

Instituto Politécnico Nacional
Escuela Superior de Cómputo
Unidad de Aprendizaje: Minería de datos
Ciclo escolar: 2022-2

Proyecto No. 4. Clúster

Grupo: 5CDM1

Equipo: 1

Nombre de los integrantes del equipo:

- 1) De Luna Ocampo Yanina
- 2) Medina Barreras Daniel Iván



INSTITUTO POLITÉCNICO NACIONAL

ESCUELA SUPERIOR DE CÓMPUTO

LICENCIATURA EN CIENCIA DE DATOS

UNIDAD DE APRENDIZAJE

MINERÍA DE DATOS

CREACIÓN DE CLUSTERS - MATRIZ

NOMBRE DE LOS ALUMNOS:

DE LUNA OCAMPO YANINA

MEDINA BARRERAS DANIEL IVÁN

PROFESOR:

OCAMPO BOTELLO FABIOLA

GRUPO:

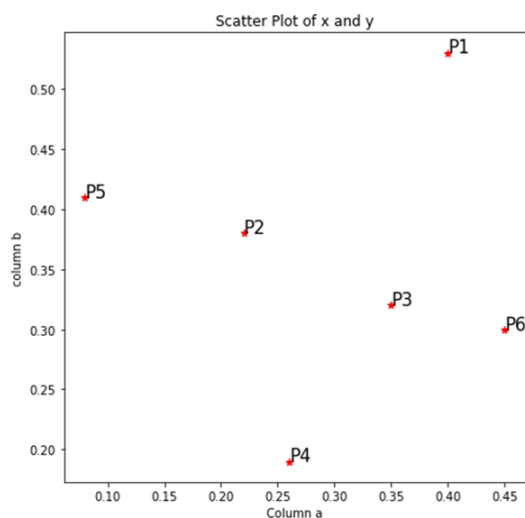
5CDM1

FECHA:

29/05/2022

PARTE 1. Diseñar una matriz con 6 puntuaciones (x, y) y desarrollar el proceso paso a paso de la creación de grupos aplicando las tres siguientes técnicas (*Inserte una portada indicando el tema que aborda*)

| Point | a | b |
|-------|------|------|
| P1 | 0.40 | 0.53 |
| P2 | 0.22 | 0.38 |
| P3 | 0.35 | 0.32 |
| P4 | 0.26 | 0.19 |
| P5 | 0.08 | 0.41 |
| P6 | 0.45 | 0.30 |





INSTITUTO POLITÉCNICO NACIONAL

ESCUELA SUPERIOR DE CÓMPUTO

LICENCIATURA EN CIENCIA DE DATOS

UNIDAD DE APRENDIZAJE

MINERÍA DE DATOS

CREACIÓN DE CLUSTERS - ENLACE SIMPLE

NOMBRE DE LOS ALUMNOS:

DE LUNA OCAMPO YANINA

MEDINA BARRERAS DANIEL IVÁN

PROFESOR:

OCAMPO BOTELLO FABIOLA

GRUPO:

5CDM1

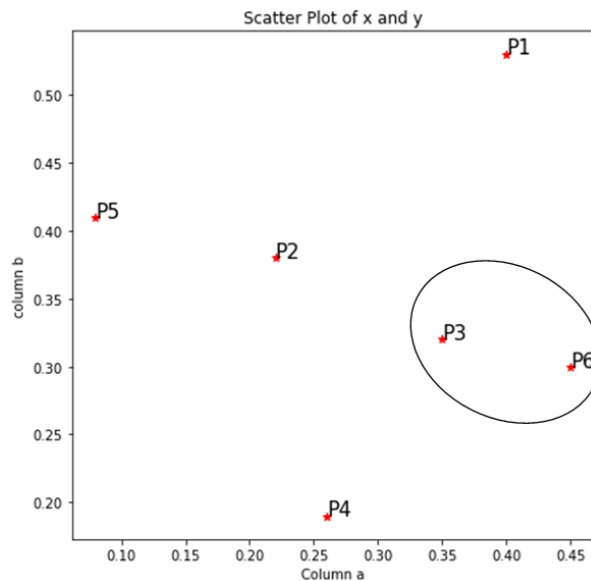
FECHA:

29/05/2022

1. Enlace simple
Matriz de distancias
Primer valor mínimo

| | P1 | P2 | P3 | P4 | P5 | P6 |
|----|------|------|------|------|------|----|
| P1 | 0 | | | | | |
| P2 | 0.23 | 0 | | | | |
| P3 | 0.22 | 0.15 | 0 | | | |
| P4 | 0.37 | 0.20 | 0.15 | 0 | | |
| P5 | 0.34 | 0.14 | 0.28 | 0.29 | 0 | |
| P6 | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0 |

Gráfico primer grupo



Actualizar matriz de distancias después crear primer grupo

$$\text{MIN}(\text{DIST}((P3,P1),(P6,P1))) = \text{MIN}[(0.22,0.23)] = 0.22$$

$$\text{MIN}(\text{DIST}((P3,P2),(P6,P2))) = \text{MIN}[(0.15,0.25)] = 0.15$$

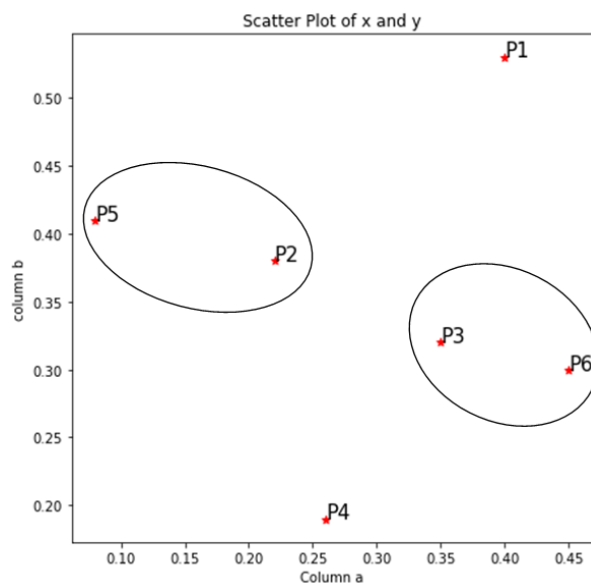
$$\text{MIN}(\text{DIST}((P3,P4),(P6,P4))) = \text{MIN}[(0.15,0.22)] = 0.15$$

$$\text{MIN}(\text{DIST}((P3,P5),(P6,P5))) = \text{MIN}[(0.28,0.39)] = 0.22$$

Matriz actualizada para P3,P6 y nuevo valor mínimo

| | P1 | P2 | P3,P6 | P4 | P5 |
|-------|------|------|-------|------|----|
| P1 | 0 | | | | |
| P2 | 0.24 | 0 | | | |
| P3,P6 | 0.22 | 0.15 | 0 | | |
| P4 | 0.37 | 0.20 | 0.15 | 0 | |
| P5 | 0.34 | 0.14 | 0.28 | 0.29 | 0 |

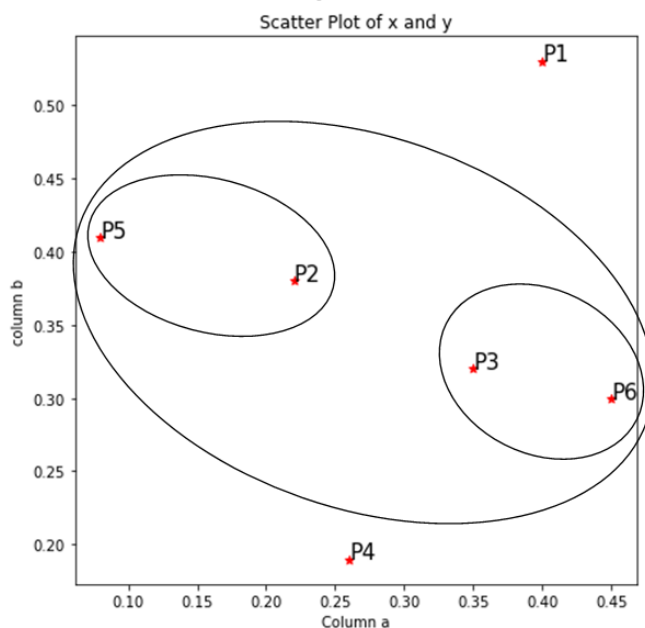
Gráfico segundo grupo



$\text{MIN}(\text{DIST}((P2,P1),(P5,P1))) = \text{MIN}[(0.23,0.34)] = 0.23$
 $\text{MIN}(\text{DIST}((P2,(P3,P6)), (P5,(P3,P6)))) = \text{MIN}[(0.15,0.28)] = 0.15$
 $\text{MIN}(\text{DIST}((P2,P4),(P5,P4))) = \text{MIN}[(0.20,0.29)] = 0.20$
 Matriz actualizada P5,P2

| | P1 | P2, P5 | P3,P6 | P4 |
|--------|------|--------|-------|----|
| P1 | 0 | | | |
| P2, P5 | 0.23 | 0 | | |
| P3,P6 | 0.22 | 0.15 | 0 | |
| P4 | 0.37 | 0.20 | 0.15 | 0 |

Gráfico para el tercer grupo



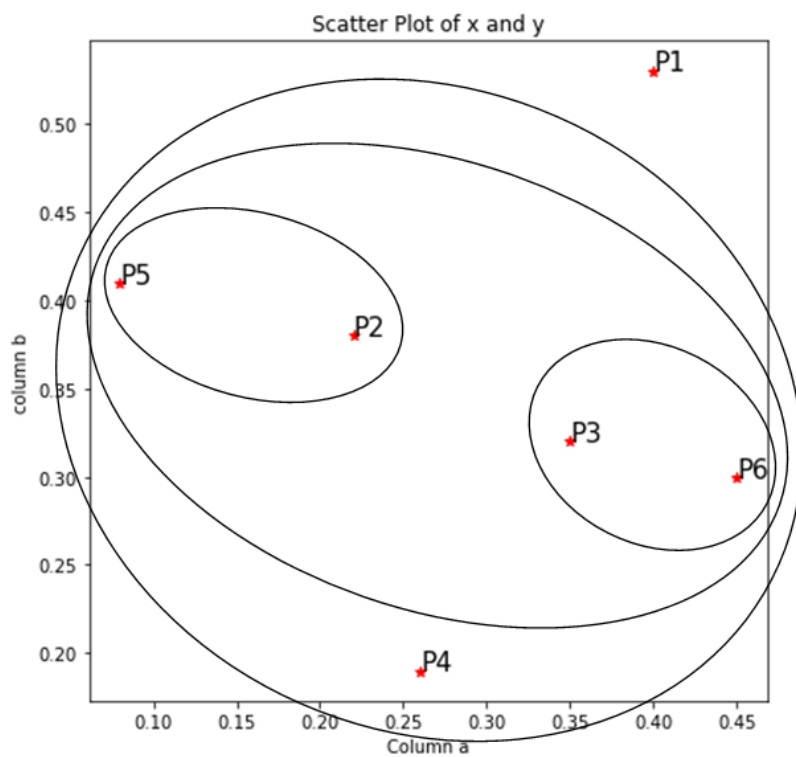
$$\text{MIN}(\text{DIST}(((P2, P5), P1), ((P3, P6), P1))) = \text{MIN}[(0.23, 0.22)] = 0.22$$

$$\text{MIN}(\text{DIST}(((P2, P5), P4), ((P3, P6), P4))) = \text{MIN}[(0.20, 0.15)] = 0.15$$

Matriz actualizada P2, P5, P3, P6

| | P1 | P2, P5, P3, P6 | P4 |
|----------------|------|----------------|----|
| P1 | 0 | | |
| P2, P5, P3, P6 | 0.22 | 0 | |
| P4 | 0.37 | 0.15 | 0 |

Gráfica 4to grupo

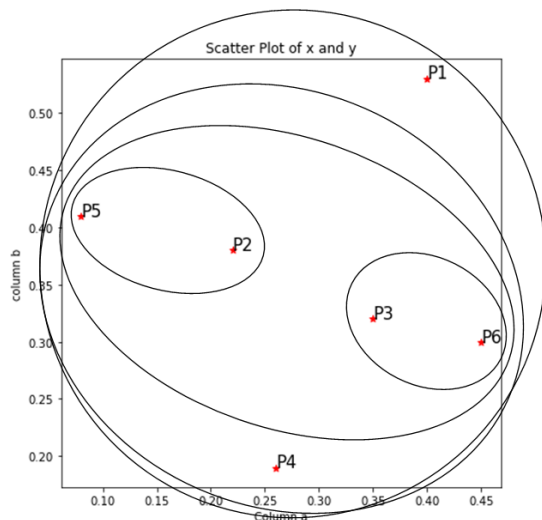


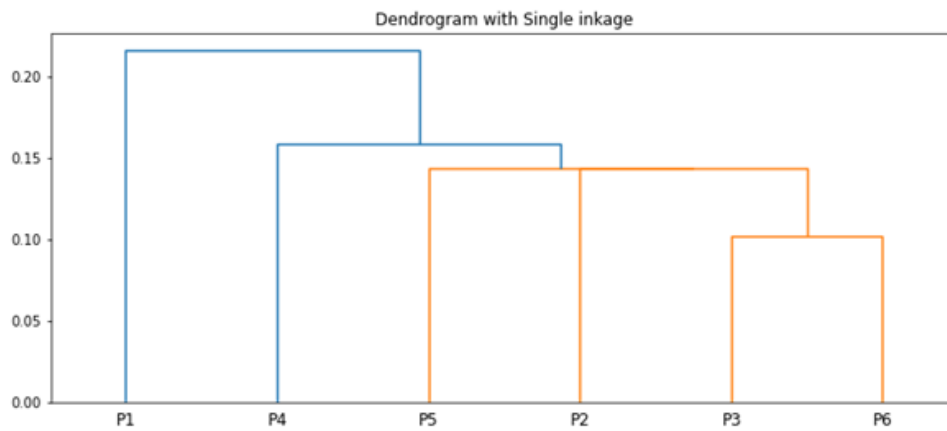
$$\text{MIN}(\text{DIST}(((p2,p5,p3,p6)),P1),((P4,P1))) = \text{MIN}[(0.22,0.37)] = 0.22$$

Matriz actualizada para P2,P5,P3,P6,P4

| | P1 | P2, P5,P3,P6,P4 |
|-------------------|------|-----------------|
| P1 | 0 | |
| P2, P5, P3, P6,P4 | 0.22 | 0 |

Gráfica final







INSTITUTO POLITÉCNICO NACIONAL

ESCUELA SUPERIOR DE CÓMPUTO

LICENCIATURA EN CIENCIA DE DATOS

UNIDAD DE APRENDIZAJE

MINERÍA DE DATOS

CREACIÓN DE CLUSTERS - ENLACE PROMEDIO

NOMBRE DE LOS ALUMNOS:

DE LUNA OCAMPO YANINA

MEDINA BARRERAS DANIEL IVÁN

PROFESOR:

OCAMPO BOTELLO FABIOLA

GRUPO:

5CDM1

FECHA:

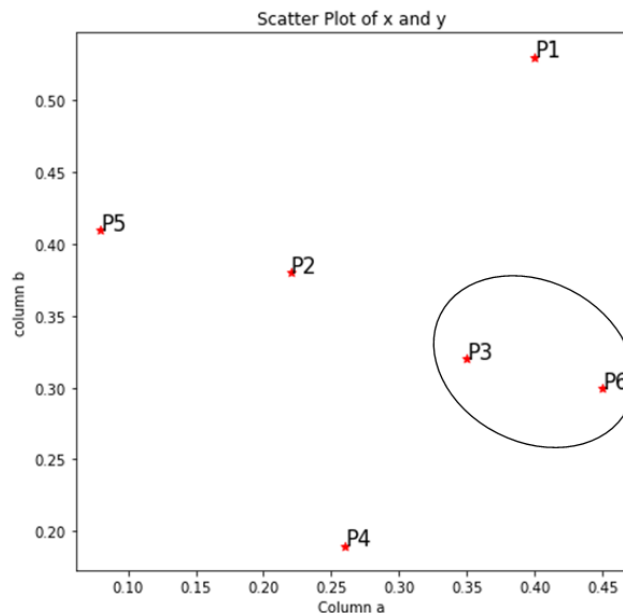
29/05/2022

2. Enlace promedio

Matriz de distancias

Primer valor mínimo

| | P1 | P2 | P3 | P4 | P5 | P6 |
|----|------|------|------|------|------|----|
| P1 | 0 | | | | | |
| P2 | 0.23 | 0 | | | | |
| P3 | 0.22 | 0.15 | 0 | | | |
| P4 | 0.37 | 0.20 | 0.15 | 0 | | |
| P5 | 0.34 | 0.14 | 0.28 | 0.29 | 0 | |
| P6 | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0 |



Actualizar matriz de distancias después crear primer grupo

$$\text{AVG}(\text{DIST}((P3, P6), P1)) =$$

$$\text{DIST}[(P3, P1), P1] = \frac{1}{2}(\text{DIST}(P3, P1) + \text{DIST}(P6, P1)) = \frac{1}{2}(0.22 + 0.23) = \frac{1}{2}(0.45) = 0.23$$

$$\text{AVG}(\text{DIST}((P3, P6), P2)) =$$

$$\text{DIST}[(P3, P6), P2] = \frac{1}{2}(\text{DIST}(P3, P2) + \text{DIST}(P6, P2)) = \frac{1}{2}(0.15 + 0.25) = \frac{1}{2}(0.4) = 0.2$$

$$\text{AVG}(\text{DIST}((P3, P6), P4)) =$$

$$\text{DIST}[(P3, P6), P4] = \frac{1}{2}(\text{DIST}(P3, P4) + \text{DIST}(P6, P2)) = \frac{1}{2}(0.15 + 0.22) = \frac{1}{2}(0.37) = 0.19$$

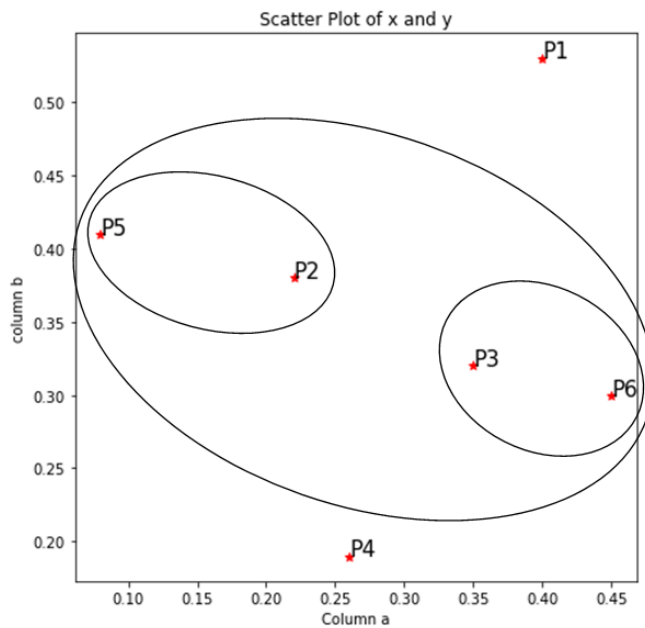
$$\text{AVG}(\text{DIST}((P3, P6), P5)) =$$

$$\text{DIST}[(P3,P6),P5] = \frac{1}{2}(\text{DIST}(P3,P5)+\text{DIST}(P6,P5))=\frac{1}{2}(0.28+0.39)=\frac{1}{2}(0.67)=0.34$$

Matriz actualizada para P3, P6, más el valor mínimo nuevo

| | P1 | P2 | P3,P6 | P4 | P5 |
|-------|------|------|-------|------|----|
| P1 | 0 | | | | |
| P2 | 0.24 | 0 | | | |
| P3,P6 | 0.23 | 0.2 | 0 | | |
| P4 | 0.37 | 0.20 | 0.19 | 0 | |
| P5 | 0.34 | 0.14 | 0.34 | 0.29 | 0 |

Gráfico para el tercer grupo



$$\begin{aligned} \text{AVG}(\text{DIST}((P2,P5),P1)) &= \\ \text{DIST}[(P2,P5),P1] &= \frac{1}{2}(\text{DIST}(P2,P1)+\text{DIST}(P5,P1))=\frac{1}{2}(0.23+0.34)=\frac{1}{2}(0.57)=0.29 \end{aligned}$$

$$\begin{aligned} \text{AVG}(\text{DIST}((P2,P5),(P3,P6))) &= \\ \text{DIST}[(P2,P5),(P3,P6)] &= \\ \frac{1}{2}(\text{DIST}(P2,(P3,P6))+\text{DIST}(P5,(P3,P6))) &= \frac{1}{2}(0.2+0.34)=\frac{1}{2}(0.54)=0.27 \end{aligned}$$

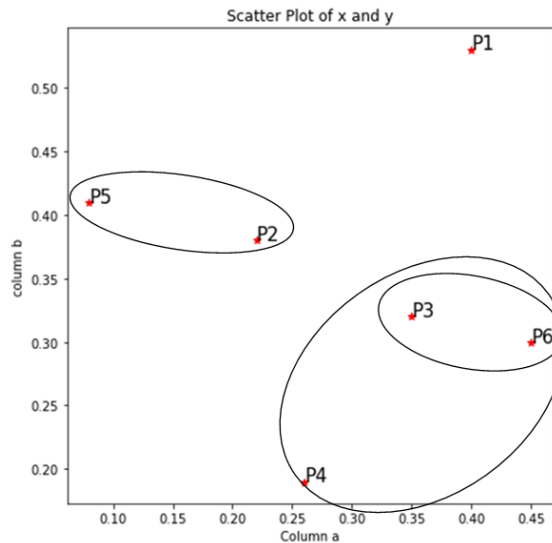
$$\text{AVG}(\text{DIST}((P2,P5),P4)) =$$

$$\text{DIST}[(P2,P5),P4] = \frac{1}{2}(\text{DIST}(P2,(P4))+\text{DIST}(P5,P4))=\frac{1}{2}(0.20+0.29)=\frac{1}{2}(0.49)=0.25$$

Matriz actualizada P2,P5

| | P1 | P2, P5 | P3,P6 | P4 |
|--------|------|--------|-------|----|
| P1 | 0 | | | |
| P2, P5 | 0.29 | 0 | | |
| P3,P6 | 0.22 | 0.27 | 0 | |
| P4 | 0.37 | 0.25 | 0.19 | 0 |

Gráfica actualizada tercer grupo



$$\text{AVG}(\text{DIST}((P3,P6,P4),P1)) =$$

$$\text{DIST}[(P3,P6,P4),P1] = \frac{1}{2}(\text{DIST}((P3,P6),P1)+\text{DIST}((P4),P1))=\frac{1}{2}(0.23+0.37)=\frac{1}{2}(0.6)=0.3$$

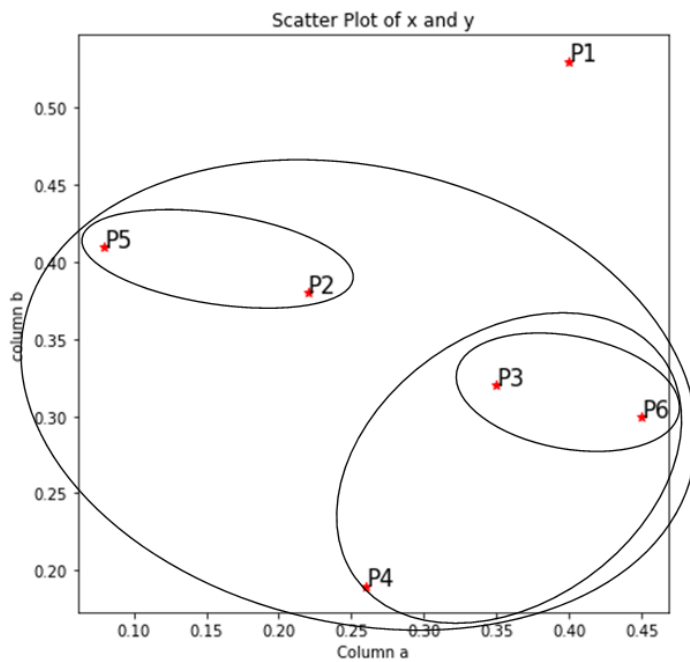
$$\text{AVG}(\text{DIST}((P3,P6,P4),(P2,P5))) =$$

$$\text{DIST}[(P3,P6,P4),(P2,P5)] =$$

$$\frac{1}{2}(\text{DIST}((P3,P6),(P2,P5))+\text{DIST}((P4),(P2,P5)))=\frac{1}{2}(0.27+0.25)=\frac{1}{2}(0.52)=0.26$$

| | | | |
|----------|------|---------|----------|
| | P1 | P2, P5, | P3,P6,P4 |
| P1 | 0 | | |
| P2, P5, | 0.29 | 0 | |
| P3,P6,P4 | 0.3 | 0.26 | 0 |

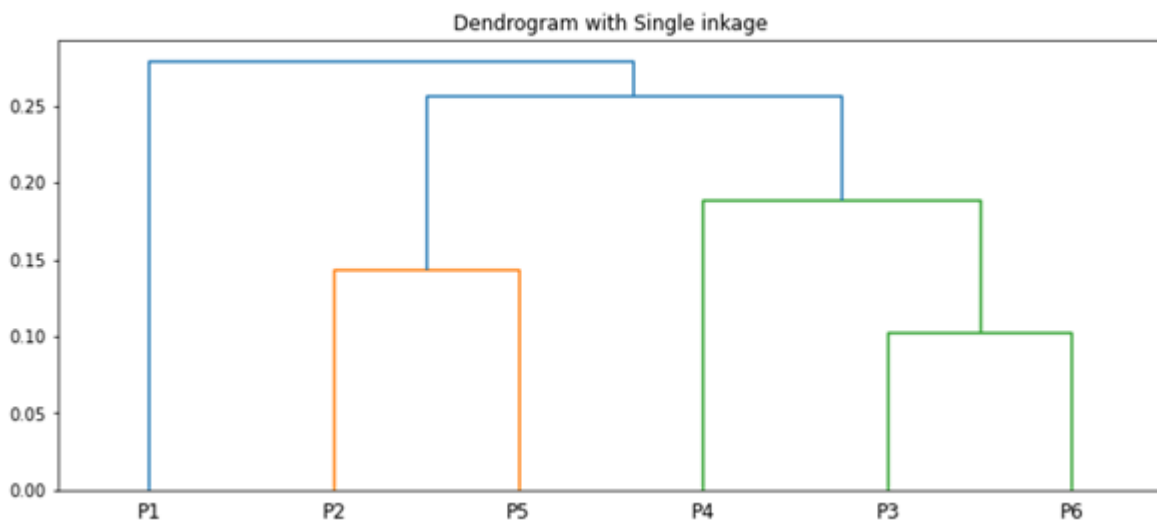
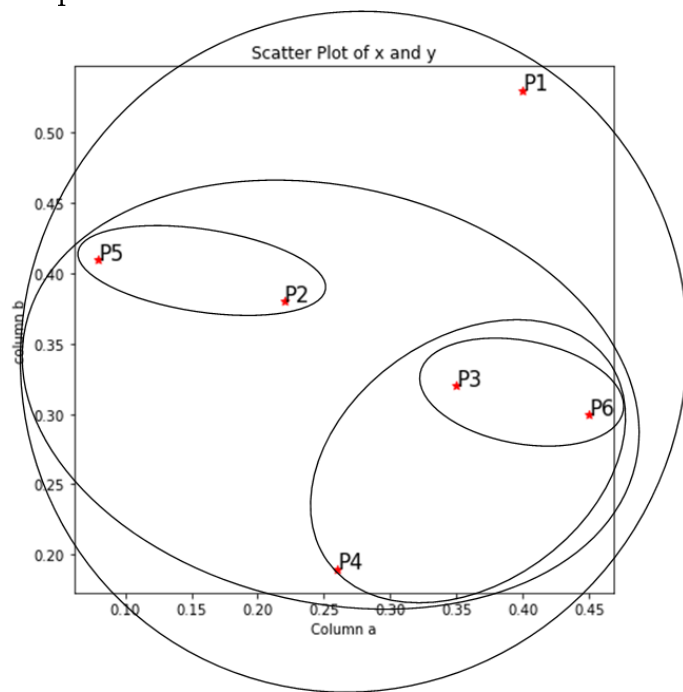
Cuarto grupo



$$\begin{aligned} \text{AVG}(\text{DIST}((P3, P6, P4), (P2, P5))) &= \\ \text{DIST}[(P3, P6, P4), (P2, P5)] &= \\ \frac{1}{2}(\text{DIST}((P3, P6, P4), P1) + \text{DIST}((P2, P5), P1)) &= \frac{1}{2}(0.3 + 0.29) = \frac{1}{2}(0.59) = 0.3 \end{aligned}$$

| | | |
|--------------------|-----|--------------------|
| | P1 | P2, P5, P3, P6, P4 |
| P1 | 0 | |
| P2, P5, P3, P6, P4 | 0.3 | 0 |

Grupo final





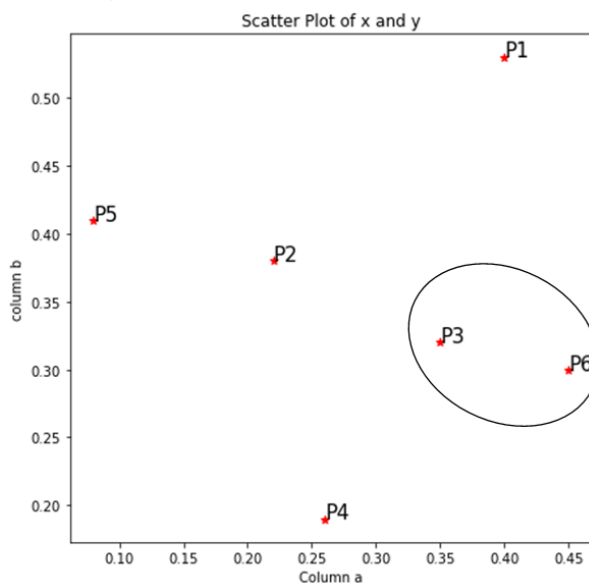
3. Enlace completo

Matriz de distancias

Primer valor mínimo

| | P1 | P2 | P3 | P4 | P5 | P6 |
|----|------|------|------|------|------|----|
| P1 | 0 | | | | | |
| P2 | 0.23 | 0 | | | | |
| P3 | 0.22 | 0.15 | 0 | | | |
| P4 | 0.37 | 0.20 | 0.15 | 0 | | |
| P5 | 0.34 | 0.14 | 0.28 | 0.29 | 0 | |
| P6 | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0 |

Primer grupo



Actualizar matriz de distancias después crear primer grupo

$$\text{MAX}(\text{DIST}((P3, P6), P1))) = \text{MAX}(\text{DIST}((P3, P1), (P6, P1))) = \text{MAX}[(0.22, 0.23)] = 0.23$$

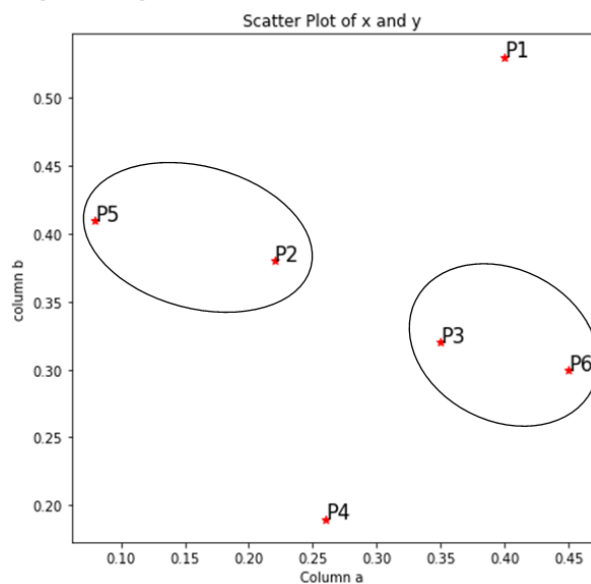
$$\text{MAX}(\text{DIST}((P3, P6), P2))) = \text{MAX}(\text{DIST}((P3, P2), (P6, P2))) = \text{MAX}[(0.15, 0.25)] = 0.25$$

$$\text{MAX}(\text{DIST}((P3, P6), P4))) = \text{MAX}(\text{DIST}((P3, P4), (P6, P4))) = \text{MAX}[(0.15, 0.22)] = 0.22$$

$$\text{MAX}(\text{DIST}((P3, P6), P5))) = \text{MAX}(\text{DIST}((P3, P2), (P6, P5))) = \text{MAX}[(0.28, 0.39)] = 0.39$$

| | P1 | P2 | P3,P6 | P4 | P5 |
|-------|------|------|-------|------|----|
| P1 | 0 | | | | |
| P2 | 0.23 | 0 | | | |
| P3,P6 | 0.23 | 0.25 | 0 | | |
| P4 | 0.37 | 0.20 | 0.2 | 0 | |
| P5 | 0.34 | 0.14 | 0.39 | 0.29 | 0 |

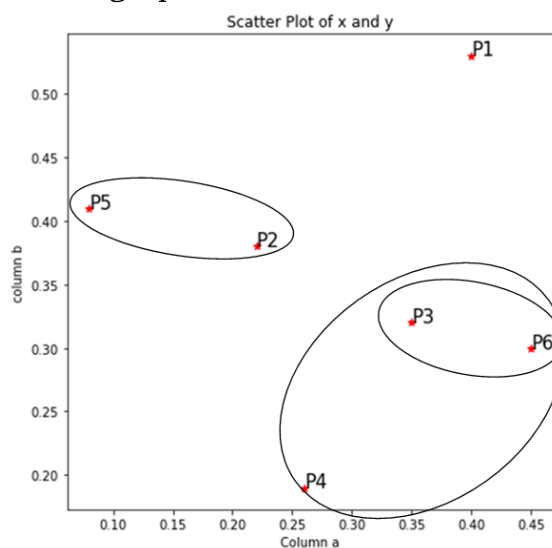
Segundo grupo



$\text{MAX}(\text{DIST}((P2,P5),P1))) = \text{MAX}(\text{DIST}((P2,P1),(P5,P1)))) = \text{MAX}[(0.23,0.34)] = 0.34$
 $\text{MAX}(\text{DIST}((P2,P5),(P3,P6)))) = \text{MAX}(\text{DIST}((P2,(P3,P6),(P5,(P3,P6)))) =$
 $\text{MAX}[(0.25,0.39)] = 0.39$
 $\text{MAX}(\text{DIST}((P2,P5),P4))) = \text{MAX}(\text{DIST}((P2,P4),(P5,P4)))) = \text{MAX}[(0.20,0.29)] = 0.29$

| | P1 | P2, P5 | P3,P6 | P4 |
|--------|------|--------|-------|----|
| P1 | 0 | | | |
| P2, P5 | 0.34 | 0 | | |
| P3,P6 | 0.23 | 0.39 | 0 | |
| P4 | 0.37 | 0.29 | 0.22 | 0 |

Cuarto grupo

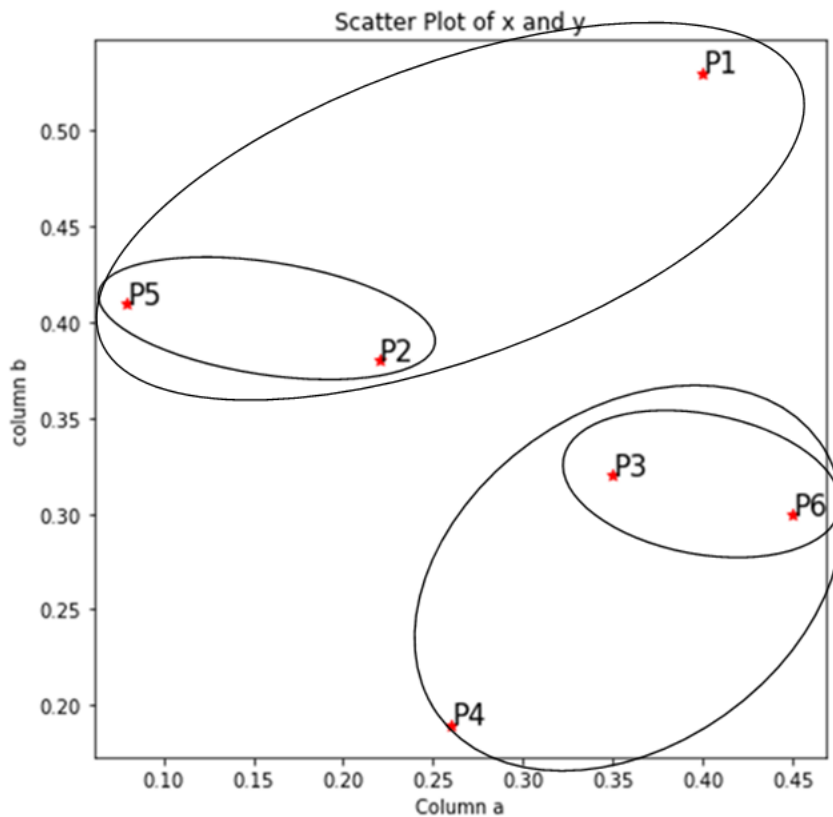


$$\text{MAX}(\text{DIST}(((P3,P6),P4),P1))) = \text{MAX}(\text{DIST}(((P3,P6),P1),(P4,P1)))) = \text{MAX}[(0.23,0.37)] = 0.37$$

$$\text{MAX}(\text{DIST}(((P3,P6),P4),(P4,P5)))) = \text{MAX}(\text{DIST}(((P3,P6),(P4,P5)),(P4,(P4,P5)))) = \text{MAX}[(0.39,0.29)] = 0.39$$

| | P1 | P2, P5, | P3,P6,P4 |
|----------|------|---------|----------|
| P1 | 0 | | |
| P2, P5 | 0.34 | 0 | |
| P3,P6,P4 | 0.37 | 0.39 | 0 |

Grupo 4



$$\text{MAX}(\text{DIST}(((P2,P5),(P3,P6,P4),(P1,(P3,P6,P4)))))) =$$

$$\text{MAX}(\text{DIST}(((P2,P5),(P3,P6,P4)),(P1,(P3,P5,P4)))) = \text{MAX}[(0.39,0.37)] = 0.39$$

| | P2,P5,P1 | P3,P6,P4 |
|----------|----------|----------|
| P2,P5,P1 | 0 | |
| P3,P6,P4 | 0.39 | 0 |

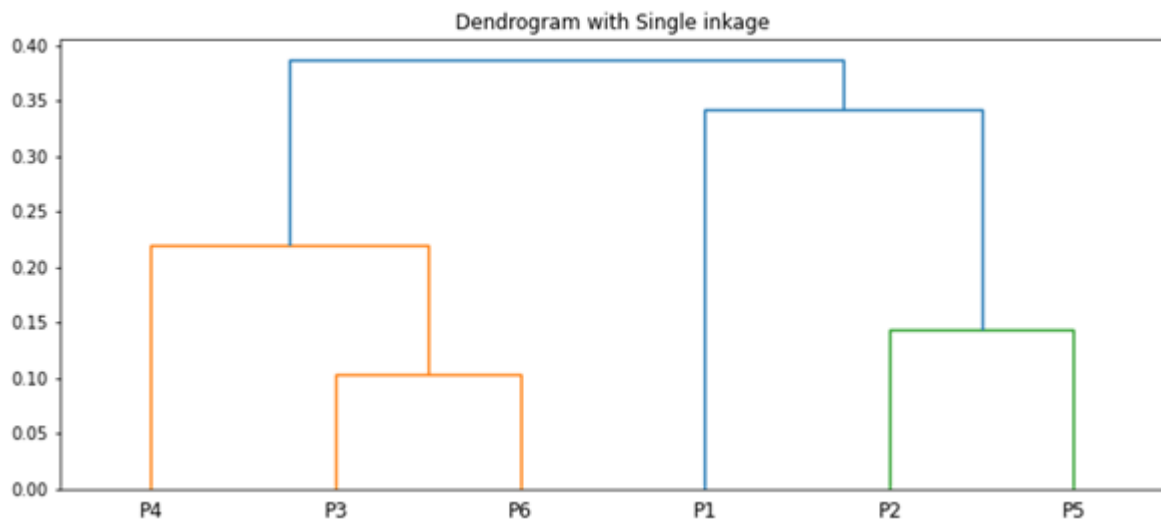
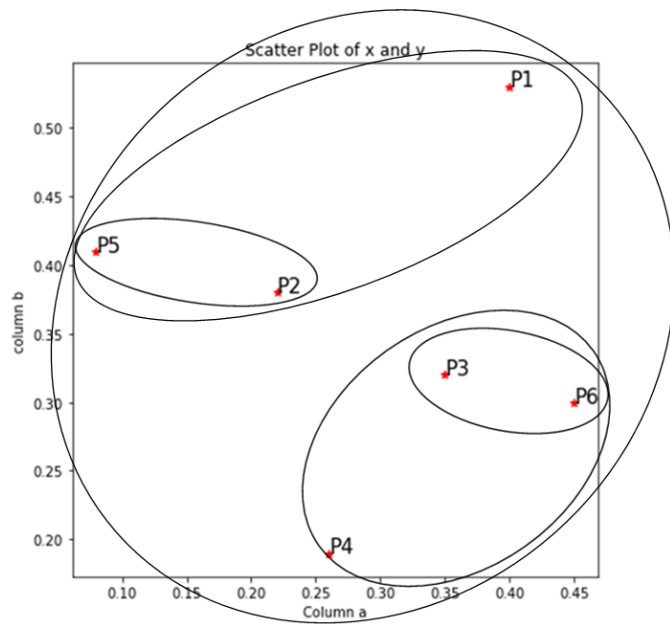
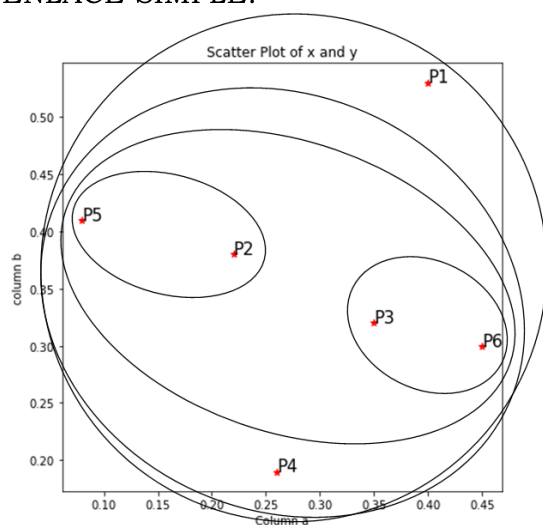
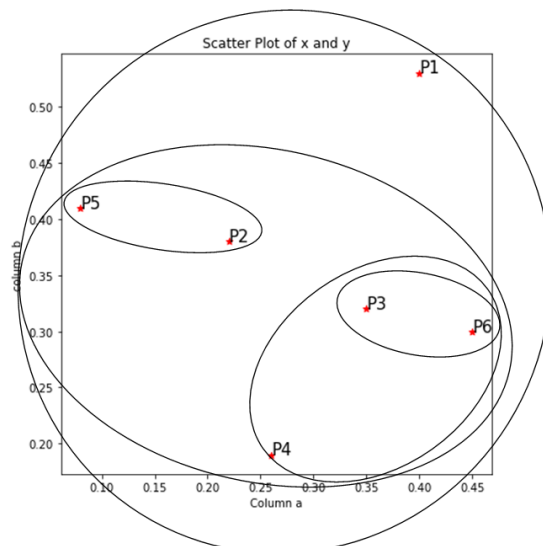


TABLA COMPARATIVA DE CLÚSTERS OBTENIDOS:
ENLACE SIMPLE:



ENLACE PROMEDIO:



ENLACE COMPLETO:

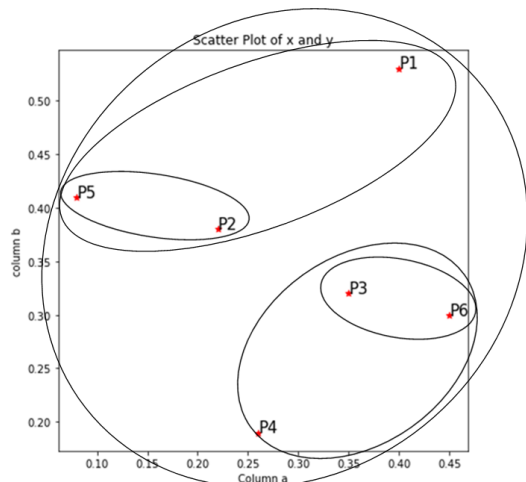
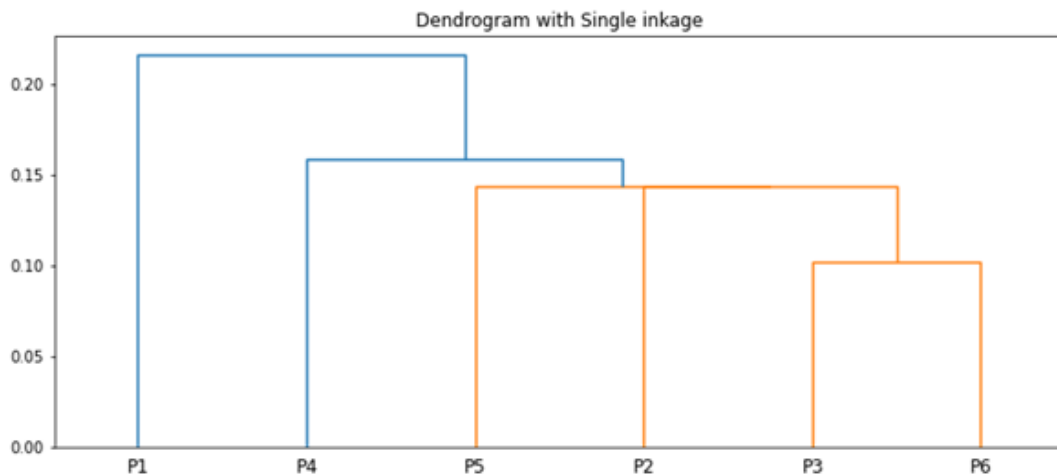
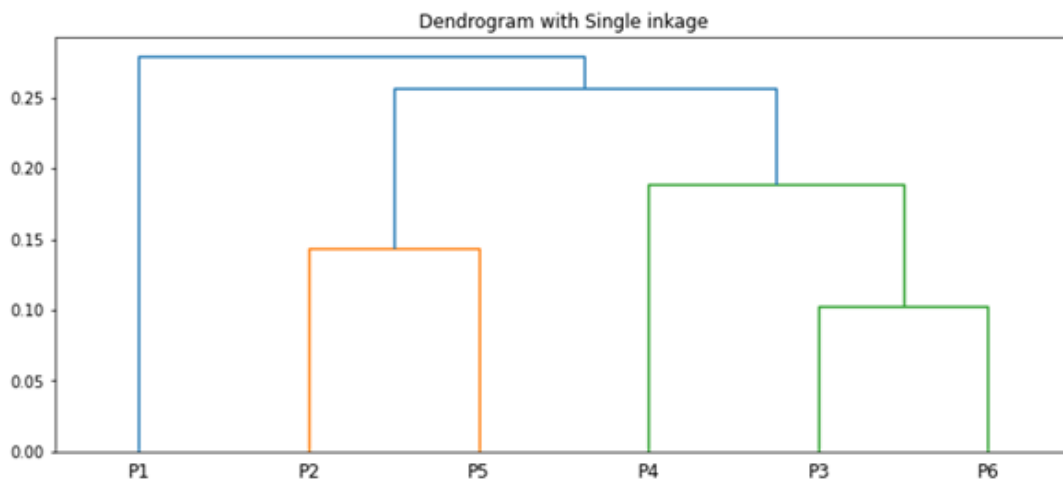


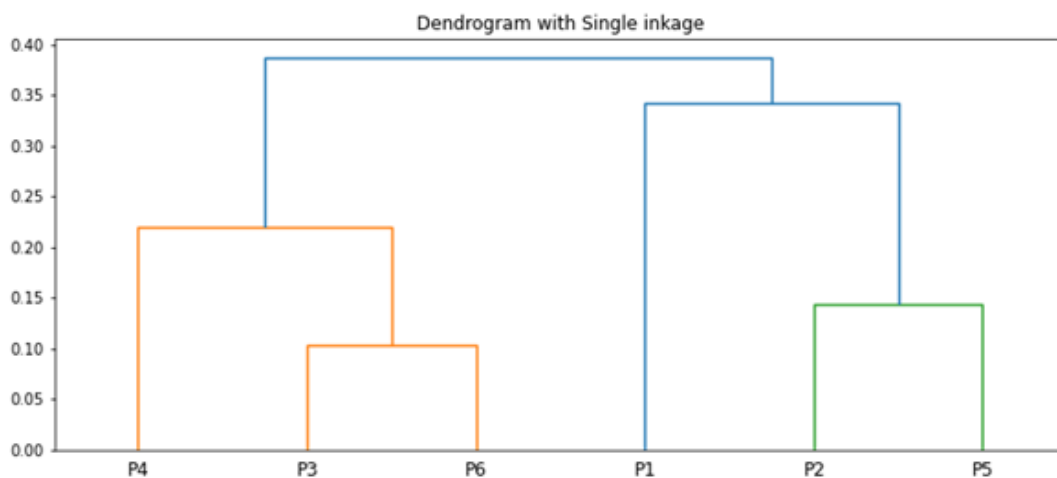
TABLA COMPARATIVA DE DENDOGRAMAS OBTENIDOS:
ENLACE SIMPLE:



ENLACE PROMEDIO:



ENLACE COMPLETO:





INSTITUTO POLITÉCNICO NACIONAL

ESCUELA SUPERIOR DE CÓMPUTO

LICENCIATURA EN CIENCIA DE DATOS

UNIDAD DE APRENDIZAJE

MINERÍA DE DATOS

CREACIÓN DE CLUSTERS - PARTE 2

NOMBRE DE LOS ALUMNOS:

DE LUNA OCAMPO YANINA

MEDINA BARRERAS DANIEL IVÁN

PROFESOR:

OCAMPO BOTELLO FABIOLA

GRUPO:

5CDM1

FECHA:

29/05/202

PARTE 2. Diseñar una matriz con datos cualitativos, cinco registros y seis características y genere los grupos.

| | Tos | Fiebre | camPeso | Hemorragia | Náuseas | dolCabeza |
|------------|-----|--------|---------|------------|---------|-----------|
| Paciente 1 | 0 | 0 | 1 | 0 | 1 | 1 |
| Paciente 2 | 1 | 1 | 0 | 1 | 0 | 0 |
| Paciente 3 | 1 | 0 | 0 | 1 | 1 | 0 |
| Paciente 4 | 0 | 0 | 1 | 1 | 1 | 0 |
| Paciente 5 | 1 | 1 | 1 | 0 | 1 | 1 |

Coeficiente de Jaccard

| | 1 | 2 | 3 | 4 | 5 |
|---|------|------|------|------|---|
| 1 | 1 | | | | |
| 2 | 0 | 1 | | | |
| 3 | 0.60 | 0.25 | 1 | | |
| 4 | 0.33 | 0.50 | 0.60 | 1 | |
| 5 | 0.17 | 0.67 | 0.17 | 0.40 | 1 |

Pasamos a la siguiente matriz, debido al umbral, que es este caso es: 0.51

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 1 | | | | |
| 2 | 0 | 1 | | | |
| 3 | 1 | 0 | 1 | | |
| 4 | 0 | 0 | 1 | 1 | |
| 5 | 0 | 1 | 0 | 0 | 1 |

Multiplicamos por sí misma, con ayuda del software obtenemos:

$$C = A \cdot B = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 2 & 0 & 1 & 0 & 0 \\ 1 & 0 & 2 & 1 & 0 \\ 0 & 2 & 0 & 0 & 1 \end{pmatrix}$$

Vemos cuántos vecinos en común hay. Obteniendo con la matriz de adyacencia lo siguiente:

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | - | | | | |
| 2 | 0 | - | | | |
| 3 | 2 | 0 | - | | |
| 4 | 1 | 0 | 2 | - | |
| 5 | 0 | 2 | 0 | 0 | - |

Podemos ver que hay grupos con el mismo número de enlaces. Escogemos uno, será el 3 con el 4.

| | 1 | 2 | (3, 4) | 5 |
|--------|-----|---|--------|---|
| 1 | - | | | |
| 2 | 0 | - | | |
| (3, 4) | 2+1 | 0 | - | |
| 5 | 0 | 2 | 0 | - |

Formamos un clúster ahora con el (1,3,4) por el número de enlaces contenidos. Nos quedaría:

| | (1, 3, 4) | 2 | 5 |
|-----------|-----------|---|---|
| (1, 3, 4) | - | | |

| | | | |
|----------|---|---|---|
| 2 | 0 | - | |
| 5 | 0 | 2 | - |

Solo nos queda el 2 y el 5 por lo que crearemos un clúster, como se muestra a continuación:

{1, 3, 4} y {2, 5}

Recordando la tabla principal, vemos que estos agrupamientos son correctos debido a que por ejemplo, el 5 tiene más cosas en común con el 2 que con los demás.

Este usa el algoritmo ROCK, que es: RObust Clustering using linKs. Sirve para conjunto de datos con atributos categóricos y booleanos. Un par de puntos se consideran vecinos si su similitud está por encima del umbral preestablecido.

- El algoritmo ROCK evalúa las distancias entre objetos utilizando el coeficiente Jaccard.
- Utiliza el parámetro θ para determinar quienes son los vecinos en cada uno de los objetos.
- Dado un punto p, un punto q es vecino de p si el coeficiente de Jaccard $\text{sim}(p, q)$ excede el valor de θ .
- Se generan los valores de la matriz de ligas (links), la cual consiste en la evaluación de $\text{links}(p, q)$ como el número de vecinos comunes entre los puntos p y q.

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$