

# Tarea del 5 de Mayo 2020 Tanagra.

## Clustering Jerárquico.

### Parte 1

1. Aplique el algoritmo HAC 1 con las variables de input: Activo Corriente/Pasivo Corriente, Activo Corriente/ Ventas Netas.
  - a. Observe el dendrograma. ¿cuántos cluster conviene generar? \_\_\_\_\_
  - b. ¿Cuántos cluster se presentan? \_\_\_\_ ¿Cuántos elementos tiene cada cluster? \_\_\_\_\_
  - c. En View Dataset, ¿qué representa la última columna? \_\_\_\_\_
  - d. Genere un diagrama de dispersión para ambas variables, e indique en el mismo los clusters definidos.
2. Aplique el algoritmo K-means con las opciones por defecto. ¿Cuántos clusters se generaron?\_\_ ¿Cuál es la cantidad máxima de iteraciones?\_\_\_\_ ¿Se presenta el dendrograma?\_\_\_\_\_ ¿Cuál es el motivo?
3. Aplique el algoritmo K-means para generar 4 clusters. ¿Cuántos elementos tiene cada cluster? \_\_\_\_\_  
\_\_\_\_\_
4. Genere un diagrama de dispersión para representar indistintamente los clusters generados por con HAC1 o con K-means. ¿Algún algoritmo de clustering le parece preferible? Fundamente.

### Parte 2

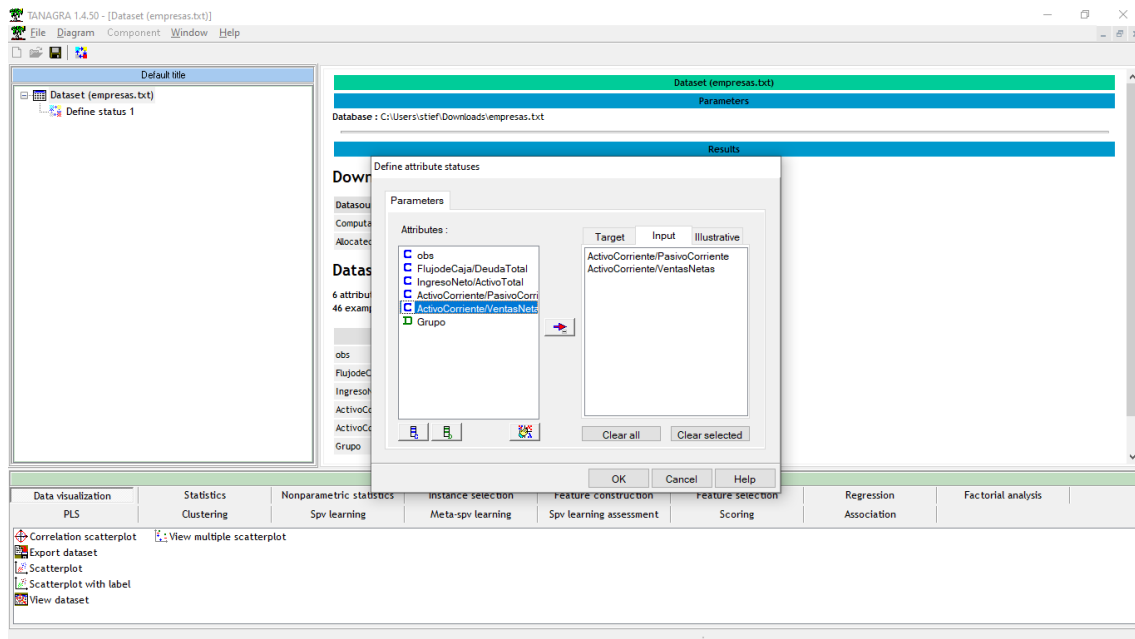
5. Repita lo realizado en la Parte 1 trabajando con las 4 variables cuantitativas como input.

1.

Para comenzar abrimos Tanagra. En “File → New File” buscamos el archivo de trabajo empresas.txt y lo cargamos. Pulsamos “Ok”.

Vamos a aplicar el algoritmo de clustering jerarquico aglomerativo . Como entradas vamos a tomar ActivoCorriente/PasivoCorriente y ActivoCorriente/VentasNetas.

Por lo tanto procedemos a crear un “Define Statues” con dichas entradas:



Como el clustering es una técnica de aprendizaje no supervisada, no hay atributo objetivo (Target).

Ahora vamos a ir a la pestaña “Clustering” y vamos a buscar el algoritmo HAC, el cual seleccionaremos y arrastraremos hasta nuestro “Define Statues”.

TANAGRA 1.4.50 - [Define status 1]

File Diagram Component Window Help

Default title

Dataset (empresas.txt)

Define status 1

HAC 1

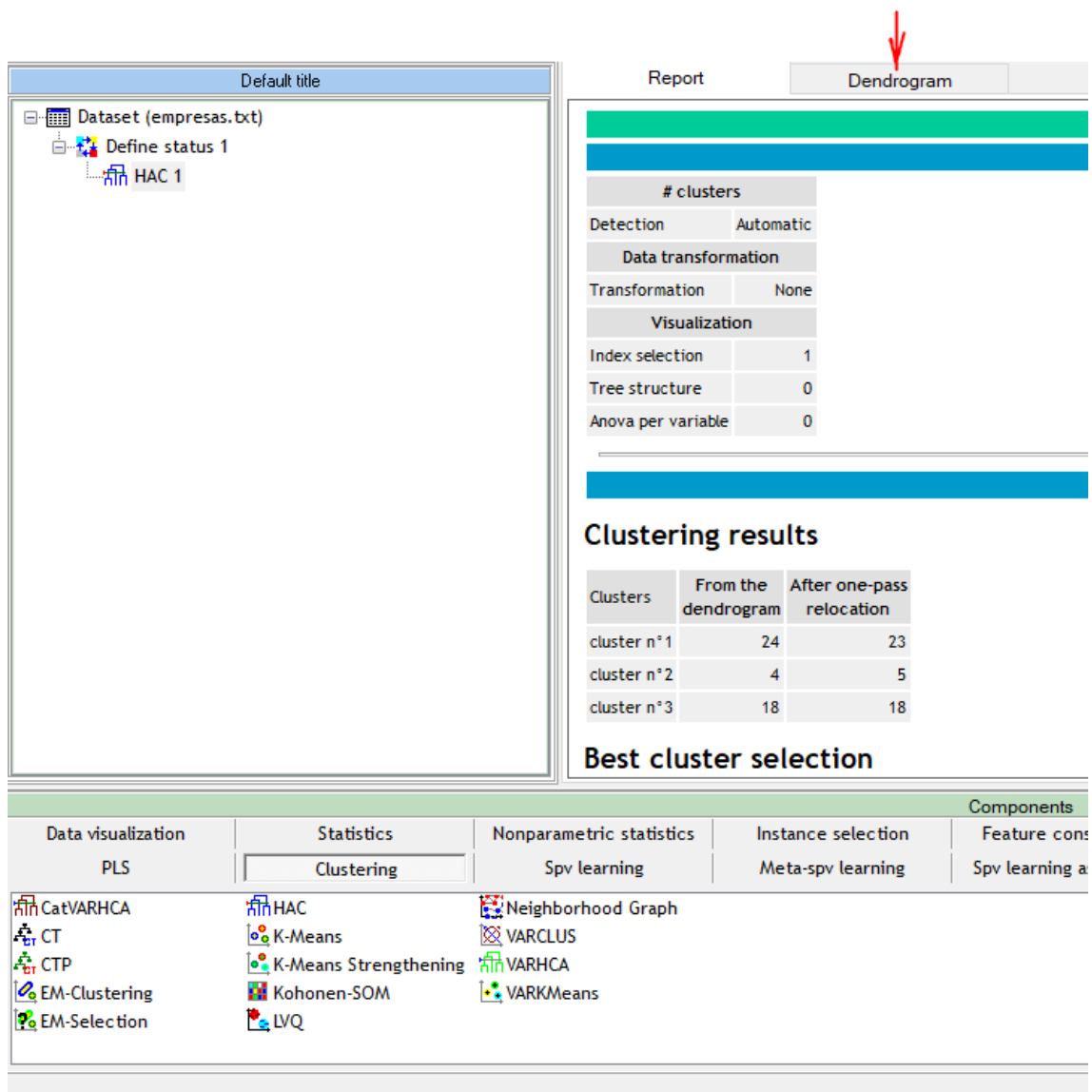
Target : 0  
Input : 2  
Illustrative : 0

Attribute	Target	Input	Illustrative
obs	-	-	-
FlujodeCaja/DeudaTotal	-	-	-
IngresoNeto/ActivoTotal	-	-	-
ActivoCorriente/PasivoCorriente	-	yes	-
ActivoCorriente/VentasNetas	-	yes	-
Grupo	-	-	-

Computation time : 0 ms.  
Created at 10/05/2020 0:46:43

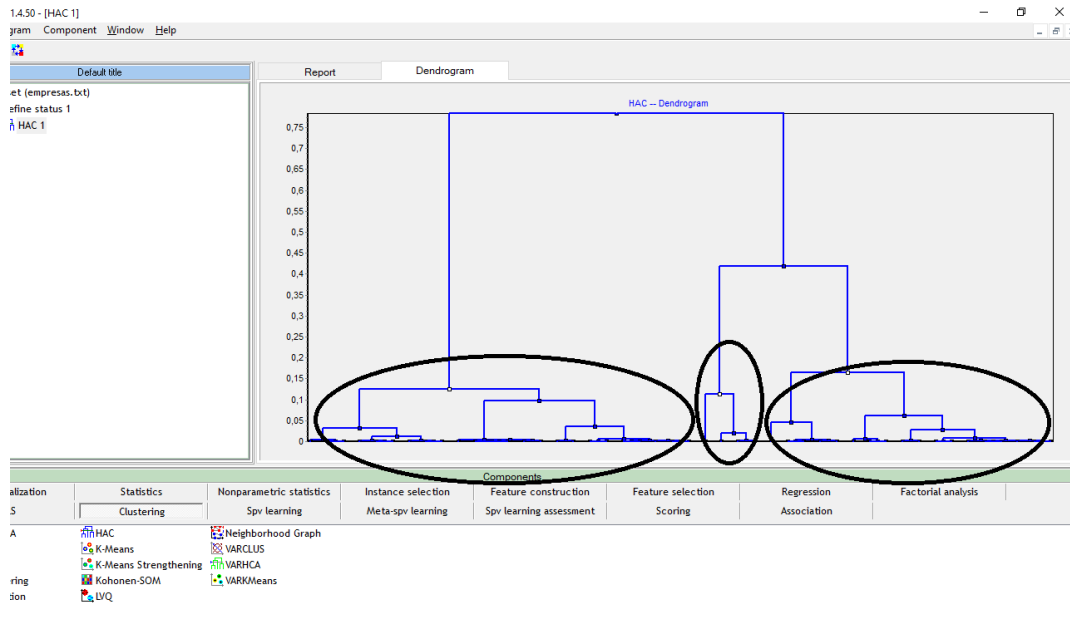
Data visualization		Statistics	Nonparametric statistics	Instance selection	Components
PLS		Clustering	Spv learning	Meta-spv learning	Spv learning a
CatVARHCA	HAC	Neighborhood Graph			
CT	K-Means	VARCLUS			
CTP	K-Means Strengthening	VARHCA			
EM-Clustering	Kohonen-SOM	VARKMeans			
EM-Selection	LVQ				

Ejecutamos el algoritmo y observamos lo siguiente:



Si pulsamos en “Dendrogram” observamos una acumulación de puntos en la parte inferior de la gráfica. Aquí se observa gráficamente que datos tienen mayor parecido entre ellos. El eje de la ordenada nos muestra la distancia entre datos. Por tanto, cuanto más cercano estén los puntos al eje de las abscisas, más parecidos serán entre ellos los datos.

Viendo dónde se concentran más los datos, podemos agrupar los mismo en tres “clusters”:



Podemos agruparlos en 3 “clusters” porque entre el “cluster” de la izquierda y el del centro, la distancia sería demasiado grande para considerarlo un solo “cluster” en vez de dos.

Para comprobar analíticamente el número de “clusters” reales, nos vamos a la pestaña “Report”, que está justo a la izquierda de la pestaña “Dendrogram”. Aquí nos vienen todos los datos generados por el algoritmo HAC:

Report

Dendrogram

HAC 1

Parameters

# clusters

Detection

Automatic

Data transformation

Transformation

None

Visualization

Index selection

1

Tree structure

0

Anova per variable

0

Results

Clustering results

Clusters	From the dendrogram	After one-pass relocation
cluster n°1	24	23
cluster n°2	4	5
cluster n°3	18	18

Best cluster selection

Aquí comprobamos que realmente se forman tres “clusters”. El algoritmo HAC detecta el mismo el número de “clusters” que se pueden formar.

El “cluster” 1 tiene 24 elementos, el “cluster” 2 tiene 4 elementos y el “cluster” 3 tiene 18 elementos.

En esta misma pantalla tenemos datos de interés como el centro de los “clusters”, que son datos importantes ya que nos muestran el dato más representativo.

TANAGRA 1.4.50 - [HAC 1]

File Diagram Component Window Help

Default title

Dataset (empresas.txt)

Define status 1

HAC 1

View dataset 1

Report Dendrogram

### Best cluster selection

Clusters	BSS ratio	Gap
1	0,0000	0,0000
2	0,3919	0,3652
3	0,6012	0,2547
4	0,6832	0,0398
5	0,7452	0,0112
6	0,8016	0,0151
7	0,8506	0,0369
8	0,8810	0,0153
9	0,9039	0,0098
10	0,9218	0,0045

### Cluster centroids

Attribute	Cluster n°1	Cluster n°2	Cluster n°3
ActivoCorriente/PasivoCorriente	1,523922	4,261280	2,067111
ActivoCorriente/VentasNetas	0,301809	0,447800	0,593378

Use GROUP CHARACTERIZATION for detailed comparisons

Computation time : 31 ms.

Components

Data visualization	Statistics	Nonparametric statistics	Instance selection	Feature construction
PLS	Clustering	Spv learning	Meta-spv learning	Spv learning assessment

Correlation scatterplot View multiple scatterplot

Export dataset

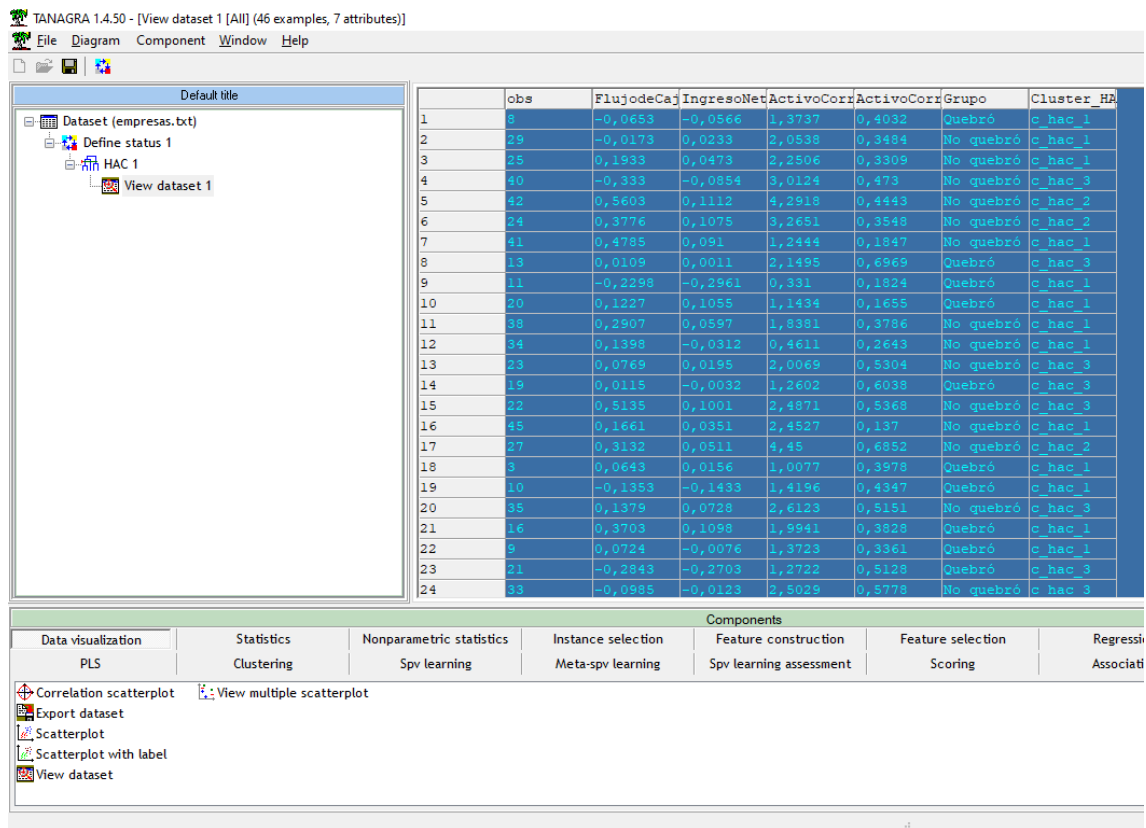
Scatterplot

Scatterplot with label

View dataset

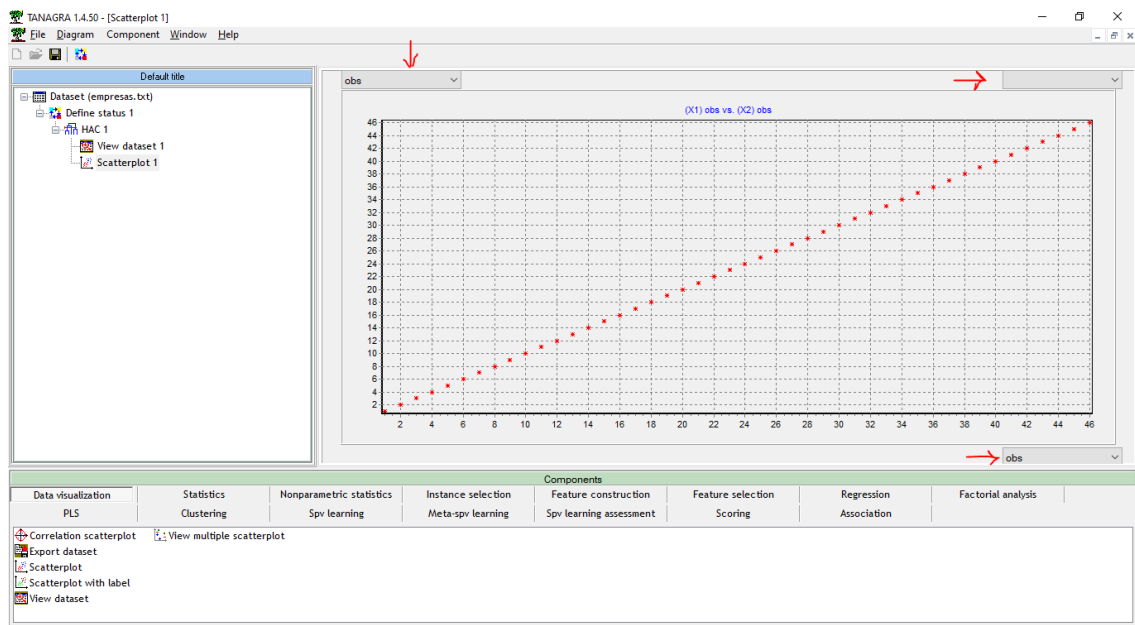
Para ver que se sitúa en la última columna nos vamos a la pestaña que dice “Data visualization” y seleccionamos “View Dataset”. Lo arrastramos hasta debajo de “HAC 1”.

Si lo ejecutamos, observamos que la última columna nos dice en que “cluster” se encuentra el dato de esa fila (cada fila es una empresa):



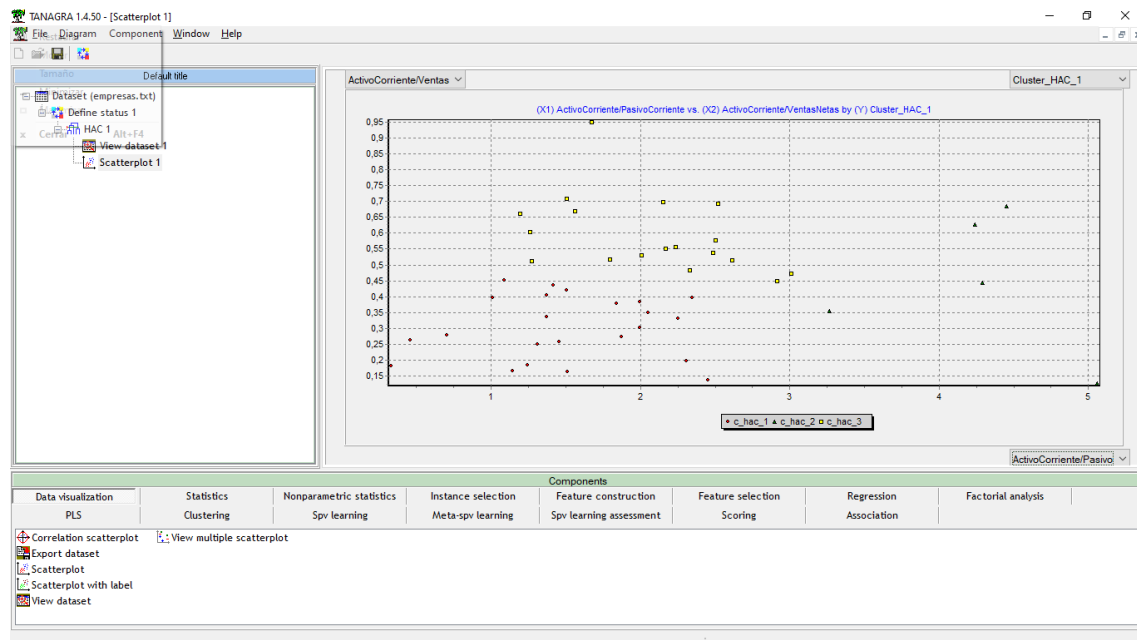
Para generar un diagrama de dispersión con esas variables, debemos ir a la pestaña “Data Visualization” y señalar la opción “Scatterplot”. Lo arrastramos hasta HAC 1.

Si lo ejecutamos observamos lo siguiente:



En cada uno de los ejes, pondremos los atributos que nosotros hemos considerado para generar el “clusters”.

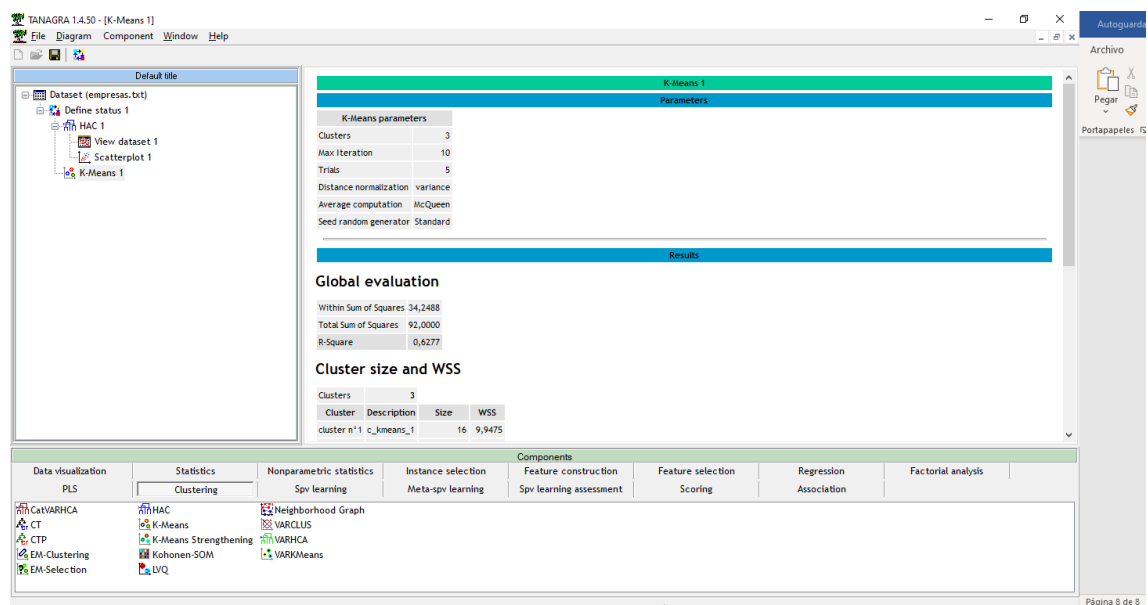
En el eje de las “x” pondremos ActivoCorriente/Pasivo corriente. En el eje de las “y” vamos a poner ActivoCorriente/VentasNetas. Por último, en la parte superior derecha seleccionaremos la opción “Cluster HAC” para generar puntos de distinto color dependiendo del “cluster” al que pertenezca.



Podemos observar que se han generado tres “clusters”. El primer “cluster” está representado por círculos rojos, el segundo con triángulos verdes y el tercero está representado por cuadrados amarillos, según reza la leyenda del gráfico.

2.

Ahora vamos a proceder a aplicar un algoritmo distinto de “clustering”, el K-Means. Para ello nos iremos a la pestaña “Clusterin” y seleccionaremos y arrastraremos hasta nuestro “Define Statues”, ejecutándolo con los parámetros por defecto.

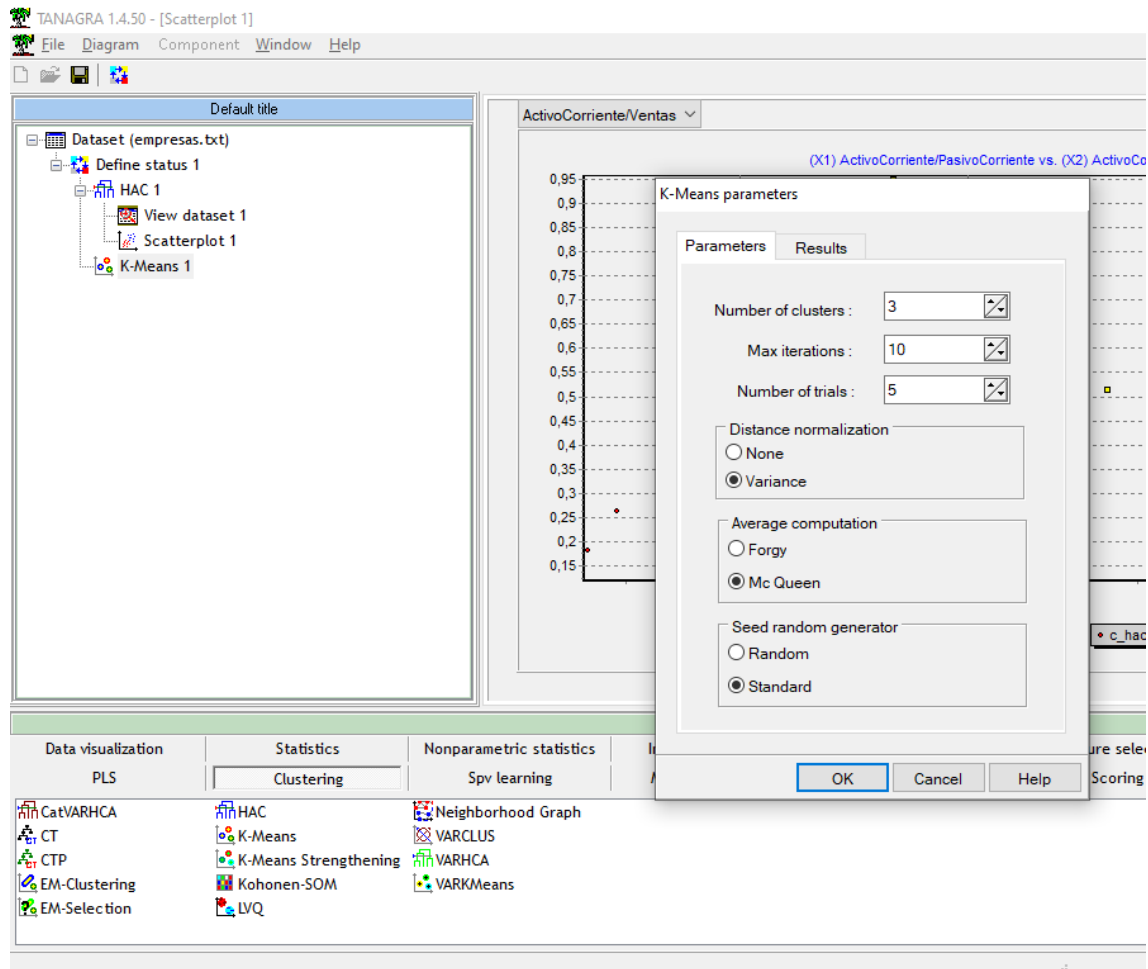




Se crean 3 “clusters”. La cantidad máxima de iteraciones son 10. Vemos que no se presenta el endrograma, ya que “K-Means” es un algoritmo de “clustering” no jerárquico y por lo tanto no hay un orden en las decisiones.

### 3.

Vamos a generar ahora 4 “clusters” con la ayuda del algoritmo “K-Means”. Para ello, nos vamos encima de “K-Means” y hacemos click derecho y pulsamos “Parameters..”.



Donde dice “Number of clusters” vamos a indicar el número de “clusters” que queremos generar (4). Tras esto pulsamos “OK” y hacemos doble click en “K-Means” para ejecutar el algoritmo:

TANAGRA 1.4.50 - [K-Means 1]

File Diagram Component Window Help

Default title

Dataset (empresas.txt)

Define status 1

HAC 1

View dataset 1

Scatterplot 1

K-Means 1

### K-Means 1

#### Parameters

K-Means parameters	
Clusters	4
Max Iteration	10
Trials	5
Distance normalization	variance
Average computation	McQueen
Seed random generator	Standard

#### Results

### Global evaluation

Within Sum of Squares	27,4808
Total Sum of Squares	92,0000
R-Square	0,7013

### Cluster size and WSS

Clusters	
4	

Cluster	Description	Size	WSS
cluster n°1 c_kmeans_1		20	8,6129

Data visualization: PLS

Statistics: Clustering

Nonparametric statistics: Spv learning

Instance selection: Meta-spv learning

Feature construction: Spv learning assessment

Feature selection: Scoring

Components: Re

CatVARHCA

CT

CTP

EM-Clustering

EM-Selection

HAC

K-Means

K-Means Strengthening

Kohonen-SOM

LVQ

Neighborhood Graph

VARCLUS

VARHCA

VARKMeans

Vemos que se han generado los 4 “clusters” que esperábamos, y que el primero de ellos se compone de 20 elementos, el segundo de 5 elementos, el tercero de 14 elementos y el cuarto de 7 elementos.

TANAGRA 1.4.50 - [K-Means 1]

File Diagram Component Window Help

Default title

Dataset (empresas.txt)

Define status 1

HAC 1

View dataset 1

Scatterplot 1

K-Means 1

### Cluster size and WSS

Clusters	
4	

Cluster	Description	Size	WSS
cluster n°1 c_kmeans_1		20	8,6129
cluster n°2 c_kmeans_2		5	7,7571
cluster n°3 c_kmeans_3		14	7,4920
cluster n°4 c_kmeans_4		7	3,6187

### R-Square for each attempt

Number of trials	
5	

Trial	R-square
1	0,626640
2	0,701296
3	0,694152
4	0,662944
5	0,663382

### Cluster centroids

Attribute	Cluster n°1	Cluster n°2	Cluster n°3	Cluster n°4
ActivoCorriente/PasivoCorriente	2,060695	4,261280	1,369800	1,695300
ActivoCorriente/VentasNetas	0,462575	0,447800	0,242279	0,711286

Data visualization: PLS

Statistics: Clustering

Nonparametric statistics: Spv learning

Instance selection: Meta-spv learning

Feature construction: Spv learning assessment

Feature selection: Scoring

Regression: Association

Factorial analysis

CatVARHCA

CT

CTP

EM-Clustering

EM-Selection

HAC

K-Means

K-Means Strengthening

Kohonen-SOM

LVQ

Neighborhood Graph

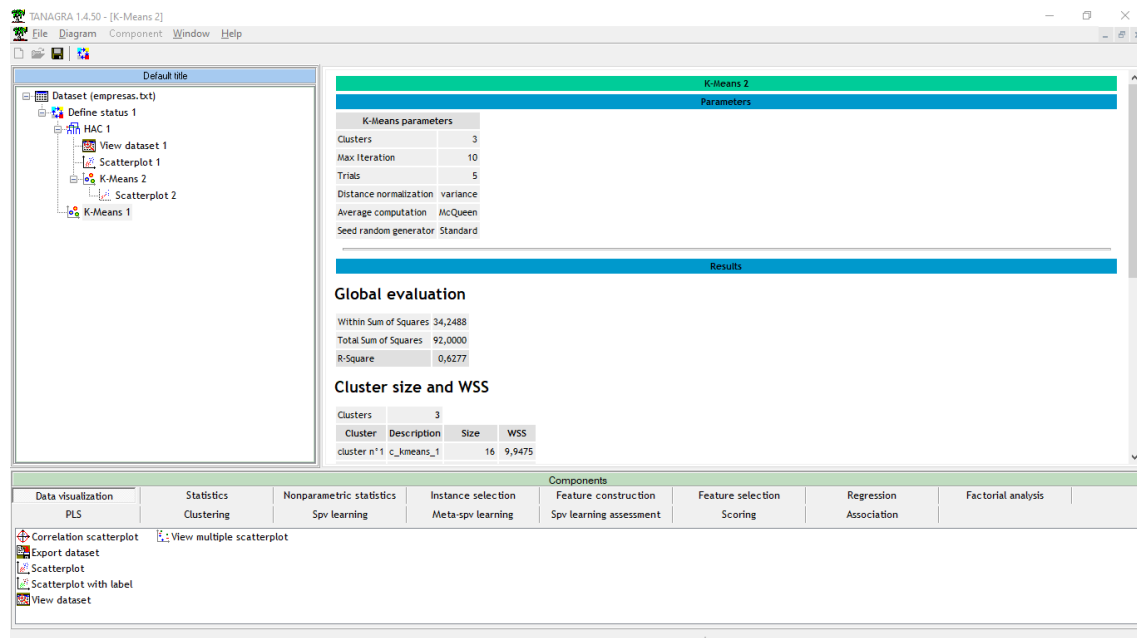
VARCLUS

VARHCA

VARKMeans

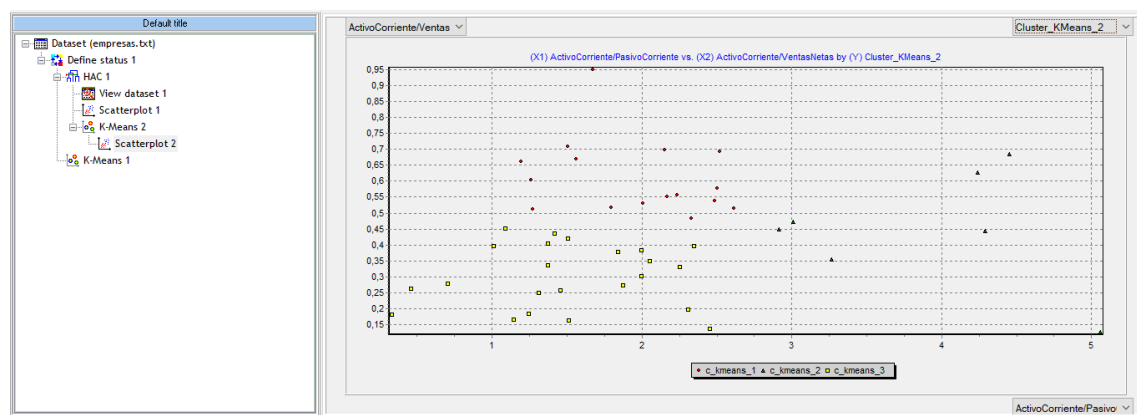
4.

Para generar de dispersión para representar indistintamente los “clusters” generados por “HAC1” o con “K-means”, debemos seleccionar en la pestaña “Clustering” el algoritmo “K-Means” y arrastrarlo hasta debajo de “HAC 1”. Tras esto lo ejecutamos haciendo doble click, y buscamos en la pestaña “Data View” el “Scatterplot”, el cual arrastramos debajo del nuevo “K-Means”.

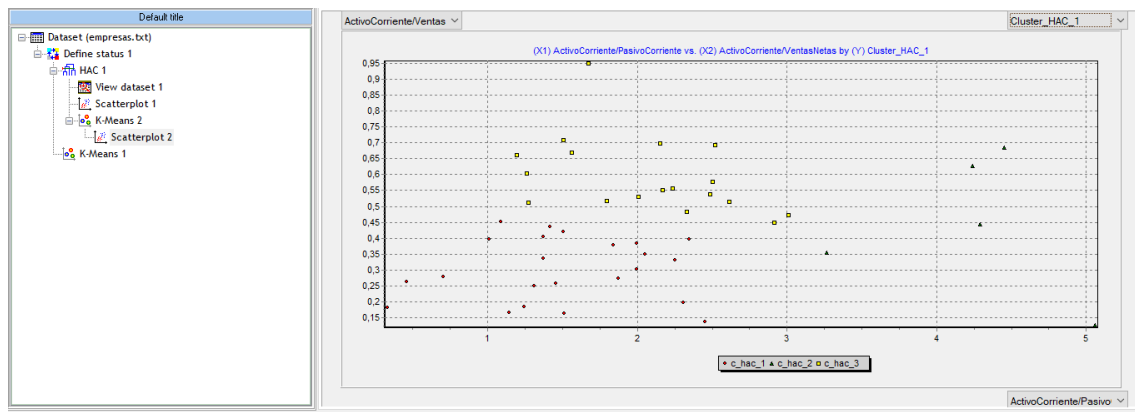


Hacemos doble click en el “Scatterplot” y procedemos a configurarlo. En el eje de ordenadas pondremos el atributo “ActivoCorriente/VentasNetas”. En el eje de las ordenadas pondremos el atributo “ActivoCorriente/PasivoCorriente”, y vemos que en la esquina superior derecha podremos elegir con cual de los dos modelos representarlo, si con el “HAC” o con el “K-Means”.

Con “K-Means”:



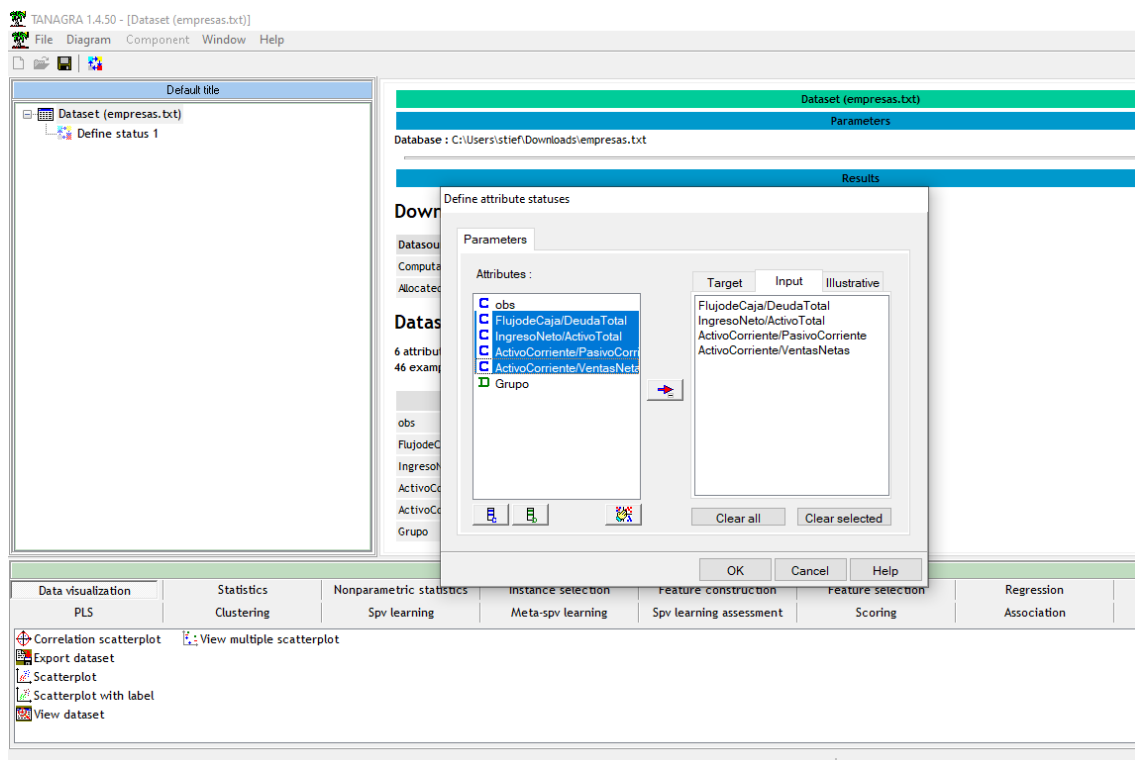
Con “HAC”:



## PARTE 2.

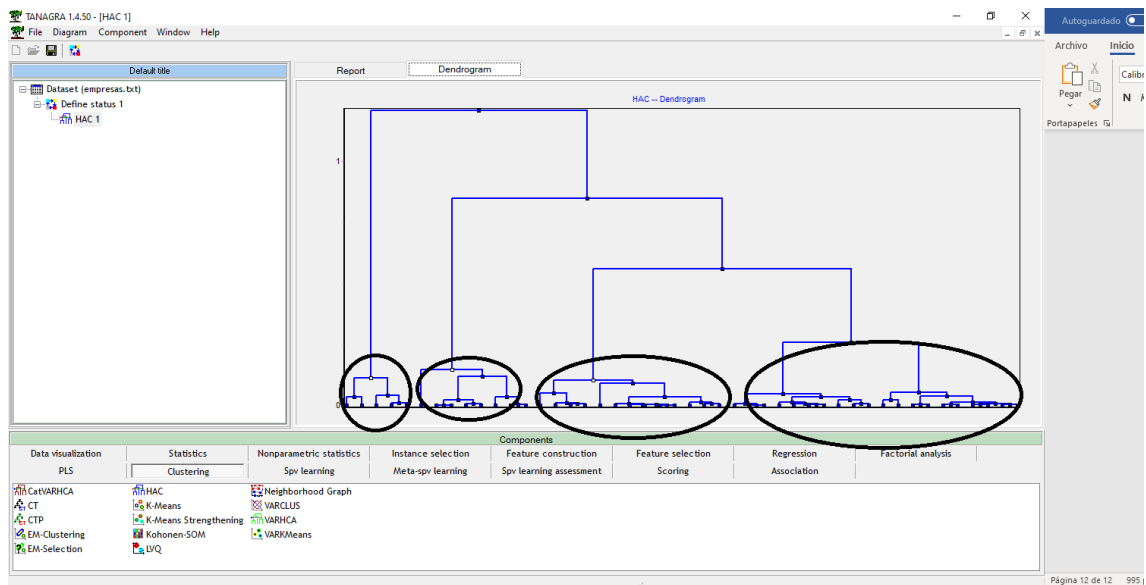
### 1.

En esta ocasión vamos a tener 4 variables como entradas. Por tanto abrimos el archivo empresas.txt, nos vamos a “Define Status” y añadimos como “Input” las 4 variables:



Tras esto vamos a la pestaña “Clustering” y seleccionamos el algoritmo HAC y arrastramos hasta nuestro “Define Status”. Ejecutamos el “HAC” y nos vamos a la pestaña “Dendogram”.

Aquí podemos observar que podríamos definir 4 “clusters”:



Si nos vamos a la pestaña “Report” podremos ver información útil. Podemos observar que realmente se han creado 4 “clusters”, el primero de ellos con 5 elementos, el segundo con 8 elementos, el tercero con 13 elementos y el cuarto con 20 elementos.

Report
Dendrogram

HAC 1

Parameters

# clusters

Automatic

Data transformation

None

Visualization

1

Index selection

0

Tree structure

0

Anova per variable

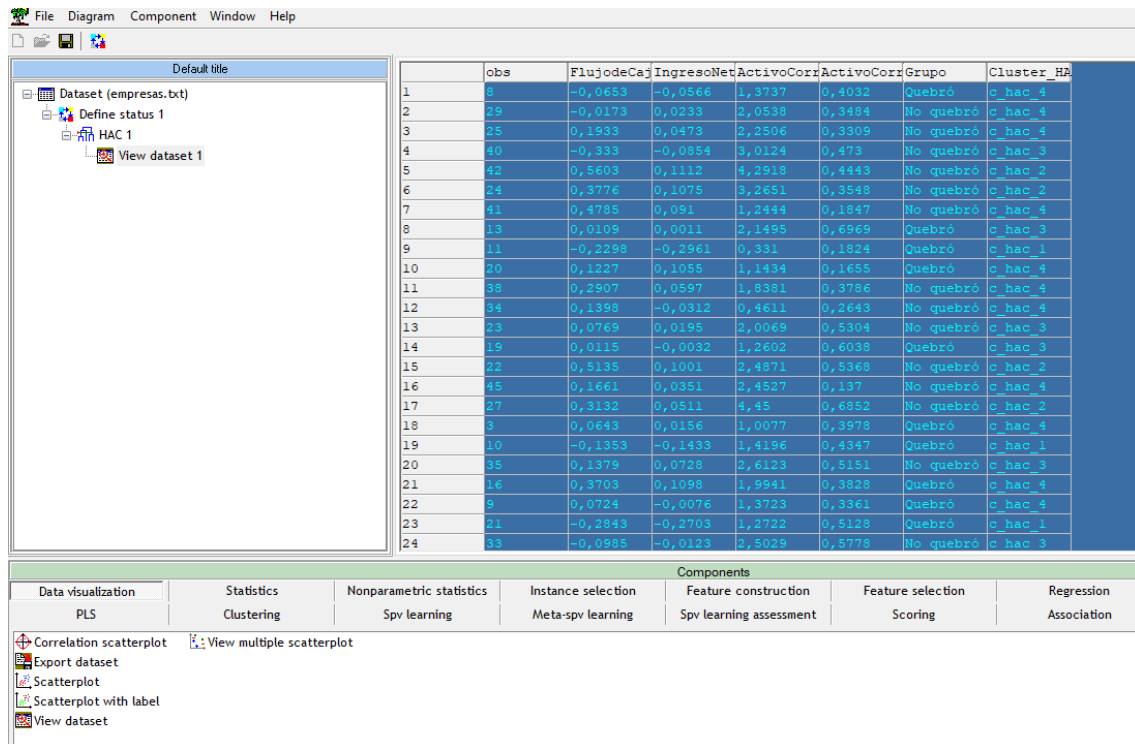
0

Results

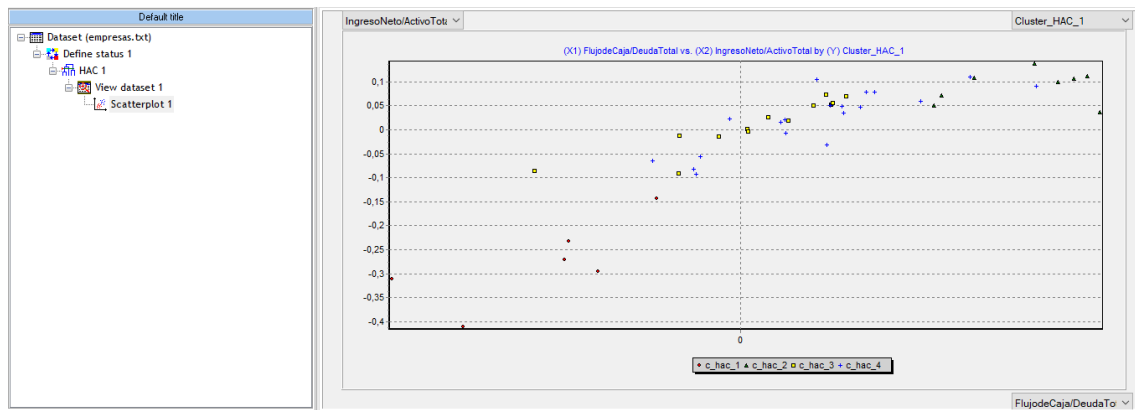
### Clustering results

Clusters	From the dendrogram	After one-pass relocation
cluster n°1	5	6
cluster n°2	8	8
cluster n°3	13	13
cluster n°4	20	19

Ahora nos vamos a la pestaña “Data Visualization”, seleccionamos “View Dataset” y lo arrastramos hasta debajo de “HAC”. Lo ejecutamos y vemos que la última columna nos indica a que “cluster” pertenece cada fila, siendo cada fila una empresa distinta.



Para generar un gráfico de dispersión usaremos el “Scatterplot” que encontramos en la pestaña “Data Visualization”. Lo seleccionamos y arrastramos hacia debajo de “HAC”. Ejecutamos y en la gráfica podremos elegir que variables queremos enfrentar para obtener un gráfico de dispersión. Por ejemplo, si seleccionamos el eje de abcisas como FlujoCaja/DeudaTotal y el eje de las ordenadas como IngresoNeto/ActivoTotal, obtendremos el siguiente gráfico de dispersión:

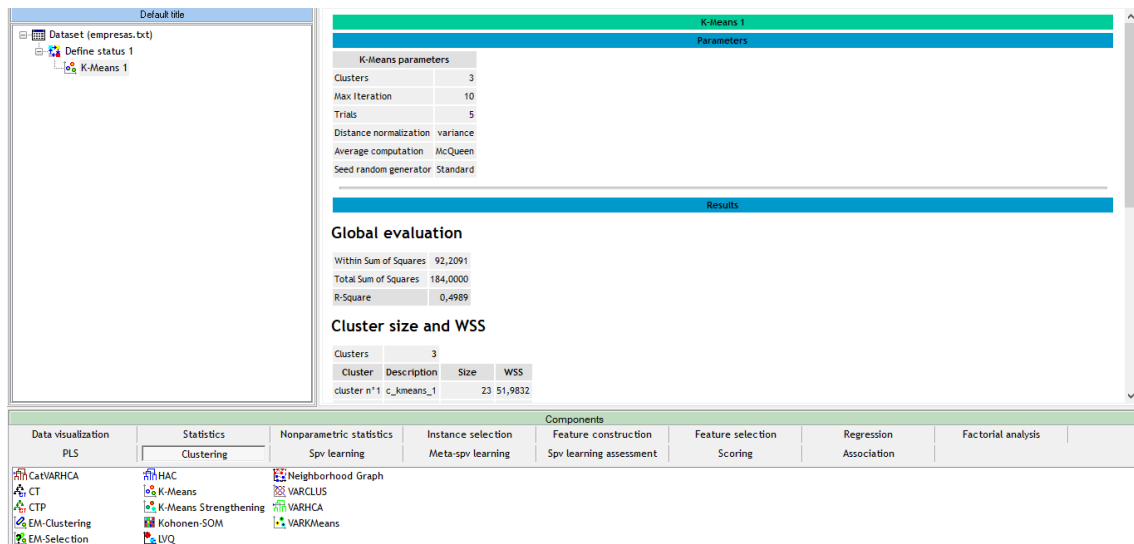


Observamos que el “cluster” uno se representa con círculos rojos, el dos con triángulos verdes, el tres con cuadrados amarillos y el cuatro con cruces azules.

Procedemos ahora a buscar el algoritmo “K-Means” en la pestaña “Clustering”. Lo seleccionamos y arrastramos debajo de nuestro “Define Status”.

## 2.

Vamos a aplicar ahora el algoritmo “K-Means”. Lo buscamos en la pestaña “Clustering”, seleccionamos y arrastramos hasta nuestro “Define Status”. Vamos a dejar los parámetros por defecto, por lo que hacemos doble click para ejecutarlo.

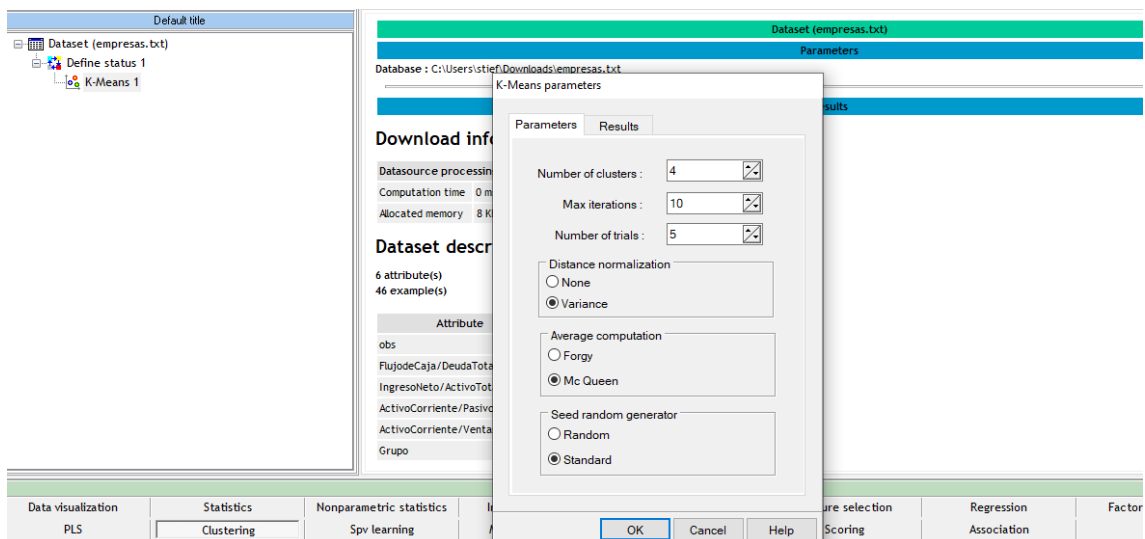


El número de “clusters” creados es 3. El máximo número de iteraciones 10.

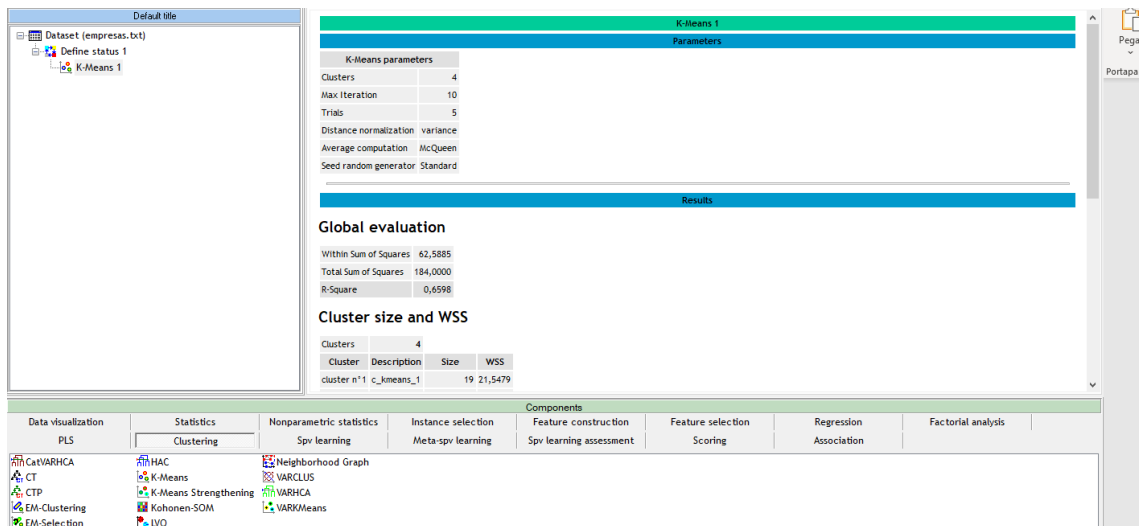
Vemos que no se presenta el endrograma, ya que “K-Means” es un algoritmo de “clustering” no jerárquico y por lo tanto no hay un orden en las decisiones.

### 3.

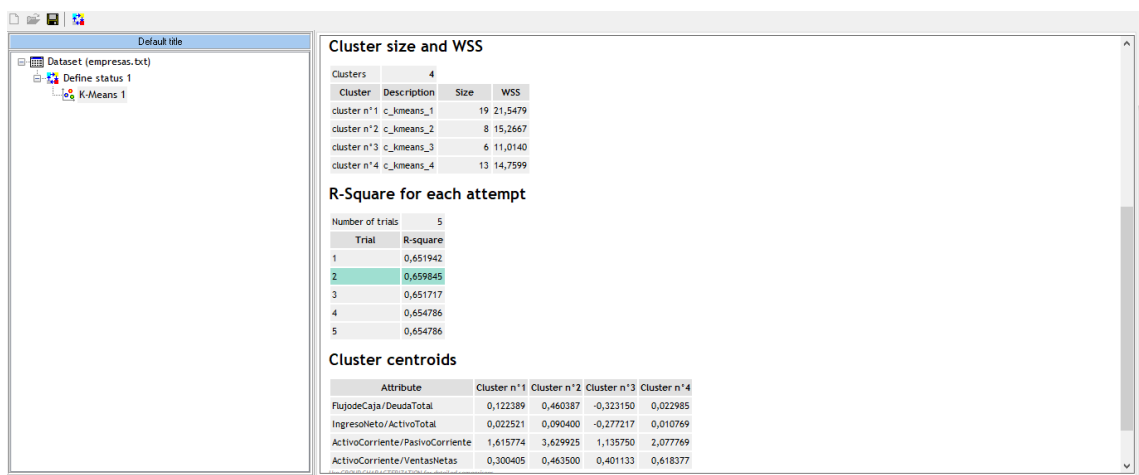
Queremos generar 4 “clusters”, por lo que hacemos click con el botón derecho sobre el algoritmo “K-Means” y seleccionamos la opción parámetros. Aquí debemos cambiar el número que indica la pestaña “Number of clusters” por 4.



Hacemos doble click para ejecutar el algoritmo y observamos la ventana emergente:



Podemos ver que se han generado 4 “clusters” tal y como habíamos determinado, y que el primero de ellos contiene 19 elementos, el segundo 8 elementos, el tercero 6 elementos y el cuarto 13 elementos.



#### 4.

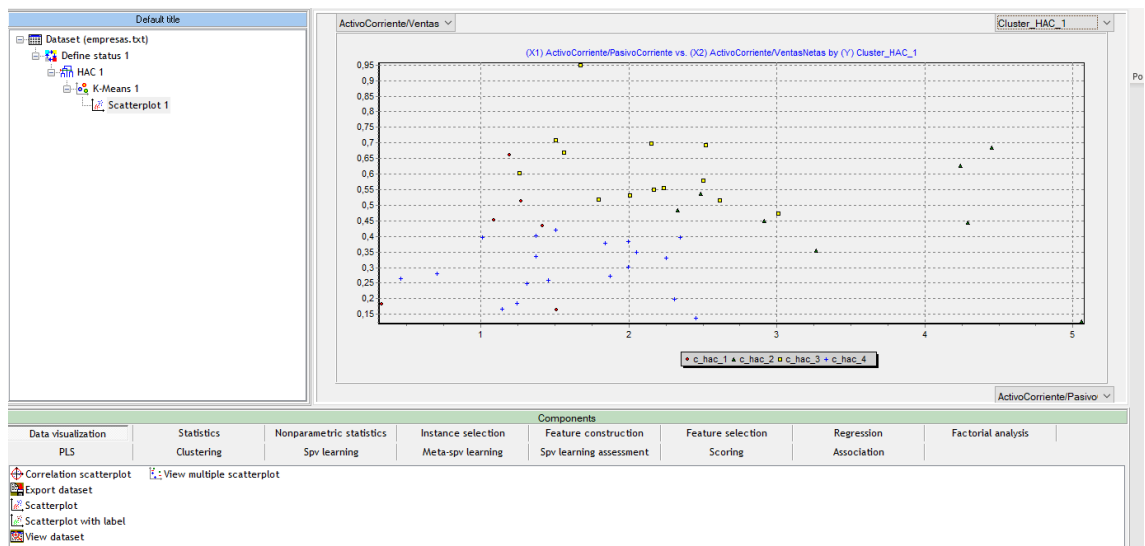
Para poder mostrar el diagrama de dispersión de cualquiera de los dos algoritmos, estos tiene que estar en la misma “rama”.

Aplicamos al “Define Status”, que contiene las 4 variables como entradas, ambos algoritmos de “clustering”: el algoritmo HAC y el algoritmo K-Means.

En el eje de ordenadas pondremos el atributo “ActivoCorriente/VentasNetas”. En el eje de las ordenadas pondremos el atributo “ActivoCorriente/PasivoCorriente”, y vemos que en la esquina superior derecha podremos elegir con cual de los dos modelos representarlo, si con el “HAC” o con el “K-Means”.

Con el algoritmo “HAC” el diagrama de dispersión sería el siguiente:





Con el algoritmo “K-Means” el diagrama de dispersión sería el siguiente:

