# Hierarchical Trees & Clustering

**Data**
$$X = [x_1, \ldots, x_p]$$
$$n \times p$$

$$x_i \in \mathbb{R}^n$$

Ex. $n = 2$
$p = 13$

$\mathbb{R}^n$



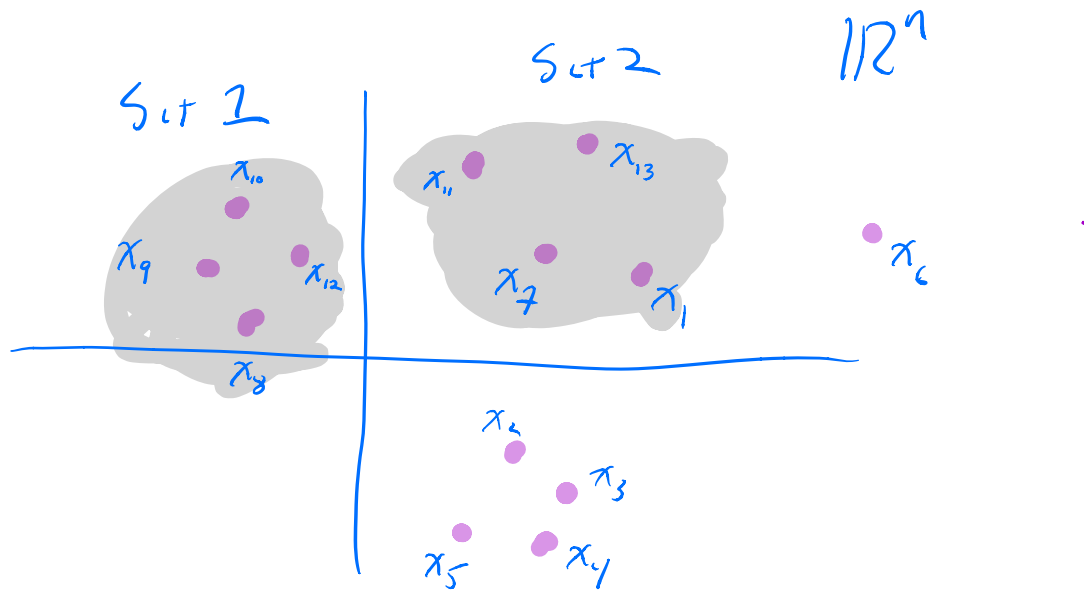**#1 Key Idea** Distance between observations

$$\boxed{dist(x_i, x_j)}$$

**Examples**

- Euclidean Distance in $\mathbb{R}^n$

- Correlation - based

correlation $\approx 1 \Rightarrow$ distance small
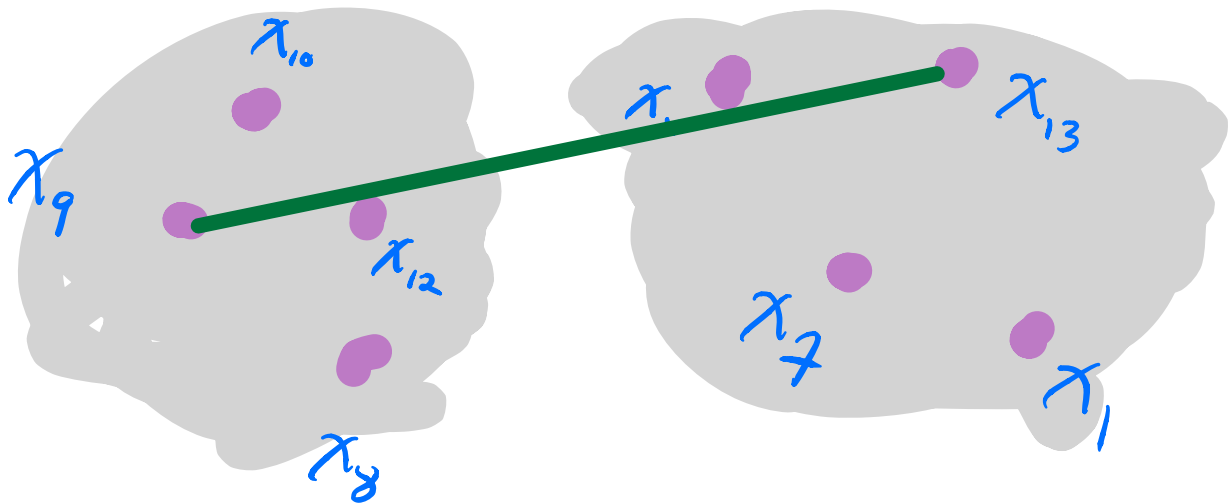correlation $\approx 0 \Rightarrow$ distance large

#2 Key Idea     Distance between sets
                of observations
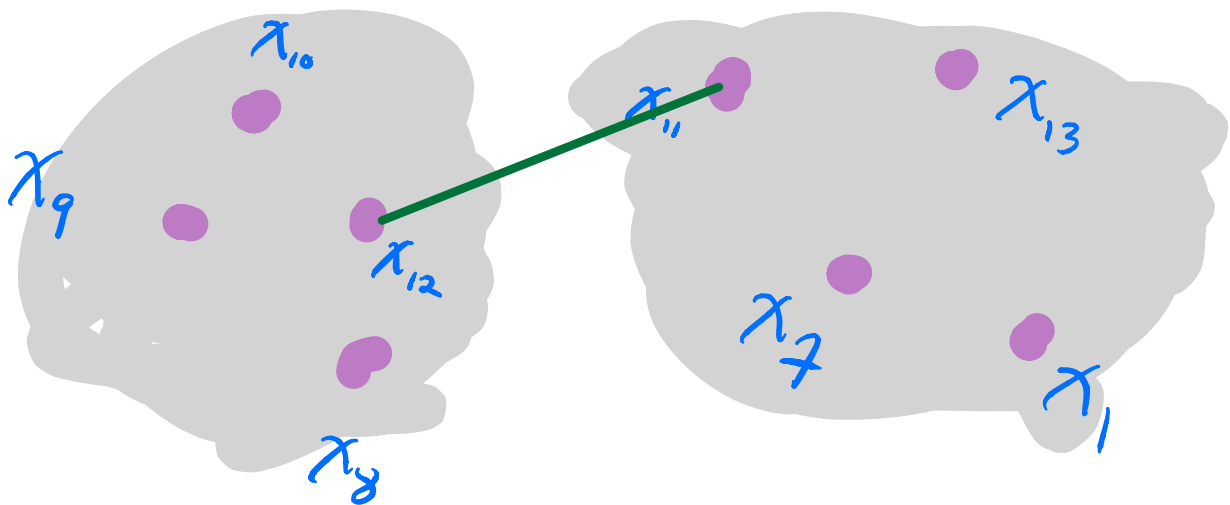
Set 1          Set 2          $\mathbb{R}^n$

$x_{10}$   $x_{11}$   $x_{13}$

$x_9$   $x_{12}$           $x_6$

         $x_7$      $x_1$

$x_8$

      $x_2$

         $x_3$

   $x_5$   $x_4$

dist ( Set 1, Set 2 )  ?

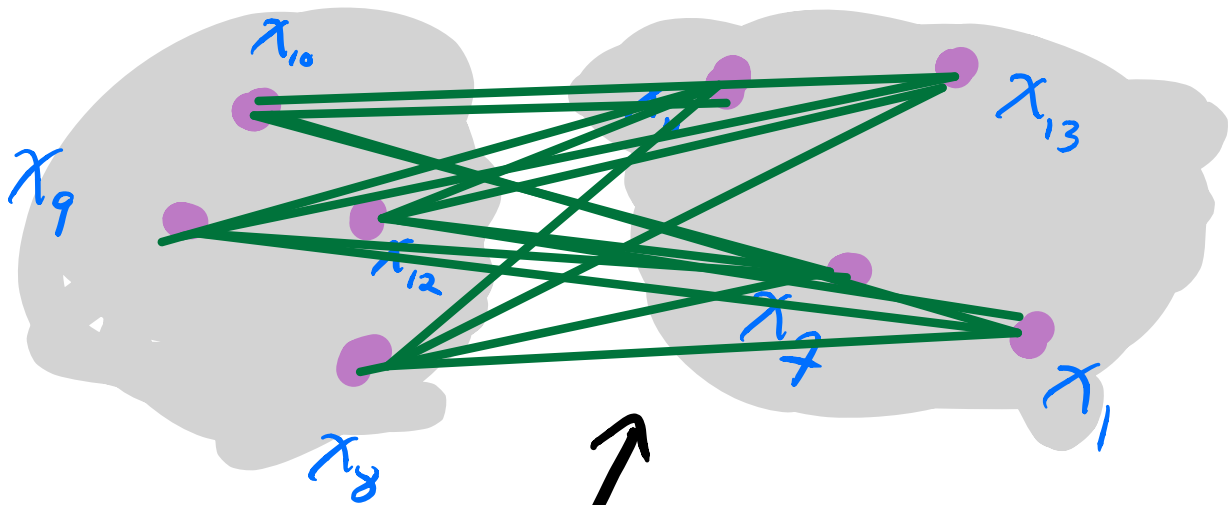# Complete

## Maximal Pair-wise Distance

# Single

## Minimal Pair-wise Distance

# Average
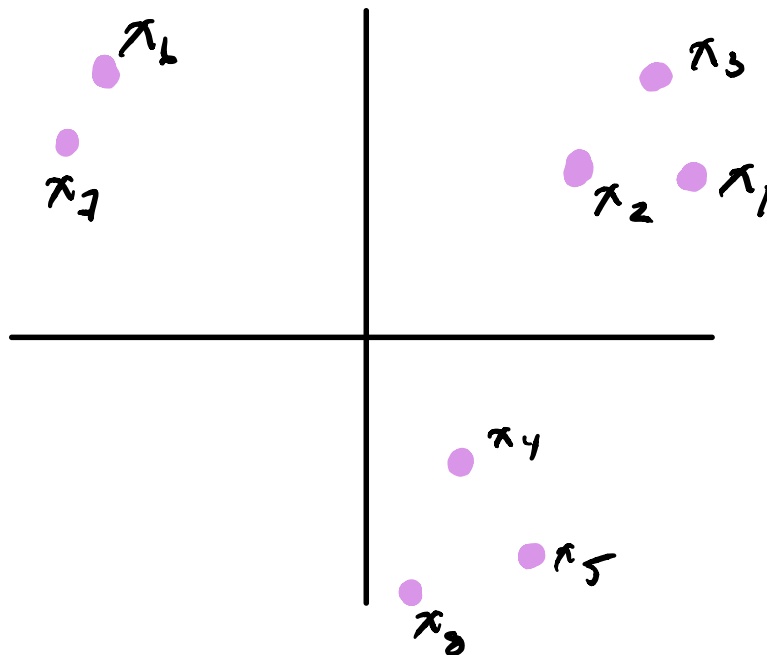
## Average of Pair-wise Distances



Average all these!

# Centroid
## Distance between Centroids
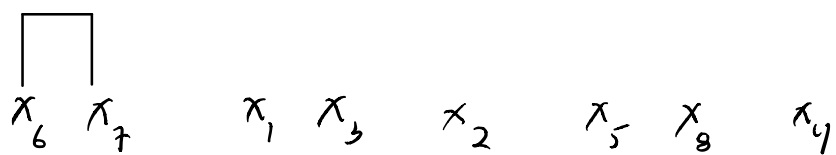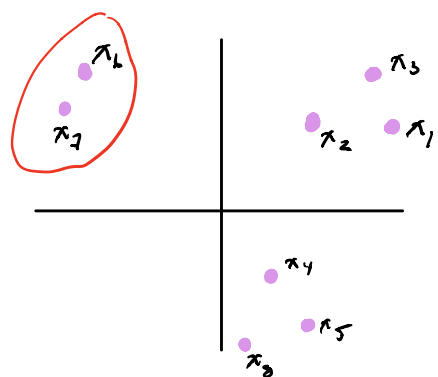


$x_{10}$

$x_9$

$x_{12}$

$x_8$

$x_{11}$

$x_{13}$

$x_7$

$x_1$

Centroid

Euclidian Distance
Central Method

$X_6$  $X_7$     $X_1$  $X_3$    $X_2$     $X_5$  $X_8$    $X_4$

$x_6$ $x_7$    $x_1$ $x_3$   $x_2$   $x_5$ $x_8$   $x_4$

$x_6$ $x_7$ $x_1$ $x_3$ $x_2$ $x_5$ $x_8$ $x_4$

$x_6$    $x_7$      $x_1$    $x_3$     $x_2$      $x_5$    $x_8$     $x_4$

$x_6$

$x_7$

$x_3$

$x_2$ $x_1$

$x_4$

$x_5$

$x_8$

$X_6$ $X_7$  $X_1$ $X_3$ $X_2$  $X_5$ $X_8$  $X_4$

$x_6$  $x_7$  $x_3$  $x_2$  $x_1$  $x_4$  $x_5$  $x_8$

$X_6$  $X_7$  $X_1$  $X_3$  $X_2$  $X_5$  $X_8$  $X_4$

$x_6$    $x_7$      $x_1$    $x_3$    $x_2$      $x_5$    $x_8$    $x_4$

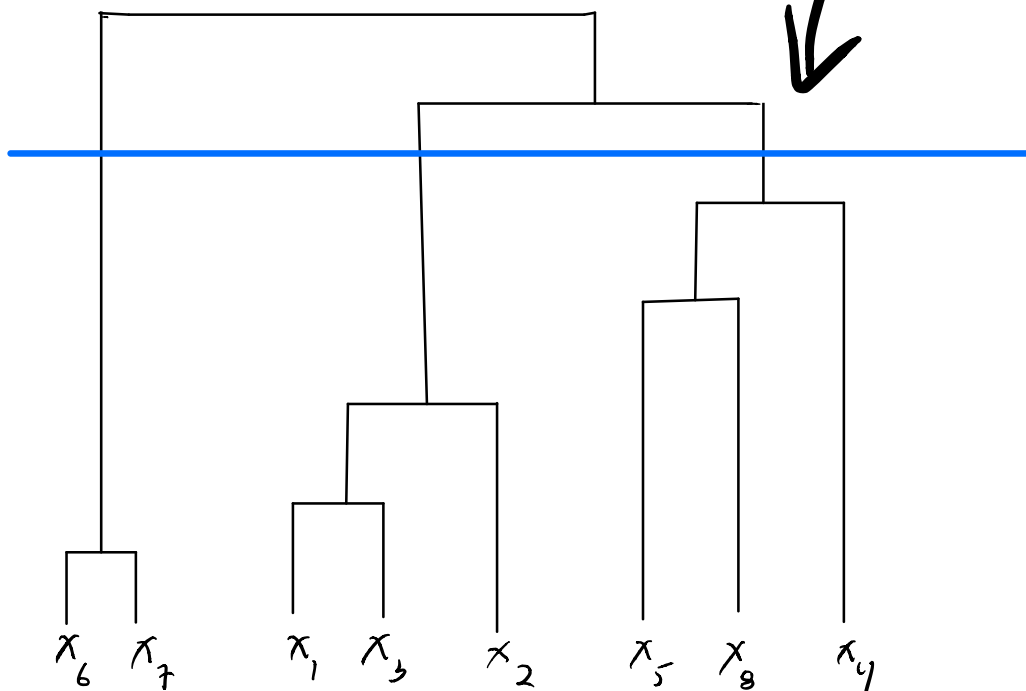$x_6$   $x_7$     $x_1$   $x_3$    $x_2$     $x_5$   $x_8$    $x_4$
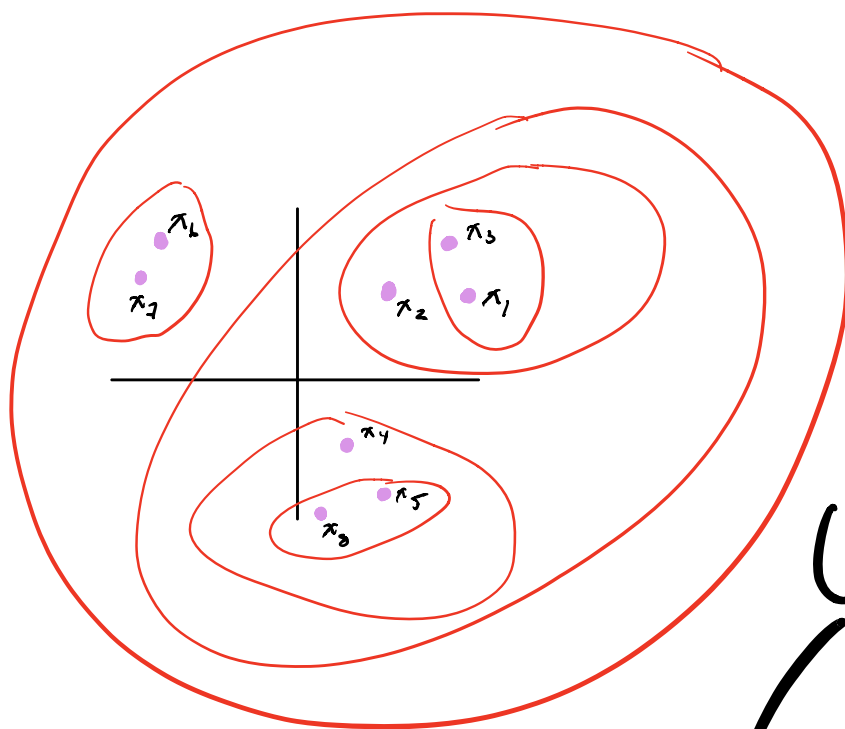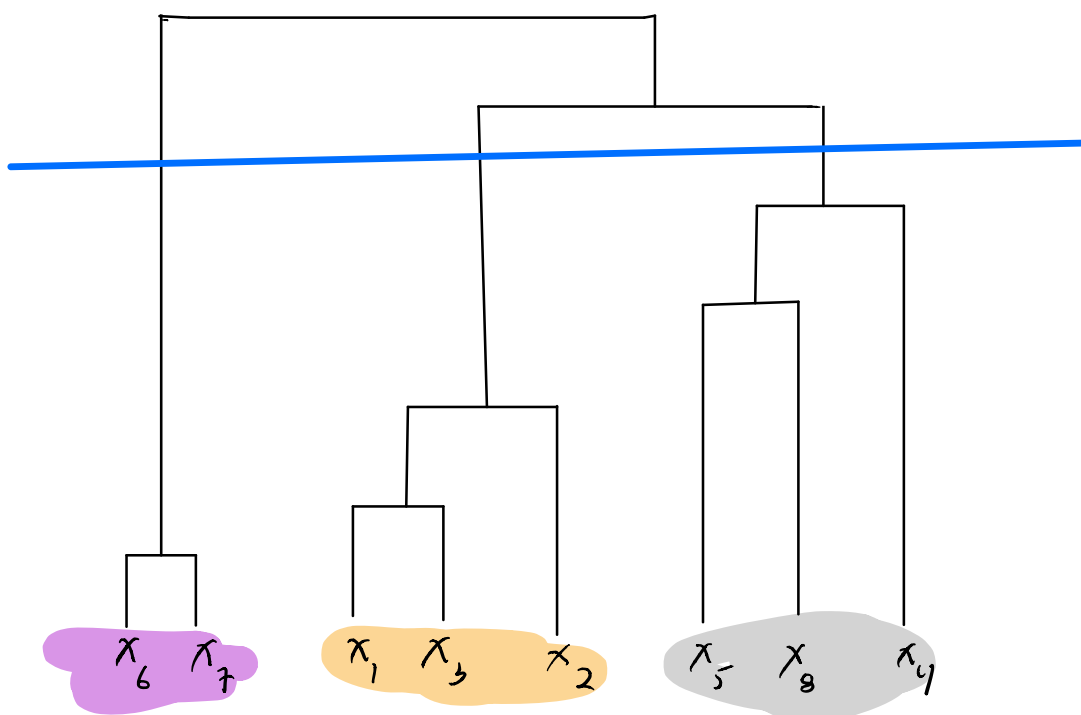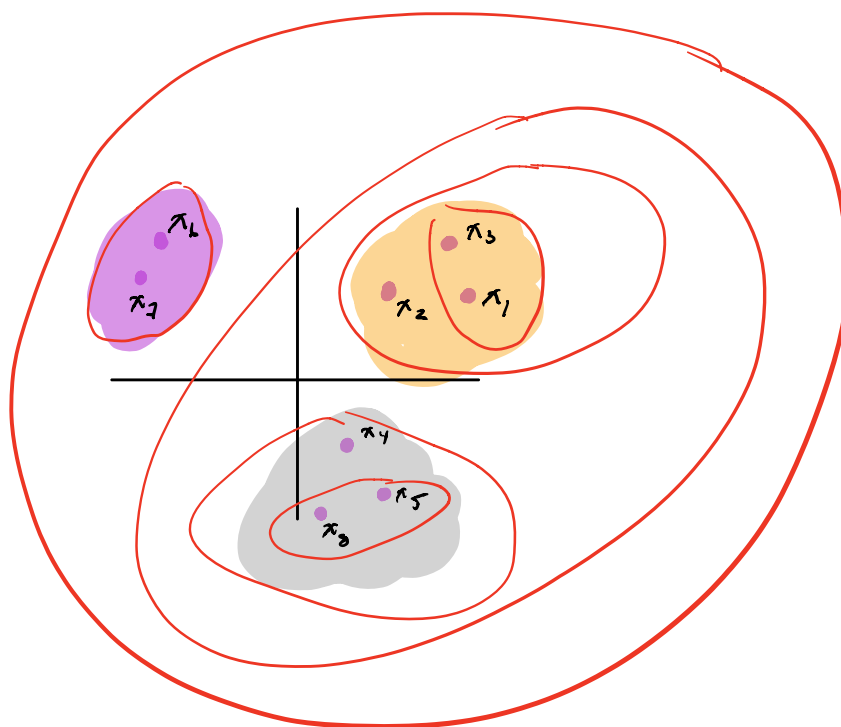
Cut

# Hierarchial Clustering

- Intuitive

- Visual

  Works in any dimension

- Sensitive to choices

  - Distance
  - Cluster type (e.g. Central vs Average)