



Wiley Interdisciplinary Reviews:
Computational Statistics

Centroid

Journal:	<i>Wiley Interdisciplinary Reviews: Computational Statistics</i>
Manuscript ID:	EOCS-019
Wiley - Manuscript type:	Focus article
Date Submitted by the Author:	14-Jul-2008
Complete List of Authors:	Abdi, Hervé
Keywords:	Inertia, Huyghens's theorem, center of gravity,, barycenter, center of mass,



view

Centroid 409

The University of Texas at Dallas

Hervé

Abdi

Keywords

Inertia, Center of Gravity, Barycenter, Center of Mass, Huyghens's Theorem

Abstract

The concept of centroid is the multivariate equivalent of the mean. Just like the mean, the centroid of a cloud of points minimizes the sum of the squared distances from the points of the cloud to a point in the space.

The notion of centroid generalizes the notion of a mean to multivariate analysis and multidimensional spaces. It applies to vectors instead of scalars, and it is computed by associating to each vector a mass which is a positive number taking values between 0 and 1 and such that the sum of all the masses is equal to 1. The centroid of a set of vectors is also called the *center of gravity*, the *center of mass*, or the *barycenter* of this set.

Notations and definition

Let \mathcal{V} be a set of I vectors with each vector being composed of J elements

$$\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_I\} \text{ with } \mathbf{v}_i = [v_{i,1}, \dots, v_{i,j}, \dots, v_{i,J}]^T. \quad (1)$$

To each vector is associated a mass denoted m_i for vector i . These masses take values between 0 and 1 and the sum of these masses is equal to 1. The set of masses is a vector denoted \mathbf{m} . The centroid of the set of vectors is denoted \mathbf{c} , it is defined as

$$\mathbf{c} = \sum_i^I m_i \mathbf{v}_i. \quad (2)$$

Examples

The mean of a set of numbers is the centroid of this set of observations. Here, the mass of each number is equal to the inverse of the number of observations: $m_i = \frac{1}{I}$.

For multivariate data, the notion of centroid generalizes the mean. For example, with the following three vectors:

$$\mathbf{v}_1 = \begin{bmatrix} 10 \\ 20 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} 8 \\ 24 \end{bmatrix}, \text{ and } \mathbf{v}_3 = \begin{bmatrix} 16 \\ 32 \end{bmatrix}, \quad (3)$$

and the following set of masses:

$$m_1 = \frac{1}{2}, \quad m_2 = \frac{1}{8}, \quad \text{and} \quad m_3 = \frac{3}{8}, \quad (4)$$

we obtain the following centroid

$$\mathbf{c} = \sum_i^I m_i \mathbf{v}_i = \frac{1}{2} \begin{bmatrix} 10 \\ 20 \end{bmatrix} + \frac{1}{8} \begin{bmatrix} 8 \\ 24 \end{bmatrix} + \frac{3}{8} \begin{bmatrix} 16 \\ 32 \end{bmatrix} = \begin{bmatrix} 12 \\ 25 \end{bmatrix}. \quad (5)$$

In this example, if we plot the vectors in a two-dimensional space, the centroid would be the center of gravity of the triangle made by these three vectors with masses assigned proportionally to their vector of mass. The notion of centroid can be used with spaces of any dimensionality.

Properties of the centroid

The properties of the centroid of a set of vectors closely parallel the more familiar properties of the mean of a set of numbers. Recall that a set of vectors defines a multidimensional space, and that to each multidimensional space is assigned a generalized Euclidean distance. The core property of the centroid is that: the centroid of a set of vectors minimizes the weighted sum of the generalized squared Euclidean distances from the vectors to any point in the space. This quantity which generalizes the notion of variance, is called the *Inertia* of the set of vectors relative to their centroid.

1
2
3
4
5
6
7
8 Of additional interest for multivariate analysis, the *theorem of Huyghens* indicates
9 that the weighted sum of the squared distances from a set of vectors to any vector
10 in the space can be decomposed as a weighted sum of the squared distances from
11 the vectors to their centroid plus the (weighted) squared distance from the centroid
12 to this point. In term of inertia, Huyghens's theorem states that the inertia of a set
13 of vectors to any point is equal to the inertia of the set of vectors to their centroid
14 plus the inertia of their centroid to this point. As an obvious consequence of this
15 theorem, the inertia of a set of vectors to their centroid is minimal. Huyghens't
16 theorem is the basis of several statistical methods such as Analysis of Variance
17 and Discriminant Analysis. In these techniques, the total Inertia of a cloud of
18 points is partitioned into the Inertia of the points to their group centroid (i.e., the
19 within group Inertia) plus the Inertia of the group centroids to their own centroid.
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60