\mathbf{C}_{k}^{*}) \approx P(\mathbf{C}_{k+1}^{*}) \oplus \min_{\substack{C_{1}, C_{2} \in \mathbf{C}_{k+1}^{*} \\ \psi(C_{1}, C_{2})}} D(C_{1}, C_{2})$$

which is the basis for the following agglomerative (binary) procedure for solving the clustering problem.

#### **Agglomerative methods**

1. 
$$k := n; \mathbf{C}(k) := \{\{\mathbf{X}\} : \mathbf{X} \in \mathbf{U}\};\$$

2. while 
$$\exists C_i, C_j \in \mathbf{C}(k)$$
:  $(i \neq j \land \psi(C_i, C_j))$  repeat

- 2.1.  $(C_p, C_q) := \operatorname{argmin}\{D(C_i, C_j): i \neq j \land \psi(C_i, C_j)\};$
- 2.2.  $C := C_p \cup C_q; k := k 1;$

2.3. 
$$\mathbf{C}(k) := \mathbf{C}(k+1) \setminus \{C_p, C_q\} \cup \{C\};$$

2.4. determine 
$$D(C, C_s)$$
 for all  $C_s \in \mathbf{C}(k)$ 

3. 
$$m := k$$

Note that, because it is based on an approximation, this procedure is not an exact procedure for solving the clustering problem.

For another, *probabilistic* view on agglomerative methods see Kamvar, Klein, Manning (2002).

*Divisive* methods work in the reverse direction. The problem here is how to efficiently find a good split  $(C_p, C_q)$  of cluster C.

#### **Some dissimilarities between clusters**

We shall use the generalized Ward's c.e.f.

$$p(C) = \frac{1}{2w(C)} \sum_{X,Y \in C} w(X) \cdot w(Y) \cdot d(X,Y)$$

and the notion of the *generalized center*  $\overline{C}$  of the cluster C, for which the dissimilarity to any cluster or unit U is defined by

$$d(U,\overline{C}) = d(\overline{C},U) = \frac{1}{w(C)} (\sum_{X \in C} w(X) \cdot d(X,U) - p(C))$$

$$\begin{aligned} \text{Minimal: } D^m(C_u,C_v) &= \min_{X \in C_u,Y \in C_v} d(X,Y) \\ \text{Maximal: } D^M(C_u,C_v) &= \max_{X \in C_u,Y \in C_v} d(X,Y) \\ \text{Average: } D^a(C_u,C_v) &= \frac{1}{w(C_u)w(C_v)} \sum_{X \in C_u,Y \in C_v} w(X) \cdot w(Y) \cdot d(X,Y) \end{aligned}$$

#### ... some dissimilarities

Gower-Bock: 
$$D^G(C_u, C_v) = d(\overline{C}_u, \overline{C}_v) = D^a(C_u, C_v) - \frac{p(C_u)}{w(C_u)} - \frac{p(C_v)}{w(C_v)}$$
  
Ward:  $D^W(C_u, C_v) = \frac{w(C_u)w(C_v)}{w(C_u \cup C_v)}D^G(\overline{C}_u, \overline{C}_v)$   
Inertia:  $D^I(C_u, C_v) = p(C_u \cup C_v)$   
Variance:  $D^V(C_u, C_v) = \operatorname{var}(C_u \cup C_v) = \frac{p(C_u \cup C_v)}{w(C_u \cup C_v)}$ 

Weighted increase of variance:

$$D^{v}(C_{u}, C_{v}) = \operatorname{var}(C_{u} \cup C_{v}) - \frac{w(C_{u}) \cdot \operatorname{var}(C_{u}) + w(C_{v}) \cdot \operatorname{var}(C_{v})}{w(C_{u} \cup C_{v})} = \frac{D^{W}(C_{u}, C_{v})}{w(C_{u} \cup C_{v})}$$

For all of them *Lance-Williams-Jambu formula* holds:

$$D(C_p \cup C_q, C_s) = \alpha_1 D(C_p, C_s) + \alpha_2 D(C_q, C_s) + \beta D(C_p, C_q) + \gamma |D(C_p, C_s) - D(C_q, C_s)| + \delta_1 v(C_p) + \delta_2 v(C_q) + \delta_3 v(C_s)$$

#### Lance-Williams-Jambu coefficients

method	$lpha_1$	$lpha_2$	eta	$\gamma$	${\delta}_t$	$v(C_t)$
minimum	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$	0	_
maximum	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$	0	—
average	$rac{w_p}{w_{pq}}$	$rac{w_{q}}{w_{pq}}$	0	0	0	—
Gower-Bock	$rac{w_p}{w_{pq}}$	$\frac{w_{q}}{w_{pq}}$	$-rac{w_pw_q}{w_{pq}^2}$	0	0	_
Ward	$\frac{w_{ps}}{w_{pqs}}$	$\frac{w_{qs}}{w_{pqs}}$	$-\frac{w_s}{w_{pqs}}$	0	0	—
inertia	$\frac{w_{ps}}{w_{pqs}}$	$\frac{w_{qs}}{w_{pqs}}$	$\frac{w_{pq}}{w_{pqs}}$	0	$-rac{w_t}{w_{pqs}}$	$p(C_t)$
variance	$\frac{w_{ps}^2}{w_{pqs}^2}$	$\frac{w_{qs}^2}{w_{pqs}^2}$	$\frac{w_{pq}^2}{w_{pqs}^2}$	0	$-\frac{w_t}{w_{pqs}^2}$	$p(C_t)$
w.i. variance	$\frac{w_{ps}^2}{w_{pqs}^2}$	$\frac{w_{qs}^2}{w_{pqs}^2}$	$-\frac{w_s w_{pq}}{w_{pqs}^2}$	0	0	_

 $w_p = w(C_p), w_{pq} = w(C_p \cup C_q), w_{pqs} = w(C_p \cup C_q \cup C_s)$ 

#### **Hierarchies**

The agglomerative clustering procedure produces a series of feasible clusterings C(n), C(n-1), ..., C(m) with  $C(m) \in Max \Phi$  (maximal elements for  $\sqsubseteq$ ).

Their union  $\mathcal{T} = \bigcup_{k=m}^{n} \mathbf{C}(k)$  is called a *hierarchy* and has the property

 $\forall C_p, C_q \in \mathcal{T} : C_p \cap C_q \in \{\emptyset, C_p, C_q\}$ 

The set inclusion  $\subseteq$  is a *tree* or *hierarchical* order on  $\mathcal{T}$ . The hierarchy  $\mathcal{T}$  is *complete* iff  $\mathbf{U} \in \mathcal{T}$ .

For  $W \subseteq U$  we define the *smallest cluster*  $C_{\mathcal{T}}(W)$  from  $\mathcal{T}$  containing W as: c1.  $W \subseteq C_{\mathcal{T}}(W)$ 

c2.  $\forall C \in \mathcal{T} : (W \subseteq C \Rightarrow C_{\mathcal{T}}(W) \subseteq C)$ 

 $C_{\mathcal{T}}$  is a *closure* on  $\mathcal{T}$  with a special property

 $\mathbf{Z} \notin C_{\mathcal{T}}(\{\mathbf{X},\mathbf{Y}\}) \Rightarrow C_{\mathcal{T}}(\{\mathbf{X},\mathbf{Y}\}) \subset C_{\mathcal{T}}(\{\mathbf{X},\mathbf{Y},\mathbf{Z}\}) = C_{\mathcal{T}}(\{\mathbf{X},\mathbf{Z}\}) = C_{\mathcal{T}}(\{\mathbf{Y},\mathbf{Z}\})$ 

#### **Level functions**

- A mapping  $h : \mathcal{T} \to \mathbb{R}_0^+$  is a *level function* on  $\mathcal{T}$  iff
- 11.  $\forall \mathbf{X} \in \mathbf{U} : h(\{\mathbf{X}\}) = 0$
- 12.  $C_p \subseteq C_q \Rightarrow h(C_p) \le h(C_q)$

A simple example of level function is h(C) = card(C) - 1.

Every hierarchy / level function determines an ultrametric dissimilarity on  ${f U}$ 

$$\delta(\mathbf{X}, \mathbf{Y}) = h(C_{\mathcal{T}}(\{\mathbf{X}, \mathbf{Y}\}))$$

The converse is also true (see Dieudonne (1960)): Let d be an ultrametric on U. Denote  $\overline{B}(X, r) = \{Y \in U : d(X, Y) \le r\}$ . Then for any given set  $A \subset \mathbb{R}^+$  the set

 $\mathbf{C}(A) = \{ \overline{B}(\mathbf{X}, r) : \mathbf{X} \in \mathbf{U}, r \in A \} \cup \{ \{ \mathbf{U} \} \} \cup \{ \{ \mathbf{X} \} : \mathbf{X} \in \mathbf{U} \}$ 

is a complete hierarchy, and  $h(C) = \operatorname{diam}(C)$  is a level function.

The pair  $(\mathcal{T}, h)$  is called a *dendrogram* or a *clustering tree* because it can be visualized as a tree.



#### Inversions

Unfortunately the function  $h_D(C) = D(C_p, C_q)$ ,  $C = C_p \cup C_q$  is not always a level function – for some Ds the *inversion*s,  $D(C_p, C_q) > D(C_p \cup C_q, C_s)$ , are possible. Batagelj (1981) showed:

**Theorem 2.3**  $h_D$  is a level function for the Lance-Williams procedure  $(\alpha_1, \alpha_2, \beta, \gamma)$  *iff:* 

- (i)  $\gamma + \min(\alpha_1, \alpha_2) \ge 0$
- $(ii) \qquad \alpha_1 + \alpha_2 \ge 0$
- $(iii) \qquad \alpha_1 + \alpha_2 + \beta \ge 1$

The dissimilarity D has the *reducibility* property (Bruynooghe, 1977) iff

 $D(C_p, C_q) \le t, \ D(C_p, C_s) \ge t, \ D(C_q, C_s) \ge t \ \Rightarrow \ D(C_p \cup C_q, C_s) \ge t$ 

**Theorem 2.4** If a dissimilarity D has the reducibility property then  $h_D$  is a level function.

#### **Adding hierarchical methods**

Suppose that we already built a clustering tree  $\mathcal{T}$  over the set of units U. To add a new unit X to the tree  $\mathcal{T}$  we start in the root and branch down. Assume that we reached the node corresponding to cluster C, which was obtained by joining subclusters  $C_p$  and  $C_q$ . There are three possibilities: or to add X to  $C_p$ , or to add X to  $C_q$ , or to form a new cluster  $\{X\}$ .

Consider again the 'greedy approximation'  $P(\mathbf{C}_{k}^{\bullet}) = P(\mathbf{C}_{k+1}^{\bullet}) + D(C_{p}, C_{q})$  where  $D(C_{p}, C_{q}) = \min_{C_{u}, C_{v} \in \mathbf{C}_{k+1}^{\bullet}} D(C_{u}, C_{v})$  and  $\mathbf{C}_{i}^{\bullet}$  are greedy solutions.

Since we wish to minimize the value of criterion P it follows from the greedy relation that we have to select the case corresponding to the maximal among values  $D(C_p \cup \{X\}, C_q), D(C_q \cup \{X\}, C_p)$  and  $D(C_p \cup C_q, \{X\}).$ 

This is a basis for the adding clustering method. We start with a tree on the first two units and then successively add to it the remaining units. The unit X is included into all clusters through which we branch it down.

20/2



# About the minimal solutions of $(P_k, SR)$

**Theorem 2.5** In the (locally with respect to transitions) minimal clustering for the problem  $(P_k, SR)$ 

SR. 
$$P(\mathbf{C}) = \sum_{C \in \mathbf{C}} \sum_{\mathbf{X} \in C} w(\mathbf{X}) \cdot d(\mathbf{X}, \overline{C})$$

each unit is assigned to the nearest representative: Let  $C^{\bullet}$  be (locally with respect to transitions) minimal clustering then it holds:

 $\forall C_u \in \mathbf{C}^{\bullet} \forall \mathbf{X} \in C_u \forall C_v \in \mathbf{C}^{\bullet} \setminus \{C_u\} : d(\mathbf{X}, \overline{C}_u) \le d(\mathbf{X}, \overline{C}_v)$ 

#### Proof

Let  $\mathbf{C}' = (\mathbf{C}^{\bullet} \setminus \{C_u, C_v\}) \cup \{C_u \setminus \{X\}, C_v \cup \{X\}\}$  be any clustering neighbouring with respect to transitions to the clustering  $\mathbf{C}^{\bullet}$ . From the theorem assumptions  $P(\mathbf{C}^{\bullet}) \leq P(\mathbf{C}')$  and the type of criterion function we have:

$$p(C_u) + p(C_v) \le p(C_u \setminus X) + p(C_v \cup X)$$

and by proposition 1.4.b:  $\leq p(C_u) - w(X).d(X, \overline{C}_u) + p(C_v \cup X).$ Therefore  $p(C_v) \leq p(C_v \cup X) - w(X).d(X, \overline{C}_u)$ , and

$$w(\mathbf{X}).d(\mathbf{X},\overline{C}_{u}) \leq p(C_{v} \cup \mathbf{X}) - p(C_{v}) =$$
  
=  $p(C_{v} \cup \mathbf{X}) - (p(C_{v}) + w(\mathbf{X}).d(\mathbf{X},\overline{C}_{v})) + w(\mathbf{X}).d(\mathbf{X},\overline{C}_{v})$   
=  $w(\mathbf{X}).d(\mathbf{X},\overline{C}_{v}) + (p(C_{v} \cup \mathbf{X}) - \sum_{\mathbf{Y} \in C_{v} \cup \mathbf{X}} w(\mathbf{Y}).d(\mathbf{Y},\overline{C}_{v}))$ 

By the definition of cluster-error function of type R the second term in the last line is negative. Therefore

$$\leq w(\mathbf{X}).d(\mathbf{X}, \overline{C}_v)$$

Dividing by w(X) > 0 we finally get

$$d(\mathbf{X}, \overline{C}_u) \le d(\mathbf{X}, \overline{C}_v)$$

## **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  $C \in \mathbf{C}$  its leader  $\overline{C}$ ; the new clustering  $\mathbf{C}$  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:

$\Lambda$	– set of <i>representatives</i>
$\mathrm{L}\subseteq\Lambda$	- representation
$\Psi$	- set of <i>feasible representations</i>
$W:\Phi\times\Psi\to\mathbb{R}_0^+$	- extended criterion function
$G:\Phi\times\Psi\to\Psi$	- representation function
$F:\Phi\times\Psi\to\Phi$	– clustering function

and

#### **Basic scheme of the dynamic clusters method**

the following conditions have to be satisfied:

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

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

W1.  $W(\mathbf{C}, G(\mathbf{C}, \mathbf{L})) \le W(\mathbf{C}, \mathbf{L})$ W2.  $W(F(\mathbf{C}, \mathbf{L}), \mathbf{L}) \le W(\mathbf{C}, \mathbf{L})$ 

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

 $\mathbf{C} := \mathbf{C}_0; \mathbf{L} := \mathbf{L}_0;$  **repeat**   $\mathbf{L} := G(\mathbf{C}, \mathbf{L});$   $\mathbf{C} := F(\mathbf{C}, \mathbf{L})$ **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

 $L_{n+1} = G(C_n, L_n)$  and  $C_{n+1} = F(C_n, L_{n+1})$ 

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 2.6** 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 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 B<sub>0</sub>: Type B with k = m + 1



The DCM sequence  $(v_n)$  is of type B if

• sets  $\Phi$  and  $\Psi$  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  $\Phi$  are finite we expect from a good dynamic clusters procedure to stabilize in finite number of steps – is of type B.

28/2

#### **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  $\Psi$  is infinite by the additional requirement:

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

With this requirement the 'symmetry' between  $\Phi$  and  $\Psi$  is distroyed. We could reestablish it by the requirement:

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

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

**Theorem 2.7** 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  $F: \Psi \to \Phi$ . In this case the condition W2 simplifies to:  $W(F(L), L) \leq W(C, L)$  which can be expressed also as  $F(L) \in \operatorname{Min}_{\mathbf{C} \in \Phi} W(\mathbf{C}, L)$ . For such, *simple* clustering functions it holds:

**Theorem 2.8** 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  $n \ge m : v_n = v_m$ .

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

**Theorem 2.9** If W3 holds and the representation function G is simple then:

a. 
$$G(\mathbf{C}) = \arg\min_{\mathbf{L}\in\Psi} W(\mathbf{C}, \mathbf{L})$$

b. 
$$\exists k, m \in \mathbb{N}, k > m \forall i \in \mathbb{N} : v_{k+i} = v_{m+i}$$

$$c. \quad \exists m \in \mathbf{N} \forall n \ge m : u_n = u_m$$

*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  $-F: \Psi \to \Phi$  and  $G: \Phi \to \Psi$ .

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

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

and

F0. 
$$F(L) \in Min_{C \in \Phi} W(C, L)$$

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.

# ... Clustering and Networks

In the next, 3. lecture we shall discuss

- clustering with relational constraint
- transforming data into graphs (neighbors)
- clustering of networks; dissimilarities between graphs (networks)
- clustering of vertices / links; dissimilarities between vertices
- clustering in large networks