

# Agglomerative Hierarchical Clustering: An Introduction to Essentials. (3) Standardization, Normalization, and Dimensionality Reduction of a Data Matrix

Refat Aljumily<sup>1</sup>

<sup>1</sup> University of Newcastle

*Received: 6 December 2015 Accepted: 2 January 2016 Published: 15 January 2016*

---

## Abstract

In a previous tutorial article I looked at a proximity coefficient and, in the light of that proximity created a vectordistance matrix and used it to construct a hierarchical tree using different hierarchical clustering methods which will be the basis for exploratory multivariate analysis. The present article deals with three topics: (i) standardization for variable scales variation, (ii) normalization for sample length variation, and (iii) dimensionality reduction or minimization of data space. These techniques reflect the author's academic background and particular area of interest and are, by necessity, not a particular purpose and are straightforwardly applicable to other kinds of data, and thus to a wide range of analysis in Linguistics. My treatment of these techniques is, necessarily, introductory and brief. I hope that this article will provide practitioners with an introductory overview of these techniques used for cluster analysis of electronic corpora of linguistic data. The assumption is that the data is in the form of an  $m \times n$  matrix  $D$  in which, may require to transform it in various ways prior to cluster analyzing it. Standardized data matrix enables practitioners to measure the variation between  $n$ -variables and to cluster the cases they describe in common scales and values, regardless of their original scales and values. Normalized data matrix enables practitioners to eliminate the effect of variation in length among  $n$ -samples and to cluster them as if they were all (about) the same length, regardless of their original length. Dimensionality-reduced space data matrix enables practitioners to select and/or extract  $n$ -most interesting variables relevant to the research question and to visualize an existing pattern, regardless of the original space. A worked example is given to illustrate the effect each transformation technique has on a given data matrix. These transformation techniques have their own strengths and weakness but are beyond the scope of

---

*Index terms*— corpus, vector, matrix, standardization, coefficient of variation, normalization, dimensionality reduction.

## 1 Introduction

language corpus typically consists of more or less numerous texts each of which is described in terms of the selected linguistic features, technically known as variables. If it is to be analyzed using clustering methods, the selected variables need to be mathematically represented. A widely used way of doing this is vector space representation. Where vector space representation is used, each text is described by a vector, and the language corpus is consequently a set of vectors. Such a set of vectors is conveniently represented as a matrix in which the rows are the texts and the columns the linguistic features (variables). Thus, language corpus consisting of  $m$

### 3 TRANSFORMATION TECHNIQUES A) VARIATION OF VARIABLE SCALES

---

41 texts each of which is described by  $n$  variables is represented by an  $m \times n$  matrix  $D$  in which  $D_i$  (for  $i = 1..m$ ) is  
42 the  $i$ 'th text,  $D_j$  (for  $j = 1..n$ ) is the  $j$ 'th variable, and  $D_{ij}$  the value of variable  $j$  for text  $i$ . Once the language  
43 corpus has been constructed in a matrix, it is important to consider the issues relevant to cluster analysis of texts.  
44 Three types of issues are considered: (i) variable scales variation, (ii) text length variation, and (iii) variables  
45 selection/ extraction. This article proposes ways to remove the effect of each of these issues: (i) normalization for  
46 variation in text length, (ii) standardization for variation in variable scales, and (iii) dimensionality reduction.  
47 These techniques can be used, if it is necessary, to transform a given data matrix prior to analyzing it.

## 2 II.

### 3 Transformation Techniques a) Variation of variable scales

50 Almost any linguistic feature in a corpus such as word-forms, sentences, grammatical sequences, parts of speech,  
51 or any other easy to count features, can be measured. We use measurements to examine these linguistic features  
52 mathematically. In general, when we measure a linguistic feature, we define or interpret its properties in relation  
53 to special scales or units of measurement, then recording its happenings. That measurement constitutes the  
54 values of the linguistic features, for example: function words usage= 3000, average word-length=3, number of  
55 punctuation marks=500, diversity of words in a text 10%, and so on. Measurement is fundamental in the creation  
56 of language data because it makes a link between a particular linguistic feature in mind and an activity that  
57 originates from an individual, and thus allows the results of cluster some scale. Scales are systems designed to  
58 tell us how much of a measurable characteristic a given variable has. Scales have different types of numerical  
59 units and ranges (scales of measurements) appropriate to them which carry different amounts of information  
60 in any given application. The variables selected for describing linguistic features involving cluster analysis may  
61 require measurement on different scalars. If variables are measured on different scales, variables with large values  
62 contribute more to the distance measure than variables with small values.

63 Given an  $m \times n$  data matrix  $M$  in which the  $m$  rows represent the  $m$  objects to be clustered, the  $n$  columns  
64 represent the  $n$  variables, and the entry at  $M_{ij}$  (for  $i = 1..m, j = 1..n$ ) represents a numerical measure of object  $i$   
65 in terms of variable  $j$ , a clustering method has no idea what the values in the data matrix mean and calculates  
66 the degrees of similarity: variables that are measured in large values will have a greater influence on the degrees  
67 of similarity between the objects than those variables measured in smaller values, and, therefore, will affect the  
68 reliability of the cluster analysis. To see this, take a look at the following data matrix which describes nine  
69 students (A, B, C, D, E, F, G, H, I) in terms of their use of three linguistic features in the academic papers, one  
70 of which represents the total number of contractions, another one function word/content word ratio, and a third  
71 function words frequency. In Table/1the first column variable represents the total number of contractions, the  
72 second FW/CW ratio in percentage, and the third FW in frequency. A hierarchical cluster analysis of the matrix  
73 rows using Squared Euclidean distance gives the following dendrogram: In other words, the clustering analysis  
74 didn't find any significant clusters; there is a clear and very strong tendency to cluster by scale of measurement.  
75 The essence of the problem now is that we need a clustering structure that reveals the proximities among the  
76 vectors independent of the variation in scaling. However, there are many standardization methods as a technique  
77 for removing the effect variation in scaling among data and making each variable receives equal contribution in  
78 the cluster analysis. Some of these methods are:

- 79 ? Standard or Z-score standardization method.
- 80 ? Standardization method based on variable mean.
- 81 ? Standardization method based on variable sum.
- 82 ? Cosine standardization method.
- 83 ? Max standardization method.
- 84 ? Range standardization method.

85 One of the reasons for this diversity is that different standardization methods are required for different purposes;  
86 for clustering or for other purposes. No one single standardization method will be suitable for all applications.  
87 Some methods can be extremely useful even if they are mathematically limited. Other methods bring different  
88 benefits, although some bring disadvantages as well. To be suitable for cluster analysis, however, a method  
89 must preserve differences in variability among variables, thereby giving a true account of the intrinsic cluster  
90 structure of the unstandardized data matrix. The emphasis is the degree to which a method preserves the pre-  
91 standardization absolute magnitudes of variability. By the intrinsic variability, we mean the amount of variability  
92 in the values of a variable expressed independently of the scale of those values and measured in statistics by the  
93 coefficient of variation, which is defined with respect to a variable  $v$  as the ratio of  $v$ 's standard deviation to its  
94 mean, and by the absolute magnitude of variability we mean the amount of variation in the values of a variable  
95 expressed in terms of the scale of those values, and is measured by the standard deviation.

96 A standardization method based on variable means does this in the sense that it has the effect of preserving  
97 intrinsic variability in the values of a variable, and it does that in the following way: individual numerical column  
98 vectors of unstandardized data matrix can be standardized in relation to their mean, where the value of a given  
99 numerical column vector  $V$  in the unstandardized matrix must be divided by the mean  $\mu_V$  of column vectors:  $V$   
100  $i \text{ std} = V_i / \mu_V$  Where:

---

101 ?  $V_i$  std is a standardized column vector in a data matrix, for  $i=1$ ?number of rows in matrix or, equivalently,  
102 the number of text files in a corpus.

103 ?  $V_i$  is an unnormalized document vector, for  $i$  as above.

104 ?  $\mu V$  is the column vector mean, or scalar, measured by the total number of values in each column vector.

105 To illustrate this, the first three students described by the total number of contractions, FW/CW ratio (in  
106 percentage), and FW (in frequency), in the data matrix of Table/1 are recalculated. In Table/2, it is clear that  
107 MEAN-standardization has made the variation magnitudes comparable and also has preserved the coefficients  
108 of variation of the unstandardized variables. This is because division by a scalar, here the column vector mean,  
109 is a linear operation that alters the scale while preserving the shape of the original value distribution. It is  
110 also clear that the standard deviations of contractions, FW/CW ratio, and FW in Table 1b are identical to the  
111 corresponding coefficients of variation. This is because, for any data vector (here representing persons), it is always  
112 the case that its coefficient of variation is identical to the standard deviation of the MEANstandardized version  
113 of vector. After standardizing the variables for the remaining persons as above, the application of a hierarchical  
114 method on the standardized data matrix in Table 1b shows sufficiently accurate clustering; the hierarchical tree  
115 in Figure/2 differs substantially, and it clusters the nine students according to the relative magnitude of values  
116 in the matrix columns, i.e. regardless of the variation in the variable scales.

## 117 4 b) Normalization for variation in sample length

118 A corpus is a collection of texts collected with a particular linguistic research project. Very often, it happens  
119 that a corpus contains texts of varying sizes; many of them can be disparate in length and not at all identical  
120 with each other. If the disparity varies greatly from text to text, a critical issue arises that must be taken into  
121 account: the data abstracted from the corpus In this figure, there is a progression from the shortest texts at the  
122 top of the tree to the longest at the bottom and this means that there is a clear and very strong tendency to  
123 cluster by length. This can easily be seen from the number to the right of each of the text names which represents  
124 the number of words in the text. The reason for this is that, in the present example, the data abstracted from a  
125 corpus is based on frequency; each vector contains frequencies of lexical types for one of the texts, and a set of  
126 vectors are stored as the rows of the data matrix. In this sense, variations in the row vector lengths are simply  
127 a result of variations in magnitudes of lexical frequencies stored on the data matrix row vectors. To understand  
128 this, assume counting the number of occurrences of some lexical type  $j$  in a corpus containing two texts, A and  
129 B. Assume that  $j$  occurs 10 times equally across those two texts. After entering the lexical frequencies into data  
130 matrix row vectors, the interpretation would obviously suggest that on the basis of their usage of  $j$ , the two texts  
131 A and B are identical and that  $j$  apparently fails to discriminate between text A from text B. If, however, one  
132 knows that text A is 5000 words long and text B 500 words long, this is no longer the case. It is clear that,  
133 although both texts have the same frequency of occurrences of  $j$ , its significance level in them is significantly  
134 different from each other. The lexical type  $j$  is relatively infrequent in text A and relatively frequent in text  
135 B and therefore this difference can be used to differentiate between those texts. If we assume again that the  
136 text B is 50000 words long instead of 500, based on its observed frequency in 500 words, then there would have  
137 been 1000 occurrences of  $j$ . In short, the longer a text, the more likely in general a given word with a specific  
138 probability of occurrence is to occur in it, and, if it occurs, the higher the frequency of occurrence is in general  
139 likely to be. These different text lengths, called variations in lengths, are inherent in all texts in collections and  
140 result in variations in the frequencies stored in the data matrix. The variation may be large or very small, but it  
141 is always present. For the cluster analysis to be accurate and reliable, weighting to compensate for variation in  
142 text length is therefore necessary to remove this effect. The common way to do so is to adjust the data matrix  
143 so that not just frequency but its significance relative to text length can be represented and thus incorporated  
144 into subsequent analysis. There are a number of normalization methods that are theoretically motivated, for  
145 example: cosine normalization ? probability normalization

146 ? normalization by mean term frequency within document

147 ? normalization by maximum term frequency within document

148 ? normalization by mean document length across collection

149 ? normalization by maximum document length across collection.

150 but, the one most easy to understand is normalization by the mean document length across collection, and  
151 the reminder of discussion will concentrate on that. In this method, to adjust the lengths of each row vector  
152 of an  $m \times n$  data matrix of lexical types frequencies, the frequency count for a given lexical type in a given  
153 text must be multiplied by the mean length of all texts then divided by the total number of frequency counts  
154 occurring in that text. The effect of this process: decreasing the values in the vectors that represent long texts,  
155 increasing them in vectors that represent short ones, and, for texts that are near or at the mean, to change the  
156 corresponding vectors little or not at all. This can be expressed as: where  $X$  here in relation to mean length of  
157 texts in a corpus:  $X_i = \frac{X_i \cdot \mu}{T}$  ?

158 ?  $X_i$  is the normalized frequency of  $i$ ' th lexical type in a row vector, for  $i=1..n$ .

159 ?  $X_i$  is unnormalized frequency of  $i$ ' th lexical type in a row vector.

160 ?  $\mu$  is the mean length of vectors across all texts (T). This obtained by dividing the sum of frequencies of  
161 matrix row vectors (T) by that of the number of texts  $n$ , for  $i=1..n: \mu(T) = \frac{\sum X_i}{n} = 1$  ?  
162 ??

163 ? Length (i) is the sum of frequencies of any row vector (i).

164 For example, let M below be a matrix having The effect of the normalization method on the data matrix  
165 shown in this example above is clear: all the values in txt.a have been substantially increased because it is  
166 significantly shorter than the mean text length: length-500 <1476 (the mean). For txt.b, the values have been  
167 slightly decreased because it is slightly longer than the average document length: length-1500 >1476. Finally,  
168 the values for txt.c have been substantially decreased because it is significantly longer than the average document  
169 length: 2430 > 1476.

170 In summary, normalization enables us to cluster and compare texts with each other irrespective of their lengths  
171 and failure to normalize for variation in text length can produce fundamentally erroneous cluster analytical results.  
172 Nevertheless, the process of normalizing data matrix column or row vectors itself has some unresolved problems  
173 and these problems are not discussed here. More on document length normalization can be found in, e.g., ??Moisl,

## 174 5 c) Dimensionality reduction

175 Dimensionality is a major issue for data analysis in any given application. Where the aim is to generate a matrix  
176 M in which the rows are the data points, the column variables are lexical types, and the value at any given matrix  
177 location M ij is the frequency of lexical type j in i, dimensionality has a particular relevance to the application of  
178 cluster analysis. In dealing with highdimensional data, however, having too much is rarely a problem. Quite the  
179 opposite –the usual situation with high-dimensional data is that there is far too little. Highdimensional spaces  
180 are inherently sparse, and, to achieve adequate definition of the data manifold, the amount of data required very  
181 rapidly becomes intractably large; this phenomenon was described as the 'curse of dimensionality' by Bellman  
182 ??1961]. The solution is that data dimensionality should be kept as low as possible consistent with the need  
183 to describe the particular research project adequately. Dimensionality reduction is the process of reducing the  
184 number of redundant variables under consideration, and can be divided into two major types: variable selection  
185 and variable extraction.

186 i

## 187 6 . Variable selection methods

188 Variable selection methods try to identify a subset of the more important user-defined variables and to remove the  
189 remainder from the analysis (given some definition of importance) without losing too much information, thereby  
190 achieving dimensionality reduction. Given that variable selection methods aim to select a subset of the more  
191 important variables, a well-defined criterion of importance is fundamental. Two of the most often used ones in  
192 the literature are variable selection based on frequency and variable selection based on variance, and these are  
193 briefly described below. Others, such as variable selection based on term frequencyinverse document frequency  
194 (TF-IDF) and measures of nonrandomness, are also available, but these give results similar to those based on  
195 frequency and variance, and the additional complexity associated with them is therefore felt not to justify their  
196 inclusion; for further information on these see [e.g. Moisl, 2015;Belew, 2000; ??alton & McGill, 1983; ??obertson,  
197 2004].

### 198 7 a. Variable selection based on frequency

199 Frequency is the simplest criterion for selecting features from a data matrix: those variables which occur most  
200 often in the research domain -in the present domain, words in text -are judged to be the most important, and  
201 lost which occur least often are taken to be least important and can therefore be discarded. With respect to  
202 clustering, the fundamental idea is that a variable should represent something which occurs often enough for it  
203 to make a significant contribution to the clustering of the data vectors. To select variables based on frequency,  
204 given an m x n frequency data matrix D; the value at Dij is the number of times variable j, for j=1?n, occurs in  
205 text i, for i=1?m. The frequency of occurrence of variable j across the entire corpus of texts is then:

206 Frequencies of for all the columns data matrix D are calculated, sorted the variables in descending order of  
207 frequency, the most useful variables are selected and the less frequent variables are eliminated from D. Substantial  
208 dimensionality reduction can be achieved by applying this criterion to a data matrix D.

209 b

## 210 8 . Variable selection based on variance

211 Variability refers to the amount of variation in the values that a variable takes. Any variable x is an interpretation  
212 of some aspect of the physical world, and a value assigned to x is a measurement of the world in terms of that  
213 interpretation. If x is to describe the ages of people, it can take different values for different persons or for the  
214 same person at different times. Unless all people are exactly the same age, or the age of the same person is fixed,  
215 the values which x takes will vary substantially, and can, therefore, contribute to the distinction of people from  
216 one another, or of the age of same person at different times (i.e. the more different people groups one tests, the  
217 more variation one will see in the ages). This possibility of variability in the values assigned to variable x gives  
218 it its descriptive utility: an identical value for x tells that what x stands for in the real world does not change,  
219 moderate variability in the value tells that aspect of the world changes only a little, and widely differing values  
220 tells that it changes substantially. In general, therefore, the possibility of variability in the values assigned to

variables is necessary to the ability of variables to describe objects and thereby to represent reality. Clustering of texts or of anything else depends on there being variability in their characteristics; identical texts having the same stylistic descriptors cannot be meaningfully clustered. When the texts to be clustered are described by variables, then the variables are only useful for the purpose if there is significant variation in the values that they take. If, for example, a large number of people were described by their weights or heights, we would expect there to be logically substantial variation in values for each of them, and any cluster analysis method could legitimately be used to cluster them. On the other hand, if a large number of people were described by variables like 'eyes', 'noses', and 'legs', there would be almost no or little variation or high correlation with other features, since, with very few exceptions, everyone has two eyes and a nose, and clustering based on these variables would be effectively useless. In any clustering application, therefore, one is looking for variables with substantial variation in their values, and can ignore variables with little or no variation. Variables with no or little variation should be removed from data matrix as they contain little information and complicate cluster analysis by making the data higher-dimensionality than it needs to be [Moisl, 2015].

Mathematically, the degree of variation in the values of a variable is described by its variance. The variance of a set of variable values is the average deviation of those values from their mean. Assume a set of  $n$  values  $\{x_1, x_2, \dots, x_n\}$  assigned to a variable  $x$ . The mean of these values  $\mu$  is  $(x_1 + x_2 + \dots + x_n)/n$ . The amount by which any given value  $x_i$  differs from  $\mu$  is then  $x_i - \mu$ . The mean difference from  $\mu$  across all values is therefore  $\sum_{i=1..n} (x_i - \mu)/n$ . This mean difference of variable values from their mean almost but not quite corresponds to the definition of variance. One more step is necessary, and it is technical rather than conceptual. Because  $\mu$  is an average, some of the variable values will be greater than  $\mu$ , and some will be less. Consequently, some of the differences  $(x_i - \mu)$  will be positive and some negative. When all the  $(x_i - \mu)$  are added up, as above, they will cancel each other out. To prevent this, the  $(x_i - \mu)$  are squared. The standard definition of variance for  $n$  values  $\{x_1, x_2, \dots, x_n\}$  assigned to a variable  $x$ , therefore, is:

To show how a variance is calculated, consider the following frequency counts of six variables (the, a, she, him, then, him) occurring in the corresponding five texts. Given a data matrix  $M$  in which the row vectors are the texts and the column vectors are lexical type variables describing the texts, and also that the aim is to cluster analyze these texts on the basis of the differences among them, the application of variance/standard deviation to dimensionality reduction is straightforward: calculate and plot the variances of the columns and, if any have variability which is low in relation to that of the others, remove them on the grounds that they contribute little to differentiation of the texts, and decide on a threshold selection (the set of retained variables from each column of the data matrix).

## 9 d) Variable extraction methods

Variable extraction methods replace the set of user-defined variables with a smaller set of variables which reduces dimensionality but captures most of the variability in the original set. These methods often achieve a greater degree of dimensionality reduction, but at a cost: the newly-defined variables are generated by mathematical procedures, and their meaning relative to the research domain is typically difficult to determine reliably. There are a wide of variable extraction methods:

- ? Singular value decomposition (SVD)
- ? Principal Components Analysis (PCA)
- ? Factor Analysis (FA)? Multi-dimensional Scaling (MDS) ? Isomap ? Self-Organizing Map (SOM)

Each one of these methods can be used for dimensionality reduction as a feature or variable extractor, and to visualize the clusters as a clustering method. The literature on these methods is extensive and this is just a brief outline that one can follow. A more comprehensive account can be found in, for example, [Moisl, 2015; Borg and Groenen, 2005; Kohonen, 2001; Tenenbaum, de Silva, and Langford, 2000; Gordon, 1999]. However, it will be useful to look briefly at one of these methods, that is, PCA, as a dimensionality reduction method, to see how it reduces the data down into basic components, removing any unnecessary variables.

Principal Components Analysis (PCA) is actually a dimensionality reduction method, which aims to transform a set of correlated variables into a—usually smaller—set of uncorrelated ones. PCA can also be used for clustering if the dimensionality is sufficiently reduced. The conceptual basis of PCA is elimination of variable redundancy. Specifically, given a matrix of  $m$  data items described by  $n$  variables, principal components analysis is a technique for redescribing the  $m$  items in terms of  $k$  variables, where  $k < n$ , such that most of the variability in the original  $n$  variables is retained. When  $k = 2$  or  $k = 3$  the  $m$  data items can be plotted in two or three dimensional space and any clusters can thereby be directly perceived. Relative to an  $n$ -dimensional data set  $D$ , the essence of PCA is this: ? An  $n$ -dimensional orthogonal basis for  $D$  is constructed, such that each axis is the leastsquares best fit to one of the  $n$  directions of variation in  $D$ . ? The axes along which there is relatively little variation are eliminated, leaving an  $m$ -dimensional basis for  $D$ , where  $m < n$ . ? The original  $n$ -dimensional data  $D$  is projected into the reduced  $m$ -dimensional space, which yields a data set  $D'$  that is dimensionality-reduced but still contains most of the variability in  $D$ .

279 10 III.

280 11 Conclusion

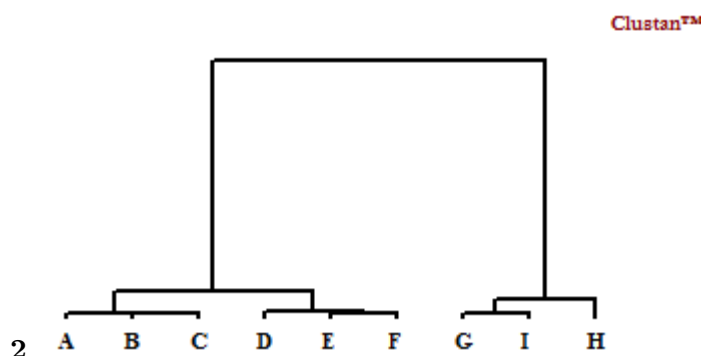
281 In this article, I discussed three techniques to adjust a data matrix before applying cluster analytical methods  
 282 to take account of the variation in scales among the variables, the variation in length among the texts, and any  
 283 superfluous variables in it using standardization, normalization, and dimensionality reduction techniques. A full  
 284 and detailed consideration of each of these techniques addressed in this article would require several articles. My  
 285 treatment of them is, necessarily, introductory and brief. Therefore, I urge interested computational linguists to  
 286 follow the more in depth sources cited in the references. The application of these techniques for cluster analysis  
 287 with specific reference to corpus linguistics is only one of many possibilities. The data items/matrix rows might  
 288 be students in a second language learning (L2) survey and the variable/matrix columns motivational factors  
 289 like learning experience, attitudes, cultural interest, and so on. n formants in a sociolinguistic or dialectological  
 290 survey and the variables/matrix columns phonetic features like voicing, and so on. The lexical frequency example  
 291 was selected because it is generic with respect to a wide range of possible applications.

IV.<sup>1</sup>



1

Figure 1: Figure 1 :



2

Figure 2: Figure 2 :

292

<sup>1</sup>© 2016 Global Journals Inc. (US)

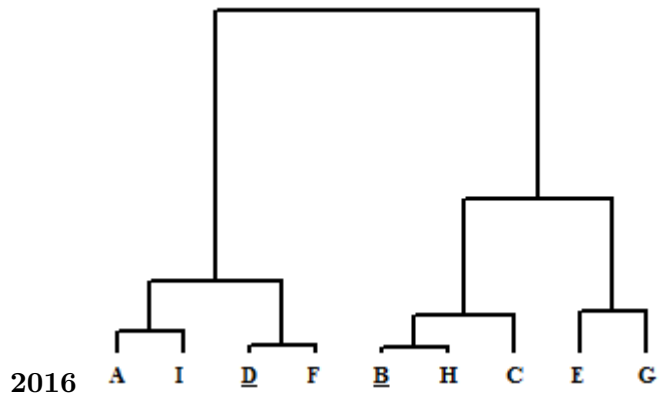


Figure 3: Year 2016 Agglomerative

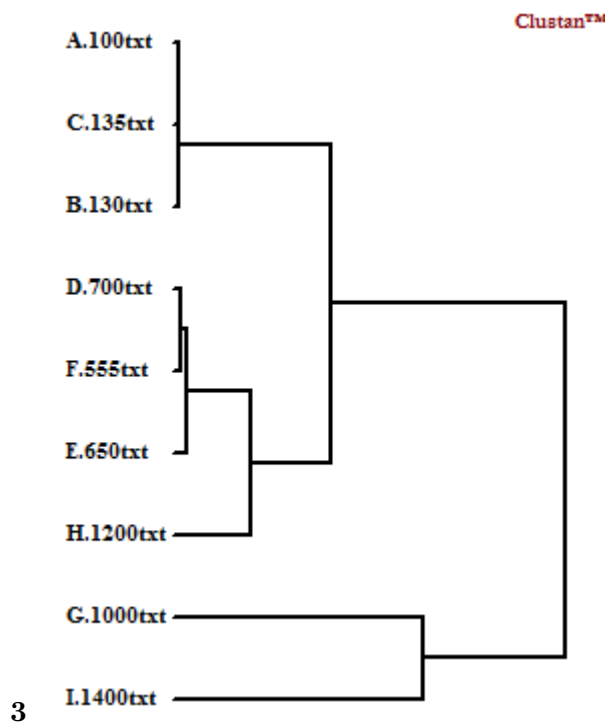


Figure 4: Figure 3 :

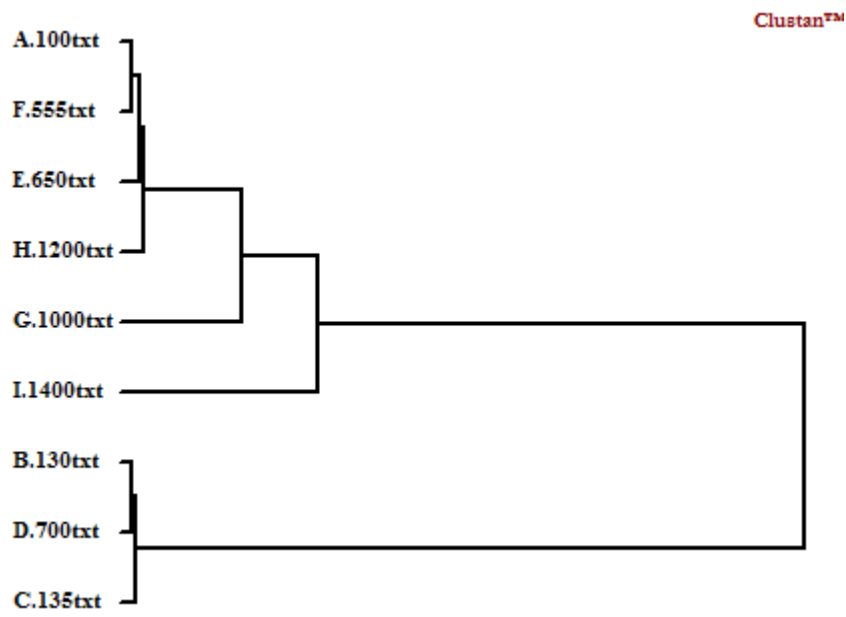


Figure 5:

$$freq(F_j) = \sum_{i=1..m} F_{i,j}$$

Figure 6:

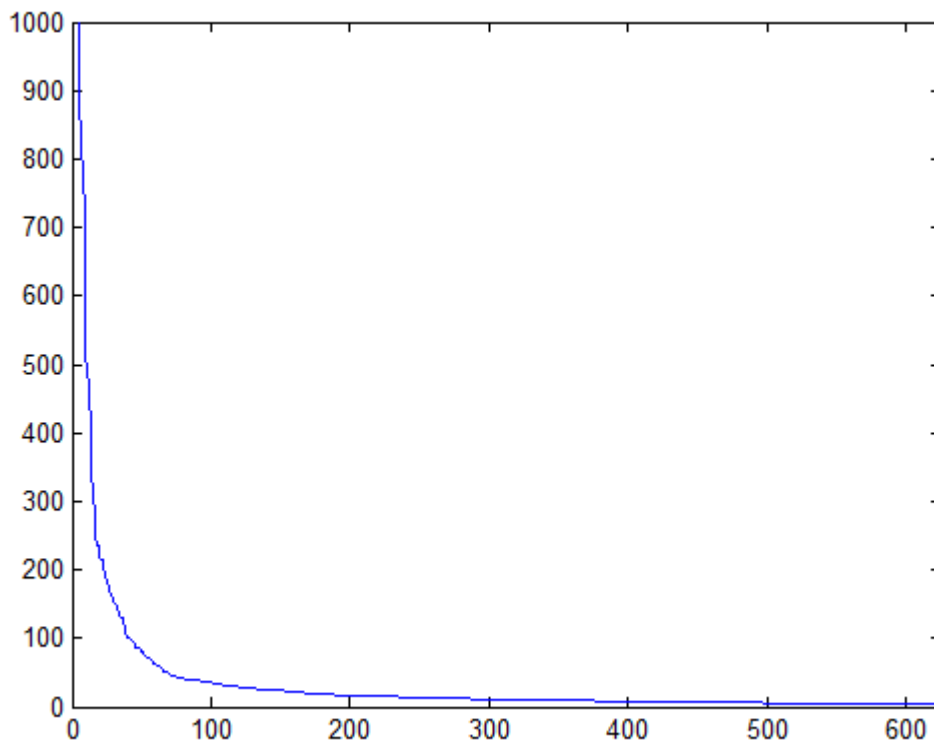


Figure 7:

---

**1**

Students	Number of contractions	FW/CW (percentage)	FW (frequency)
A	187	40	27000
B	185	35	25000
C	184	33	26000
D	170	29	23500
E	166	25	22000
F	164	26	21000
G	160	60	15000
H	150	53	10000
I	159	61	14500

Figure 8: Table 1 :

**2**

students	Contraction	FW/CW	FW	Contraction	FW/CW	FW
A	187	40	27000	1.01	1.11	1.03
B	185	35	25000	1	0.97	0.96
C	184	33	26000	0.99	0.91	1
Std	1.247	2.943	816.496	0.084	0.022	0.028
CV	0.006	0.081	0.0314	0.084	0.022	0.028

[Note: a. unSTD matrix of Table(1) b. Mean STD matrix of Table(1) ]

Figure 9: Table 2 :



## 293 .1 Acknowledgments

294 The author wishes to thank all those who dedicated their time answering my queries and providing me with  
295 valuable comments during the preparation of this study.

## 296 .2 Conflicts of Interest

297 The author declares no conflict of interest.

298 [Tenenbaum and Langford ()] ‘A global geometric framework for nonlinear dimensionality reduction’. J Tenen-  
299 baum , V , J Langford . *Science* 2000. 290 p. .

300 [Milligan and Cooper ()] ‘An examination of procedures for determining the number of clusters in a data set’. G  
301 Milligan , Cooper . *Psychometrika* 1985. 50 p. .

302 [Priddy and Keller ()] *Artificial Neural Networks: An Introduction*, K L Priddy , P E Keller . 2005. USA: Spie  
303 Press.

304 [Belew ()] R Belew . *Finding Out About: A Cognitive Perspective in Search Engine Technology and the WWW*,  
305 (Cambridge) 2000. Cambridge University Press.

306 [Borg and Groenen ()] I Borg , P Groenen . *Modern Multidimensional Scaling. 2 nd*, (Berlin) 2005. Springer.

307 [Moisl ()] *Cluster Analysis for Corpus Linguistics*, Hermann Moisl . 2015. Berlin: De Gruyter Mouton.

308 [Gordon et al. ()] ‘Data clustering: a review’. A Gordon , Chapman , A Halljain , M Murty , P Flynn . *ACM*  
309 *Computing Surveys* 1999. 1999. 31 p. . (Classification 2 nd)

310 [Singhal et al. ()] ‘Document Length Normalization’. A Singhal , G Salton , C Buckley . *Information Processing*  
311 *and Management* 1995. 32 p. .

312 [Chu et al. ()] ‘Effect of data standardization on chemical clustering and similarity searching’. C Chu , J Holliday  
313 , P Willett . *Journal of Chemical Information and Modeling* 2009. 49 p. .

314 [Dy and Bodley ()] ‘Feature selection for unsupervised learning’. J Dy , C Bodley . *Journal of Machine Learning*  
315 *Research* 2004. 5 p. .

316 [Dy (ed.) ()] *In: Computational Methods of Feature Selection*, J Dy . H. Liu and H. Motada. London: Chapman  
317 and Hall CRC (ed.) 2008. p. . (Unsupervised feature selection)

318 [Singhal et al. ()] ‘Pivoted document length normalization’. A Singhal , C Buckley , M Mitra . *Proceedings of the*  
319 *19th ACM Conference on Research and Development in Information Retrieval (SIGIR-96)*, (the 19th ACM  
320 Conference on Research and Development in Information Retrieval (SIGIR-96) 1996. p. .

321 [Kohonen ()] *Self-Organizing Maps*, T Kohonen . 2001. Berlin: Springer. (3rd ed)

322 [Gnanandesikan et al. ()] ‘Weighting and selection of variables for cluster analysis’. R Gnanandesikan , S Tsao ,  
323 J Kettenring . *Journal of Classification* 1995. 12 p. .