

Suppose our dataset D is made up of datapoints with p attributes (i.e. p-dimensional vectors). Suppose $\mathbf{x}' = \langle x'_1, x'_2, \dots, x'_p \rangle$ is the global centroid for the dataset, i.e. the component-wise average of the all of the datapoints. The total sum of squares for the dataset D is defined as

$$TSS = \sum_{x \in D} \sum_{i=1}^p (x_i - x'_i)^2$$

In other words, it is the Euclidian distances between the data points and the centroid squared, summed over all the data points in the dataset.

For the total within sum of squares, we do the above calculation for each cluster C_i (now \mathbf{x}' is the cluster centroid):

$$TSS_i = \sum_{x \in C_i} \sum_{i=1}^p (x_i - x'_i)^2$$

We then sum over all of the clusters:

$$TWSS = \sum_{i=1}^k TSS_i$$