# Affinity of $ODL_I^3$ classes

Hierarchical clustering techniques based on the concept of *affinity*.

*Affinity coefficients* to determine the level of similarity between two classes in different source schemas.

An affinity function *A()* is defined on top of the Common Thesaurus to evaluate the affinity of two terms.

    a strength σr is assigned to each type of relationship in the Common Thesaurus, with σr(SYN) ≥ σr(BT/NT) ≥ σr(RT)

    The affinity *A(t,t')* of two terms *t* and *t'* is equal to the highest-strength path of relationships between them, if at least one path exist, and is zero otherwise.

~ denote that two terms have affinity.

# *Name Affinity coefficient - NA*

Is the measure of the affinity between name of classes $c,c$ , if this measure exceeds a specified threshold

| Coefficient | Value | Condition |
|---|---|---|
| $NA(c,c')$ | $A(n_c, n_{c'})$ | if $A(n_c, n_{c'}) \geq \alpha$ |
| | $0$ | if $A(n_c, n_{c'}) < \alpha$ |

**Legend:**

$n_c, n_{c'}$ denote the name of $c$ and $c'$, respectively.

$\alpha$ is a threshold used to select high values of $NA(c, c')$.

# Structural Affinity Coefficient - SA

Is the measure of the level of matching of two classes *c,c* based on attribute relationships in the Common Thesaurus:

$$SA(c, c') = \frac{2 \cdot |\{(a_t, a_q) \mid a_t \in A(c), a_q \in A(c'), n_t \sim n_q\}|}{|A(c)| + |A(c')|} \cdot F_c$$

$$F_c \quad = \quad \frac{|\{x \in C \mid flag(x) = 1\}|}{|C|}$$

$$C \quad = \quad \{(a_t, a_q) \mid a_t \in A(c), a_q \in A(c'), \langle a_t \text{ SYN } a_q \rangle \text{ or } \langle a_t \text{ BT } a_q \rangle \text{ or } \langle a_t \text{ NT } a_q \rangle\}$$

**Legend**

*A(C)* is  the set of attributes in *C*

*Flag(x)=1* stands for a valid relationship in the Common Thesaurus

# *Global Affinity Coefficient - GA*

Is the measure of the affinity between classes $c,c$
computed as the weighted sum of the *NA* and *SA* coefficients

| Coefficient | Value | Condition |
|---|---|---|
| $GA(c, c')$ | $w_{NA} \cdot NA(c, c') + w_{SA} \cdot SA(c, c')$ | in all cases |

**Legend:**

$w_{NA}$ and $w_{SA}$, with $w_{NA}, w_{SA} \in [0, 1]$ and $w_{NA} + w_{SA} = 1$, are introduced

to assess the relevance of each coefficient in computing the global affinity value.

# Clustering of $ODL_I^3$ classes

The clustering phase produces an affinity tree when classes are the leaves and nodes have an associated affinity value.