

A Self-Organizing Principle

Erik Kemperman

November, 2003

About this document

This document describes an idea regarding the self-organization of data. The idea first took hold in the context of the DATA CLOUD (DC) project¹, but perhaps it can be made to be of use in other projects as well. For this reason, the document is set up in general terms. DC serves as an example, rather than as the framework of discourse.

The document is divided in three parts, three views on the principle. First, it will be explained in 'normal' English. Second, in terms of formal math and algorithms. And finally, the view will be illustrated by means of a simple example.

The principle in words

- **Data**
Unit of information. For example, a number, text, image, a/v clip, etc.
- **Metadata**
Information about information. Specifically, we will assume that data is annotated with key/value pairs called attributes. For example, metadata could include attributes such as title, author, content-type, keywords, description, etc.
- **Affinity**
A quantity (metric, number) reflecting a structural relationship between two data elements. E.g., "similarity." Typically, affinities are associated with certain attributes. E.g., "similarity in keywords." Affinity relationships are assumed to be symmetrical., i.e. $\text{affinity}(x,y) = \text{affinity}(y,x)$.
- **Dataset**
A collection of data. This is the unit of information we will consider for self-organization. For simplicity, we will assume the data and the dataset

¹DATA CLOUD is a project by Archined (<http://www.archined.nl/>), V2_Lab (<http://lab.v2.nl/>) and STEALTH group. At the time of writing, we are working on version 2.5.

do not change while organizing, but provisions can easily be made to allow dynamic changes to the dataset.

- **Self-organization**

A process in which elements (individuals) organize themselves according to certain criteria, rules or principles applied to the whole set (collective). Specifically, we'll consider a process in which each element is represented by an avatar in a virtual N -dimensional space. From random starting positions, the process of self-organization aims for these avatars to move around in this space such that the *spatial* relationships (proximity/distance) between avatars may eventually be, in some way, a meaningful reflection of *structural* relationships in the dataset.

The principle formalized

Consider a set of n data elements, $D = \{d_i | 1 \leq i \leq n\}$. Each data element d_i is represented by an avatar a_i in a virtual N -dimensional Euclidean space L , $a_i \in L \subseteq E^N$, where L is a sphere with radius $r_L : L = \{l \in E^N | \sqrt{l_1^2 + \dots + l_N^2} \leq r_L\}$. At time $t = 0$, all avatars are positioned at random coordinates, uniformly distributed over the space, $pos_i(0) = random_L$. At subsequent discrete time steps t , the position of avatar a_i is denoted as $pos_i(t)$. Let the scalar $m = \frac{n(n-1)}{2}$ be the number of distinct pairs i and j . Consider the matrix of affinities between each pair of elements $A_{i,j} = affinity(d_i, d_j)$ where $1 \leq i \leq j \leq n$, $A_{i,j} \in [0, 1]$, and $A_{i,j} = A_{j,i}$. Finally, consider a scalar constant $0 < rate \leq 1$.

At time t , the Euclidean distance in avatar space between a_i and a_j is given

by $dist_{i,j}(t) = dist_{j,i}(t) = \sqrt{\sum_{k=1}^N (pos_i(t)_k - pos_j(t)_k)^2}$. The target distance between a_i and a_j is time-invariant and is given by $dist_{i,j}^* = dist_{j,i}^* = 2r_L(1 - A_{i,j})$.

The aim of organization is to minimize the squared error $(dist_{i,j}(t) - dist_{i,j}^*)^2$ as t increases, for each pair i and j . Collectively, this is equivalent to minimizing the mean squared error $mse(t) = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n (dist_{i,j}(t) - dist_{i,j}^*)^2}{m}$ as time progresses.

At time t , let $diff_{i,j}(t) = pos_i(t) - pos_j(t)$ be the N -dimensional difference vector between the pair of avatars a_i and a_j . That is, $pos_i(t) + \frac{diff_{i,j}(t)}{2}$ is the position in the middle of the straight line between avatars a_i and a_j . Let the scalar $c = rate \frac{dist_{i,j}(t) - dist_{i,j}^*}{2r_L m}$. In other words, $c > 0$ if the current distance between the avatars is greater than the target distance, $c < 0$ if the distance is less than the target distance and $c = 0$ if the distance is just right. Now, we can move a_i and a_j such that at the next time step, c will be closer to zero:

$$pos_i(t+1) = pos_i(t) + \frac{c}{2} diff_{i,j}(t) \text{ and}$$

$$pos_j(t+1) = pos_j(t) - \frac{c}{2} diff_{i,j}(t)$$

with the limitation that avatars can not travel outside the sphere L . By repeating this step for each pair i and j in each time step t , the collective may

gradually minimize $mse(t)$. In the ideal goal state at time $t = t^*$ the mean squared error is zero, $mse(t^*) = 0$. If this state is achieved, the spatial relationships of the avatars are a perfect reflection of the structural data relationships in terms of the affinity metric. This state is stable, as $pos_i(t^* + 1) = pos_i(t^*)$ for all $1 \leq i \leq n$.

The principle at work

In this section, the principle is illustrated by a simple example. Suppose we have a set of $n = 3$ colors: $d_1 = rgb(0, 0, 0)$, $d_2 = rgb(1, 0, 0)$, and $d_3 = rgb(1, 1, 1)$. Let's define an affinity-metric for colors as $aff = 1 - \sqrt{\frac{\delta r^2 + \delta g^2 + \delta b^2}{3}}$. In other words, the $m = n(n-1)/2 = 3$ pairs have the following affinity values:

$$A_{1,2} = 1 - \sqrt{\frac{1^2 + 0^2 + 0^2}{3}} = 1 - \sqrt{\frac{1}{3}} = 0.423$$

$$A_{1,3} = 1 - \sqrt{\frac{1^2 + 1^2 + 1^2}{3}} = 0$$

$$A_{2,3} = 1 - \sqrt{\frac{0^2 + 1^2 + 1^2}{3}} = 1 - \sqrt{\frac{2}{3}} = 0.184$$

Let's use $rate = \frac{1}{2}$ and $r_L = 1$ and organize in an avatar space of $N = 2$ dimensions. At time $t = 0$ the avatars are assigned to random positions inside the sphere (circle) L . Suppose that

$$pos_1(0) = (-0.5, 0.5)$$

$$pos_2(0) = (0.5, 0.5)$$

$$pos_3(0) = (0.8, 0.1)$$

This gives us the distances

$$dist_{1,2}(0) = 1.414 \text{ (target is } dist_{1,2}^* = 2r_L(1 - A_{1,2}) = 1.155)$$

$$dist_{1,3}(0) = 1.360 \text{ (target is } dist_{1,3}^* = 2r_L(1 - A_{1,3}) = 2.000)$$

$$dist_{2,3}(0) = 0.500 \text{ (target is } dist_{2,3}^* = 2r_L(1 - A_{2,3}) = 1.633)$$

The mean square error at time $t = 0$ is equal to $mse(0) = 0.442$. By updating the positions of the avatars for each subsequent timestep t as explained in the previous section, the real distances will eventually correspond with the target distances, at time $t = t^*$:

$$dist_{1,2}(t^*) = dist_{1,2}^*$$

$$dist_{1,3}(t^*) = dist_{1,3}^*$$

$$dist_{2,3}(t^*) = dist_{2,3}^*$$

$$mse(t^*) = 0$$