

CSCI 347: Introduction to Data Mining

Hierarchical Clustering Example

AGGLOMERATIVE CLUSTERING ALGORITHM

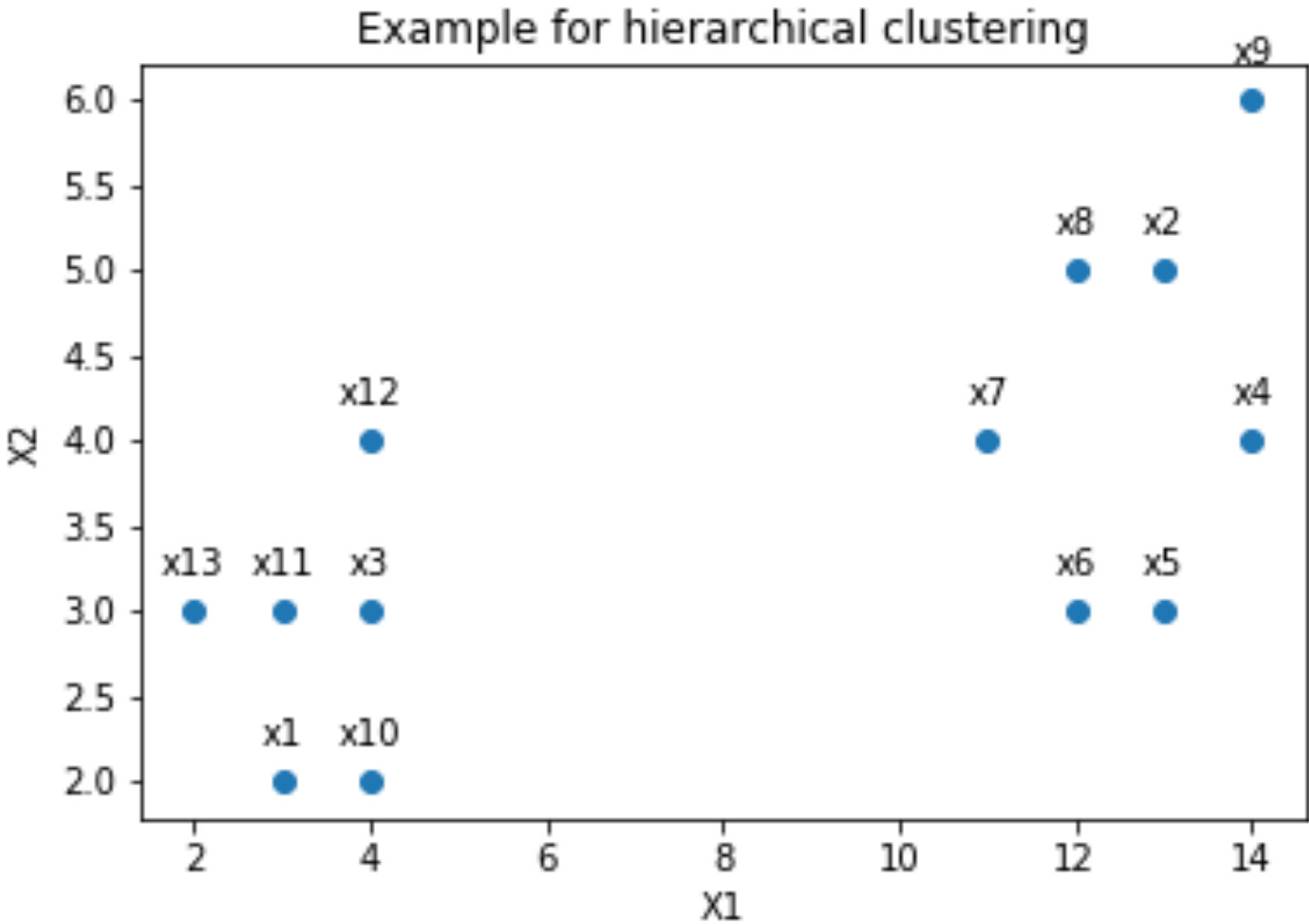
AgglomerativeClustering(D, k) :

1. $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$
2. $\Delta = \{\delta(x_i, x_j) : x_i, x_j \in D\}$
3. repeat:
 4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$
 5. $C_{i,j} = C_i \cup C_j$
 6. $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$
 7. Update distance matrix Δ to reflect new clustering
4. Until $|\mathcal{C}| = k$

EXAMPLE

$D =$

	X_1	X_2
x_1	3	2
x_2	13	5
x_3	4	3
x_4	14	4
x_5	13	3
x_6	12	3
x_7	11	4
x_8	12	5
x_9	14	6
x_{10}	4	2
x_{11}	3	3
x_{12}	4	4
x_{13}	2	3



AGGLOMERATIVE CLUSTERING ALGORITHM

AgglomerativeClustering(D, k) :

1. $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$
2. $\Delta = \{\delta(x_i, x_j) : x_i, x_j \in D\}$
3. repeat:
 4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$
 5. $C_{i,j} = C_i \cup C_j$
 6. $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$
 7. Update distance matrix Δ to reflect new clustering
4. Until $|\mathcal{C}| = k$

$$D =$$

	X_1	X_2
x_1	3	2
x_2	13	5
x_3	4	3
x_4	14	4
x_5	13	3
x_6	12	3
x_7	11	4
x_8	12	5
x_9	14	6
x_{10}	4	2
x_{11}	3	3
x_{12}	4	4
x_{13}	2	3

 =

[illegible]

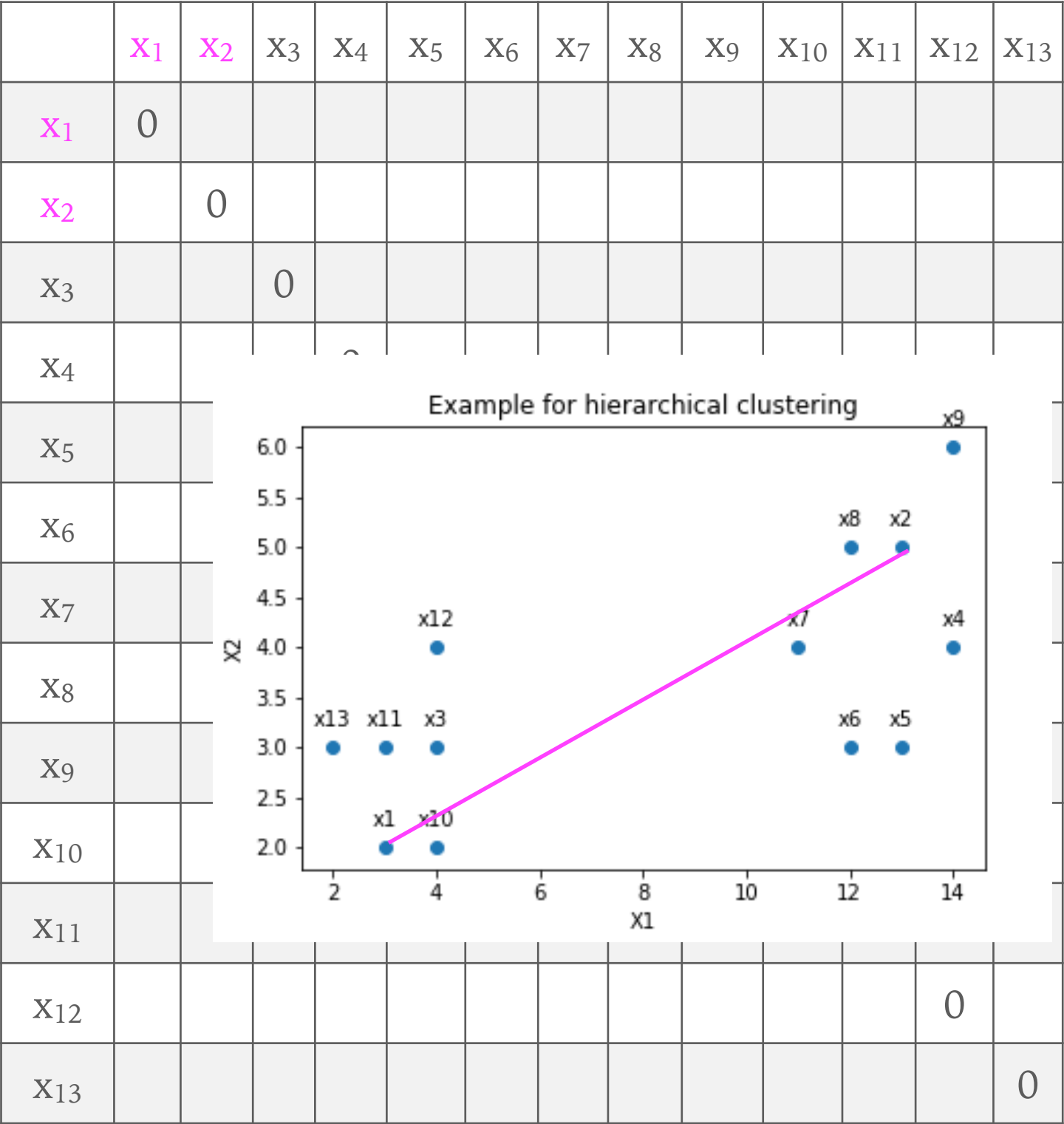
EXAMPLE

$$\delta(x_1, x_2) = \sqrt{(3 - 13)^2 + (2 - 5)^2} = \sqrt{100 + 9} = \sqrt{109} = 10.44$$

D =

	X ₁	X ₂
x ₁	3	2
x ₂	13	5
x ₃	4	3
x ₄	14	4
x ₅	13	3
x ₆	12	3
x ₇	11	4
x ₈	12	5
x ₉	14	6
x ₁₀	4	2
x ₁₁	3	3
x ₁₂	4	4
x ₁₃	2	3

Δ =



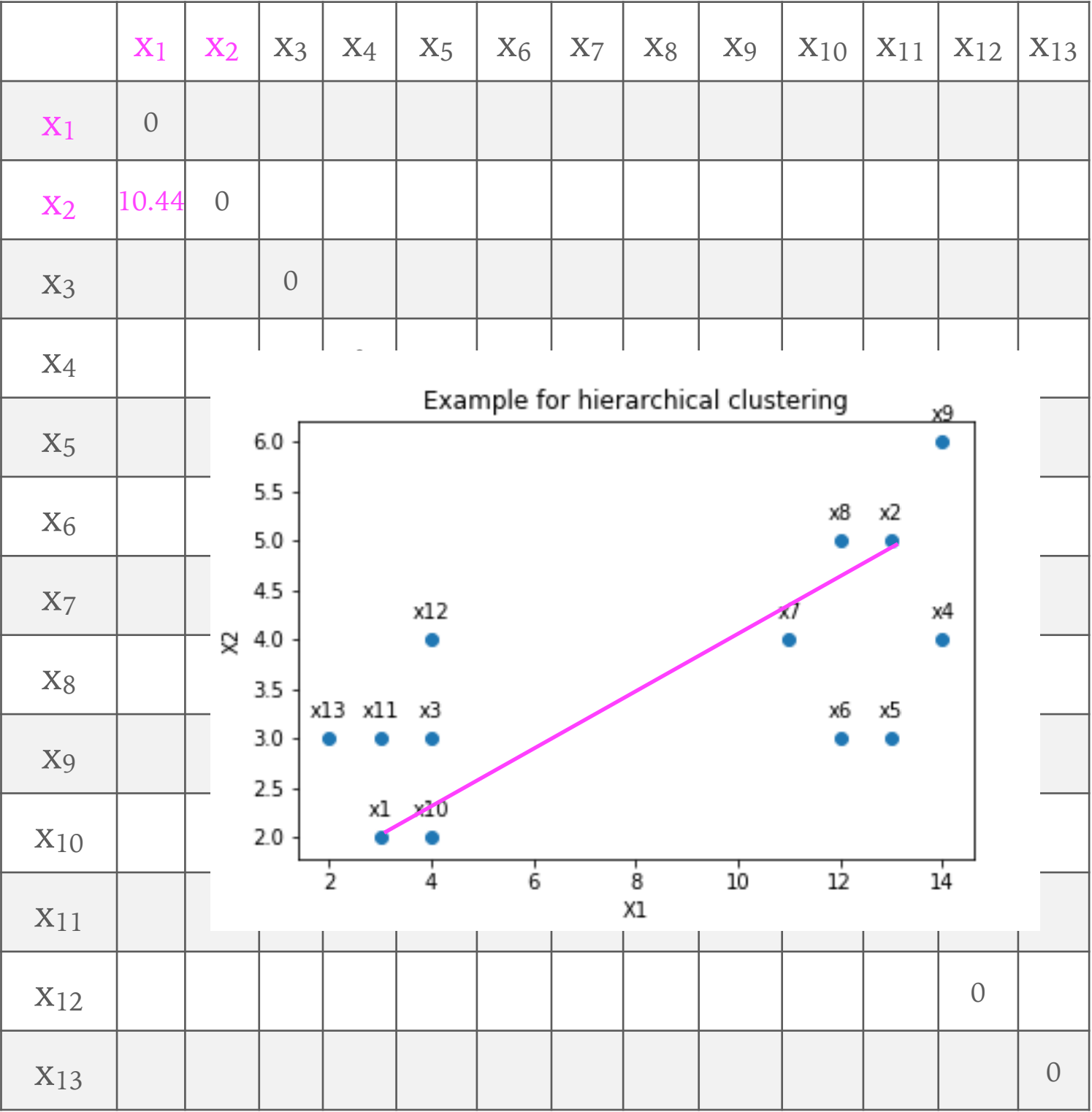
EXAMPLE

$$\delta(x_1, x_2) = \sqrt{(3 - 13)^2 + (2 - 5)^2} = \sqrt{100 + 9} = \sqrt{109} = 10.44$$

D =

	X ₁	X ₂
x ₁	3	2
x ₂	13	5
x ₃	4	3
x ₄	14	4
x ₅	13	3
x ₆	12	3
x ₇	11	4
x ₈	12	5
x ₉	14	6
x ₁₀	4	2
x ₁₁	3	3
x ₁₂	4	4
x ₁₃	2	3

Δ =



EXAMPLE

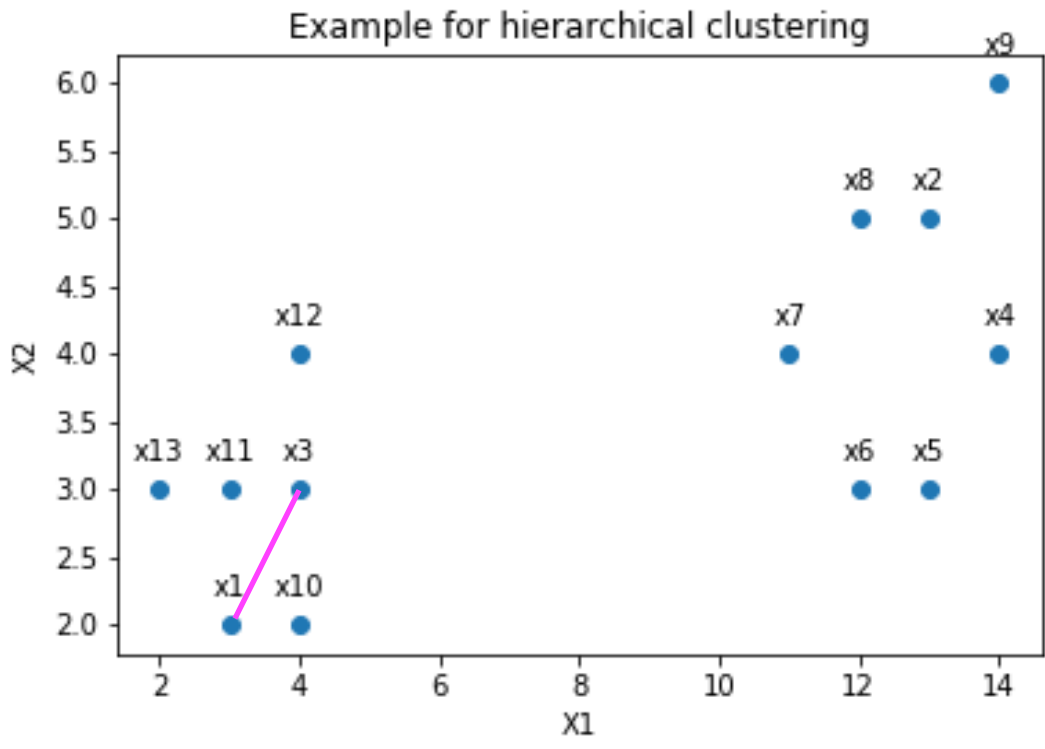
$\delta(x_1, x_3) = \sqrt{(3 - 4)^2 + (2 - 3)^2} = \sqrt{1 + 1} = \sqrt{2} = 1.41$

D =

	X ₁	X ₂
x ₁	3	2
x ₂	13	5
x ₃	4	3
x ₄	14	4
x ₅	13	3
x ₆	12	3
x ₇	11	4
x ₈	12	5
x ₉	14	6
x ₁₀	4	2
x ₁₁	3	3
x ₁₂	4	4
x ₁₃	2	3

Δ =

	x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈	x ₉	x ₁₀	x ₁₁	x ₁₂	x ₁₃
x ₁	0												
x ₂	10.44	0											
x ₃	1.41		0										
x ₄													
x ₅													
x ₆													
x ₇													
x ₈													
x ₉													
x ₁₀													
x ₁₁													
x ₁₂												0	
x ₁₃													0



EXAMPLE

$D =$

	X_1	X_2
x_1	3	2
x_2	13	5
x_3	4	3
x_4	14	4
x_5	13	3
x_6	12	3
x_7	11	4
x_8	12	5
x_9	14	6
x_{10}	4	2
x_{11}	3	3
x_{12}	4	4
x_{13}	2	3

$\Delta =$

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}	X_{11}	X_{12}	X_{13}
x_1	0												
x_2	10.44	0											
x_3	1.41	9.21	0										
x_4	11.18	1.41	10.05	0									
x_5	10.05	2	9	1.41	0								
x_6	9.06	2.24	8	2.24	1	0							
x_7	8.25	2.24	7.07	3	2.24	1.41	0						
x_8	9.49	1	8.25	2.23	2.24	2	1.41	0					
x_9	11.70	1.41	10.44	2	3.16	3.61	3.61	2.24	0				
x_{10}	1	9.49	1	10.20	9.06	8.06	7.28	8.54	10.77	0			
x_{11}	1	10.20	1	11.05	10	9	8.06	9.22	11.40	1.41	0		
x_{12}	2.24	9.06	1	10	9.06	8.06	7	8.06	10.20	2	1.41	0	
x_{13}	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	2.24	1	2.24	0

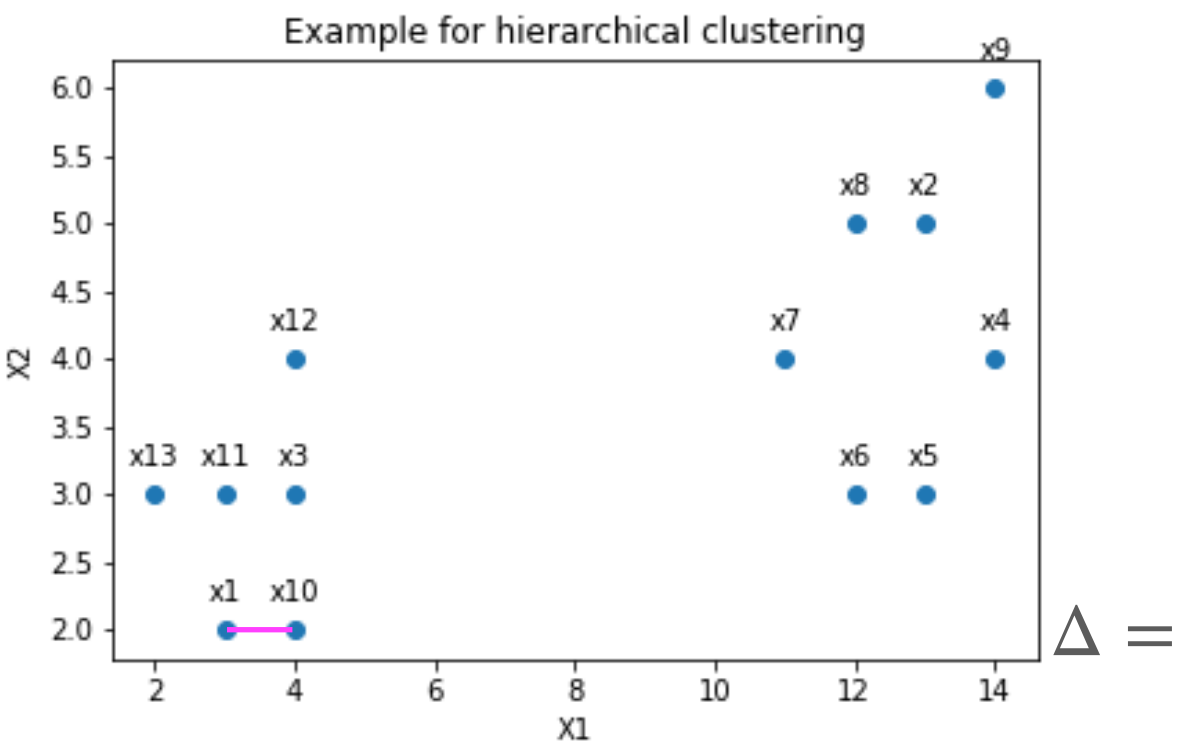
AGGLOMERATIVE CLUSTERING ALGORITHM

AgglomerativeClustering(D, k) :

1. $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$
2. $\Delta = \{\delta(x_i, x_j) : x_i, x_j \in D\}$
3. repeat:
 4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$
 5. $C_{i,j} = C_i \cup C_j$
 6. $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$
 7. Update distance matrix Δ to reflect new clustering
4. Until $|\mathcal{C}| = k$

EXAMPLE

4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$



	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂	X ₁₃
X ₁	0												
X ₂	10.44	0											
X ₃	1.41	9.21	0										
X ₄	11.18	1.41	10.05	0									
X ₅	10.05	2	9	1.41	0								
X ₆	9.06	2.24	8	2.24	1	0							
X ₇	8.25	2.24	7.07	3	2.24	1.41	0						
X ₈	9.49	1	8.25	2.23	2.24	2	1.41	0					
X ₉	11.70	1.41	10.44	2	3.16	3.61	3.61	2.24	0				
X ₁₀	1	9.49	1	10.20	9.06	8.06	7.28	8.54	10.77	0			
X ₁₁	1	10.20	1	11.05	10	9	8.06	9.22	11.40	1.41	0		
X ₁₂	2.24	9.06	1	10	9.06	8.06	7	8.06	10.20	2	1.41	0	
X ₁₃	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	2.24	1	2.24	0

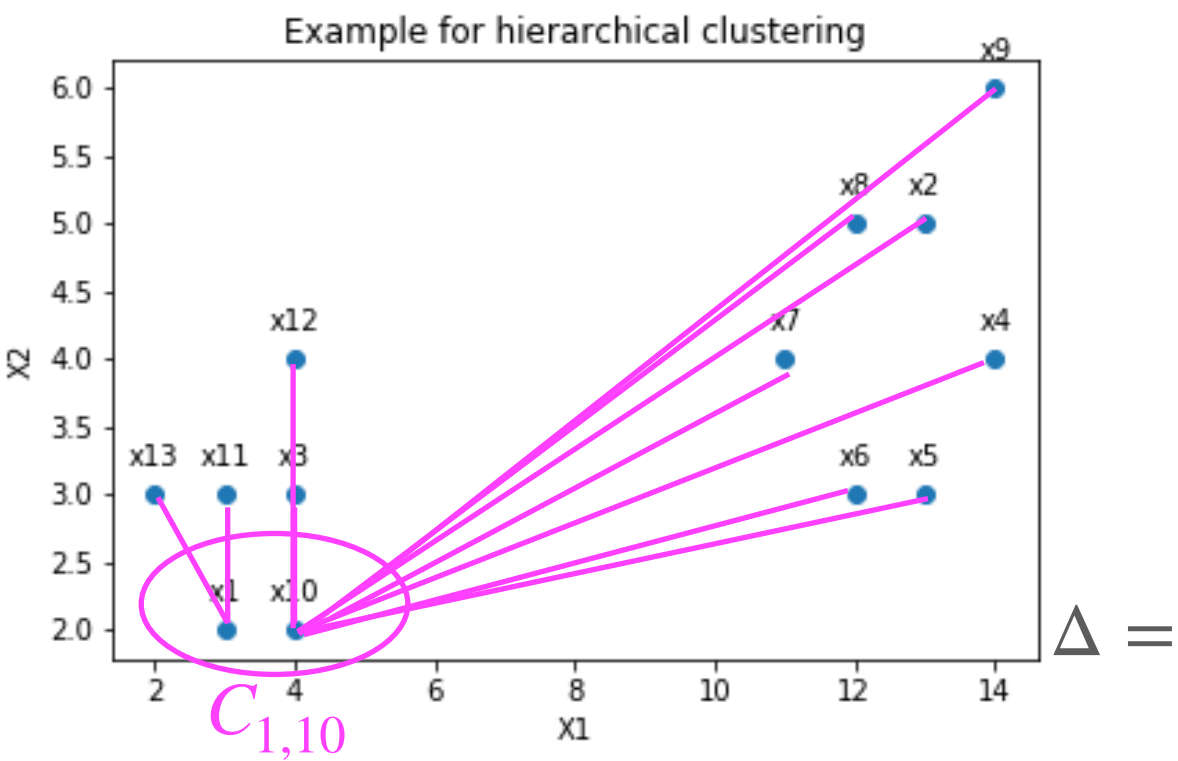
AGGLOMERATIVE CLUSTERING ALGORITHM

AgglomerativeClustering(D, k) :

1. $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$
2. $\Delta = \{\delta(x_i, x_j) : x_i, x_j \in D\}$
3. repeat:
 4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$
 5. $C_{i,j} = C_i \cup C_j$
 6. $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$
 7. Update distance matrix Δ to reflect new clustering
8. Until $|\mathcal{C}| = k$

EXAMPLE

Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$



	$\{x_1, x_{10}\}$	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{11}	x_{12}	x_{13}
$\{x_1, x_{10}\}$	0											
x_2	9.49	0										
x_3	1	9.21	0									
x_4	10.20	1.41	10.05	0								
x_5	9.06	2	9	1.41	0							
x_6	3.61	2.24	8	2.24	1	0						
x_7	7.28	2.24	7.07	3	2.24	1.41	0					
x_8	8.54	1	8.25	2.23	2.24	2	1.41	0				
x_9	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0			
x_{11}	1	10.20	1	11.05	10	9	8.06	9.22	11.40	0		
x_{12}	2	9.06	1	10	9.06	8.06	7	8.06	10.20	1.41	0	
x_{13}	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	1	2.24	0

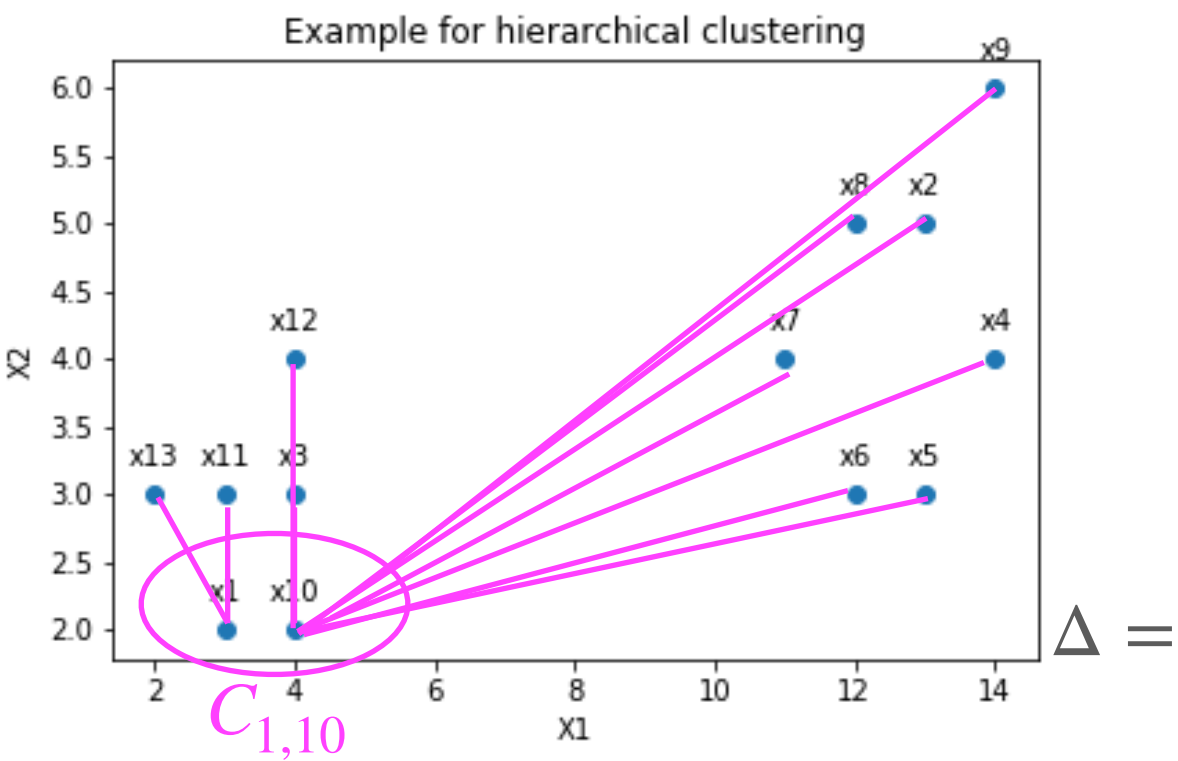
$$C_{i,j} = C_i \cup C_j$$

$$\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$$

Update Δ

EXAMPLE

Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$



$\Delta =$

	$\{x_1, x_{10}\}$	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{11}	x_{12}	x_{13}
$\{x_1, x_{10}\}$	0											
x_2	9.49	0										
x_3	1	9.21	0									
x_4	10.20	1.41	10.05	0								
x_5	9.06	2	9	1.41	0							
x_6	3.61	2.24	8	2.24	1	0						
x_7	7.28	2.24	7.07	3	2.24	1.41	0					
x_8	8.54	1	8.25	2.23	2.24	2	1.41	0				
x_9	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0			
x_{11}	1	10.20	1	11.05	10	9	8.06	9.22	11.40	0		
x_{12}	2	9.06	1	10	9.06	8.06	7	8.06	10.20	1.41	0	
x_{13}	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	1	2.24	0

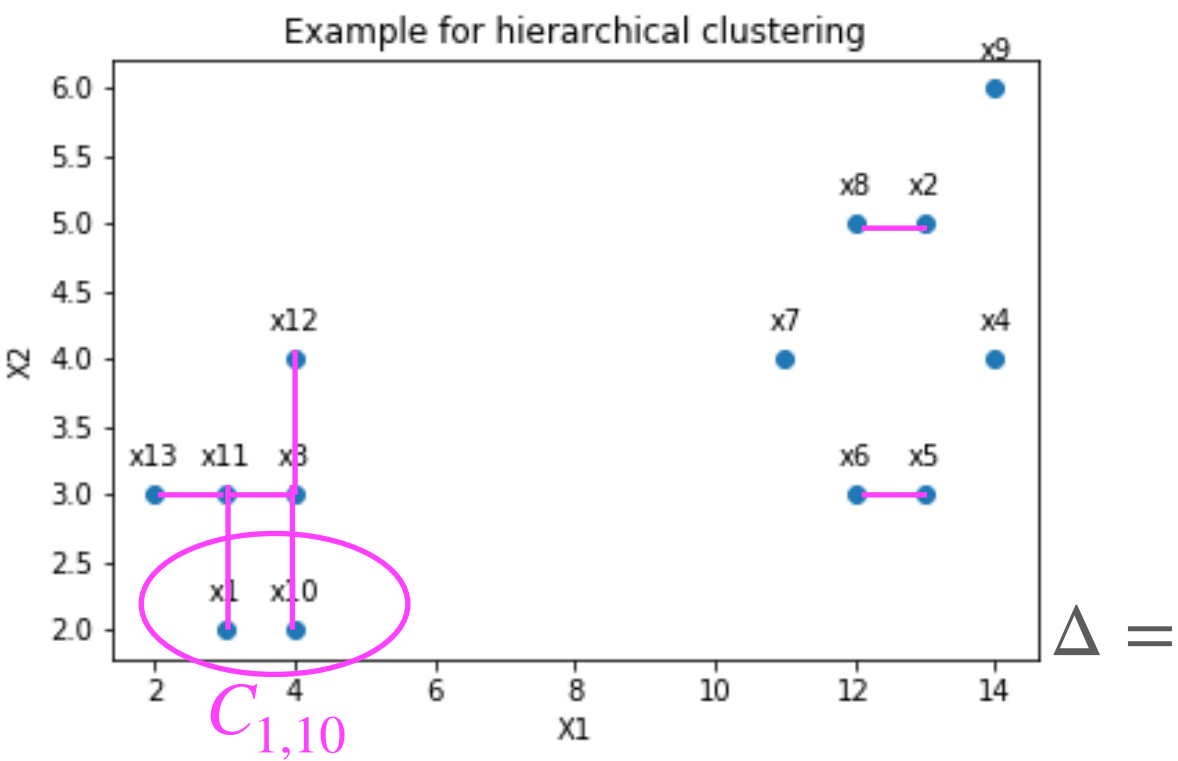
AGGLOMERATIVE CLUSTERING ALGORITHM

AgglomerativeClustering(D, k) :

1. $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$
2. $\Delta = \{\delta(x_i, x_j) : x_i, x_j \in D\}$
3. repeat:
 4. Find the closest pair of clusters $C_i, C_j \in \mathcal{C}$
 5. $C_{i,j} = C_i \cup C_j$
 6. $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{i,j}\}$
 7. Update distance matrix Δ to reflect new clustering
8. Until $|\mathcal{C}| = k$

EXAMPLE

Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$

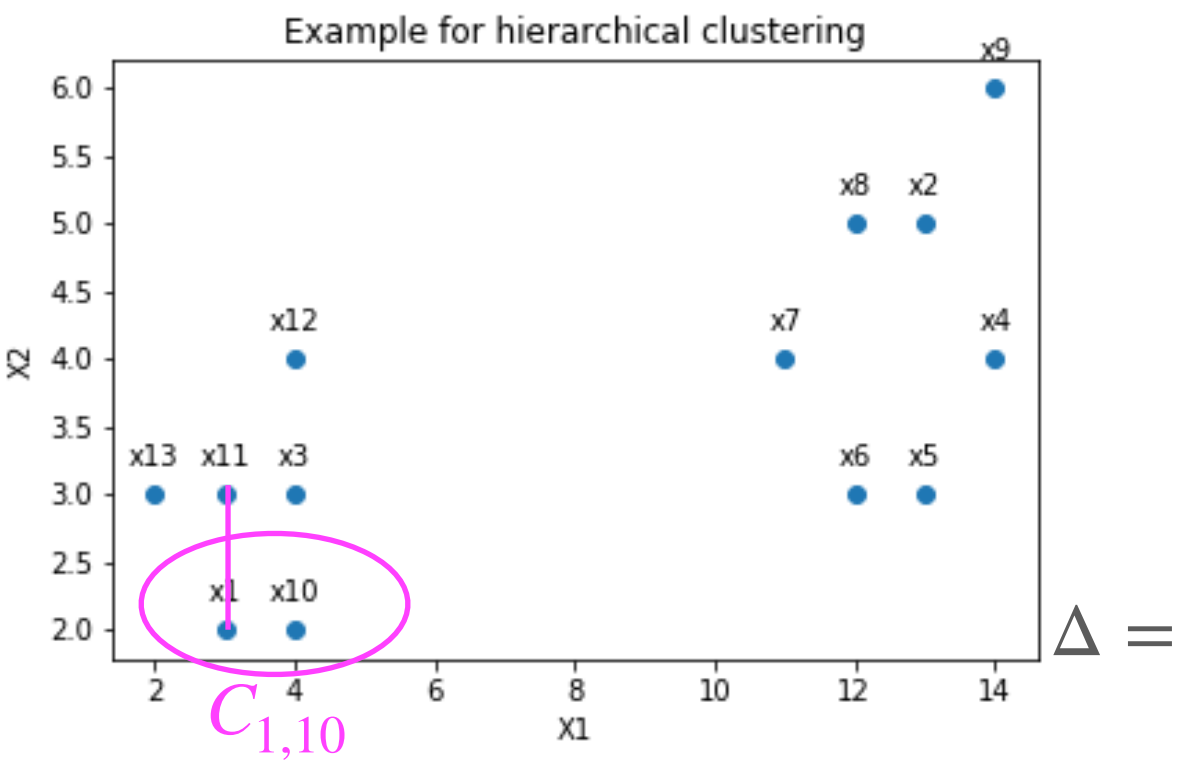


$\Delta =$

	$\{X_1, X_{10}\}$	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{11}	X_{12}	X_{13}
$\{X_1, X_{10}\}$	0											
X_2	9.49	0										
X_3	1	9.21	0									
X_4	10.20	1.41	10.05	0								
X_5	9.06	2	9	1.41	0							
X_6	3.61	2.24	8	2.24	1	0						
X_7	7.28	2.24	7.07	3	2.24	1.41	0					
X_8	8.54	1	8.25	2.23	2.24	2	1.41	0				
X_9	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0			
X_{11}	1	10.20	1	11.05	10	9	8.06	9.22	11.40	0		
X_{12}	2	9.06	1	10	9.06	8.06	7	8.06	10.20	1.41	0	
X_{13}	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	1	2.24	0

EXAMPLE

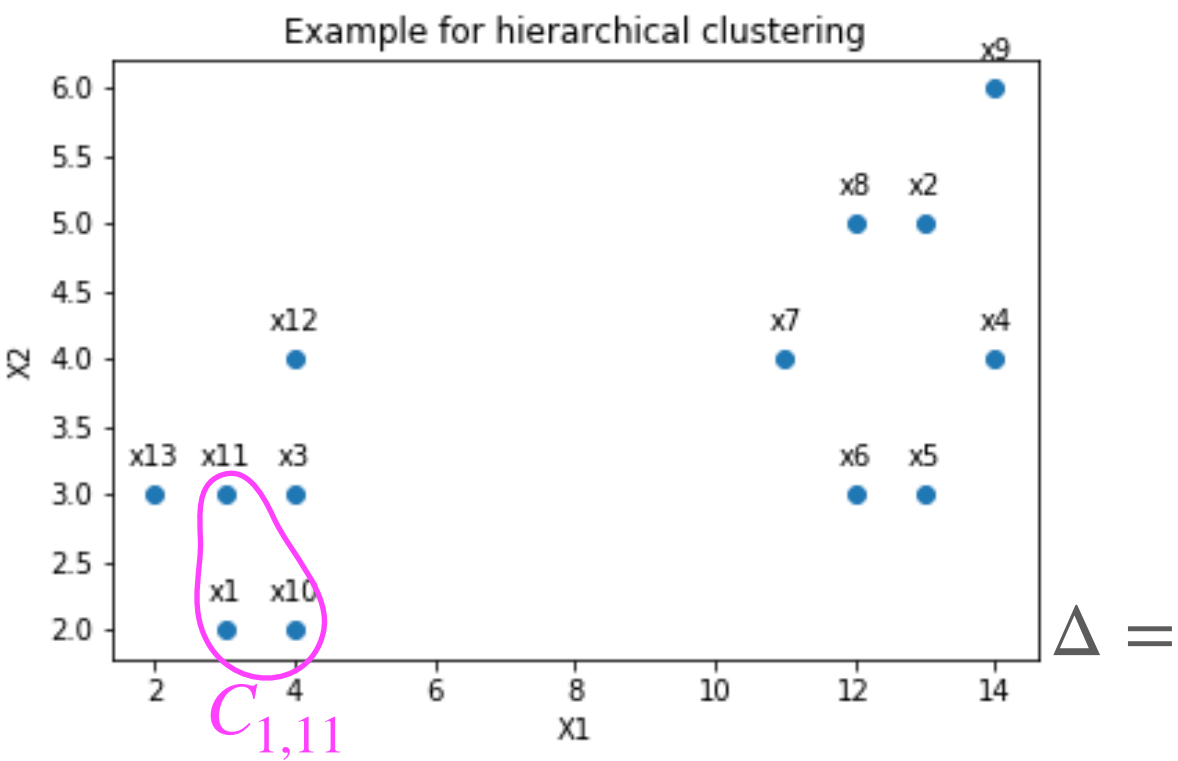
Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$



	$\{X_1, X_{10}\}$	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{11}	X_{12}	X_{13}
$\{X_1, X_{10}\}$	0											
X_2	9.49	0										
X_3	1	9.21	0									
X_4	10.20	1.41	10.05	0								
X_5	9.06	2	9	1.41	0							
X_6	3.61	2.24	8	2.24	1	0						
X_7	7.28	2.24	7.07	3	2.24	1.41	0					
X_8	8.54	1	8.25	2.23	2.24	2	1.41	0				
X_9	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0			
X_{11}	1	10.20	1	11.05	10	9	8.06	9.22	11.40	0		
X_{12}	2	9.06	1	10	9.06	8.06	7	8.06	10.20	1.41	0	
X_{13}	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	1	2.24	0

EXAMPLE

Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$

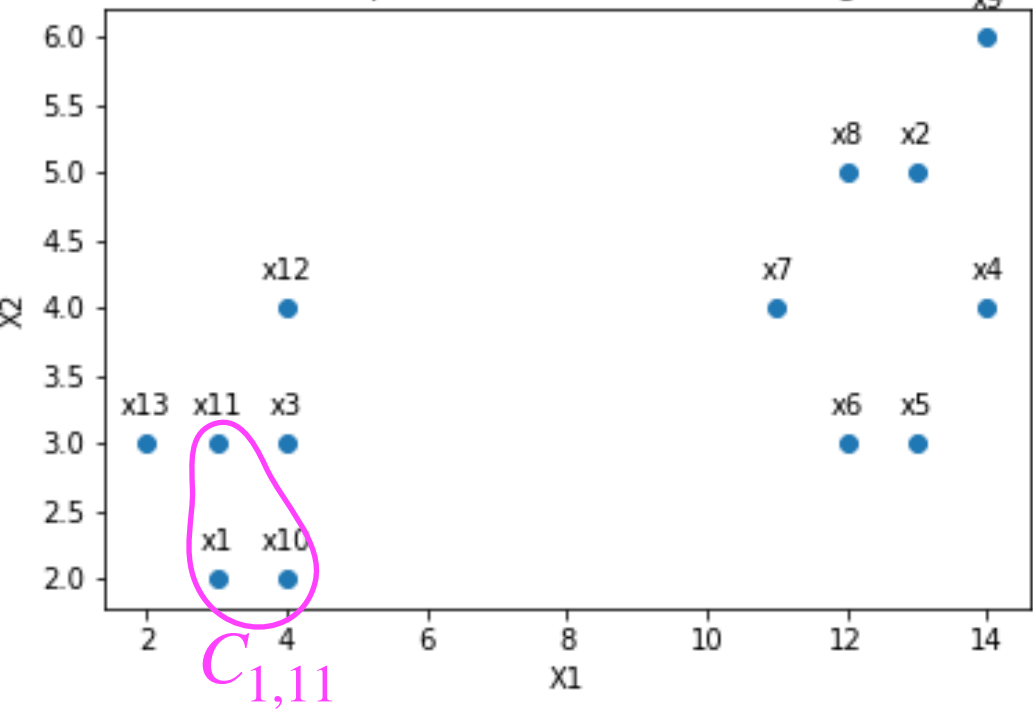


	{X ₁ ,X ₁₀ }	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₁	X ₁	X ₁₃
{X ₁ ,X ₁₀ }	0											
X ₂	9.49	0										
X ₃	1	9.21	0									
X ₄	10.20	1.41	10.05	0								
X ₅	9.06	2	9	1.41	0							
X ₆	3.61	2.24	8	2.24	1	0						
X ₇	7.28	2.24	7.07	3	2.24	1.41	0					
X ₈	8.54	1	8.25	2.23	2.24	2	1.41	0				
X ₉	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0			
X ₁₁	1	10.20	1	11.05	10	9	8.06	9.22	11.40	0		
X ₁₂	2	9.06	1	10	9.06	8.06	7	8.06	10.20	1.41	0	
X ₁₃	1.41	11.18	2	12.04	11	10	9.06	10.20	12.37	1	2.24	0

EXAMPLE

Single linkage: $\delta(C_i, C_j) = \min\{\delta(x, y) \mid x \in C_i, y \in C_j\}$

Example for hierarchical clustering



$\Delta =$

	$\{x_1, x_{10}, x_{11}\}$	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{12}	x_{13}
$\{x_1, x_{10}, x_{11}\}$	0										
x_2	9.49	0									
x_3	1	9.21	0								
x_4	10.20	1.41	10.05	0							
x_5	9.06	2	9	1.41	0						
x_6	3.61	2.24	8	2.24	1	0					
x_7	7.28	2.24	7.07	3	2.24	1.41	0				
x_8	8.54	1	8.25	2.23	2.24	2	1.41	0			
x_9	10.77	1.41	10.44	2	3.16	3.61	3.61	2.24	0		
x_{12}	1.41	9.06	1	10	9.06	8.06	7	8.06	10.20	0	
x_{13}	1	11.18	2	12.04	11	10	9.06	10.20	12.37	2.24	0

DIFFERENT DISTANCE MEASURES WILL AFFECT RESULTS

- L_1 norm, L_2 norm, etc...
- Single Linkage
- Complete Linkage
- Group Average
- Mean Distance
- Minimum Variance/Ward's Method