

Introduction to cluster analysis

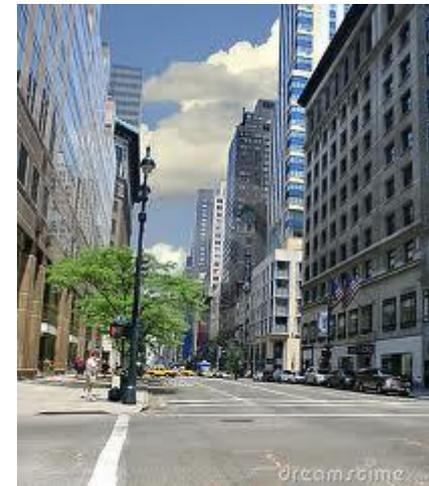
Lecture 05.01

What is Cluster Analysis?

Finding groups of objects such that the objects in each group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

Labeling objects with group label

- Classes – conceptually meaningful groups of objects that share common characteristics
- Humans are skilled at dividing objects into groups (clustering) and assigning new objects to one of the groups (classification)
- Clusters are potential classes, and **cluster analysis** is a technique for **automatically discovering classes from unlabeled data**

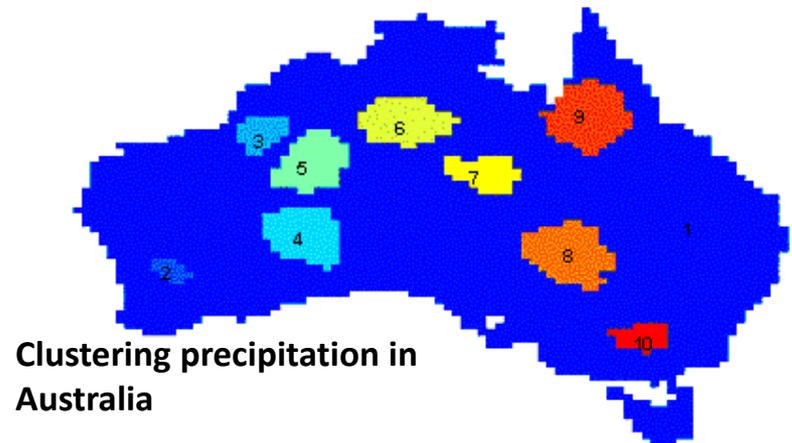


Applications of Cluster Analysis

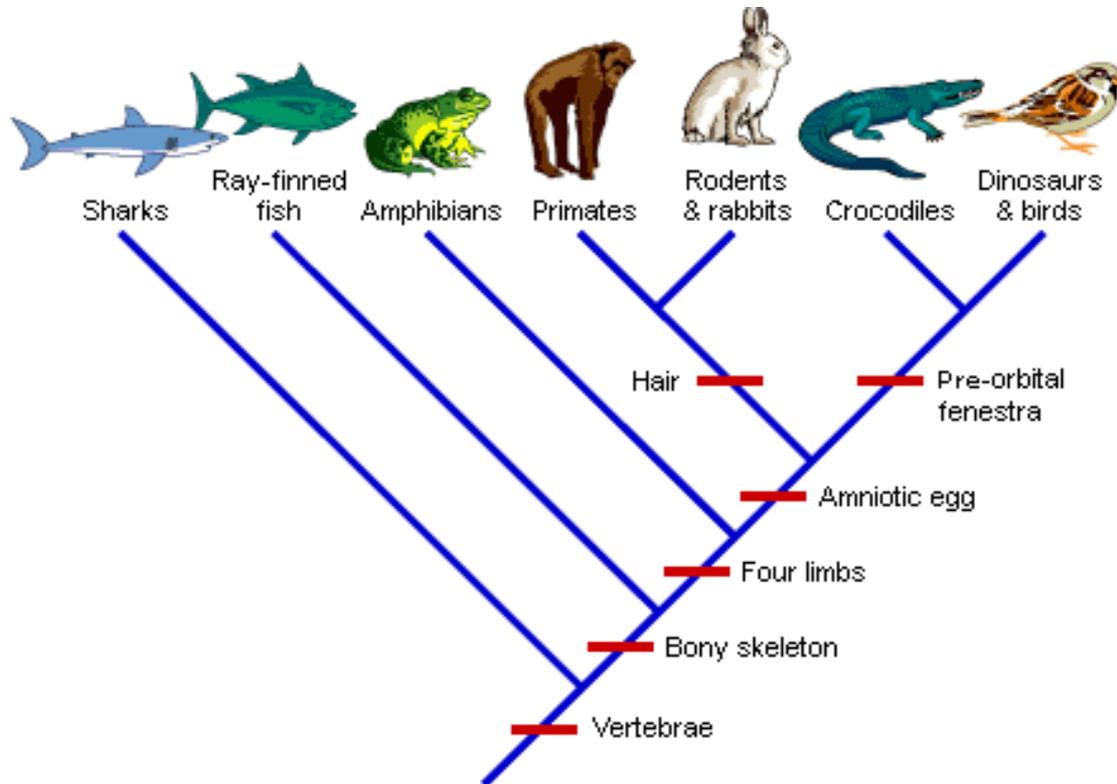
- **Clustering for Understanding**
 - Group related documents for browsing
 - Group genes and proteins that have similar functionality
 - Group stocks with similar price fluctuations
 - Segment customers into a small number of groups for additional analysis and marketing activities.

	<i>Discovered Clusters</i>	<i>Industry Group</i>
1	Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN, Sun-DOWN	Technology1-DOWN
2	Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN	Technology2-DOWN
3	Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN	Financial-DOWN
4	Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP	Oil-UP

- **Clustering for Summarization**
 - Reduce the size of large data sets



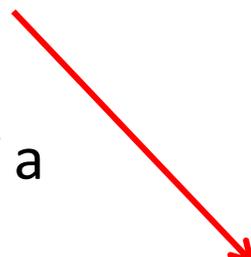
Grouping animals into clusters: biological systematics



- Grouping animals into hierarchical groups to better understand evolution

Grouping documents into clusters: information retrieval

- Grouping WEB query results into small number of clusters, each capturing a particular aspect of a query



Search for *tiger*

[Giant Tiger - Main Page](#)

www.gianttiger.com/

Welcome to Giant **Tiger**, your all Canadian family

Searches related to **tiger**

[tiger pictures](#)

[tiger woods](#)

[tiger animal](#)

[tiger tiger](#)

[tiger beer](#)

[tiger facts](#)

[tiger direct](#)

[tiger information](#)



1 2 3 4 5 6 7 8

[Advanced search](#)

[Search Help](#)

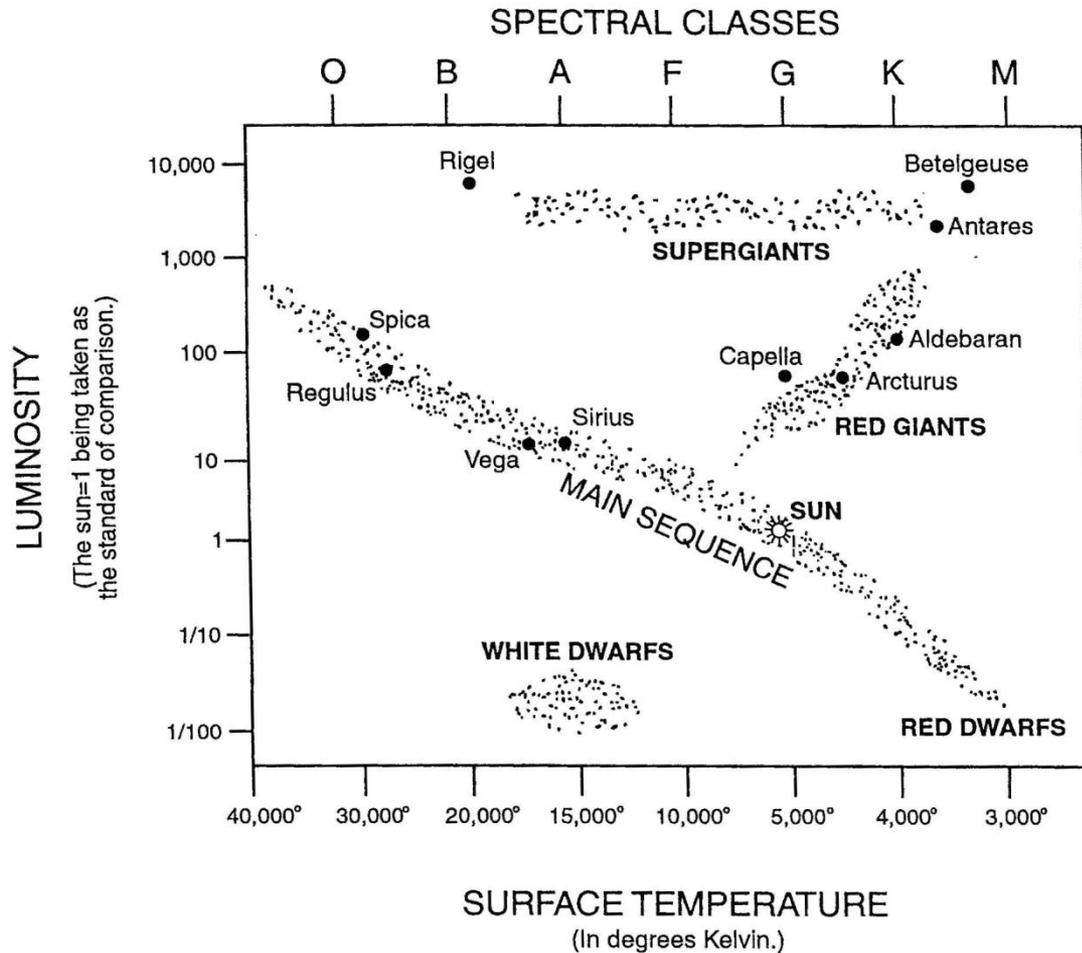
[Give us](#)

[Google Home](#)

[Advertising Programs](#)

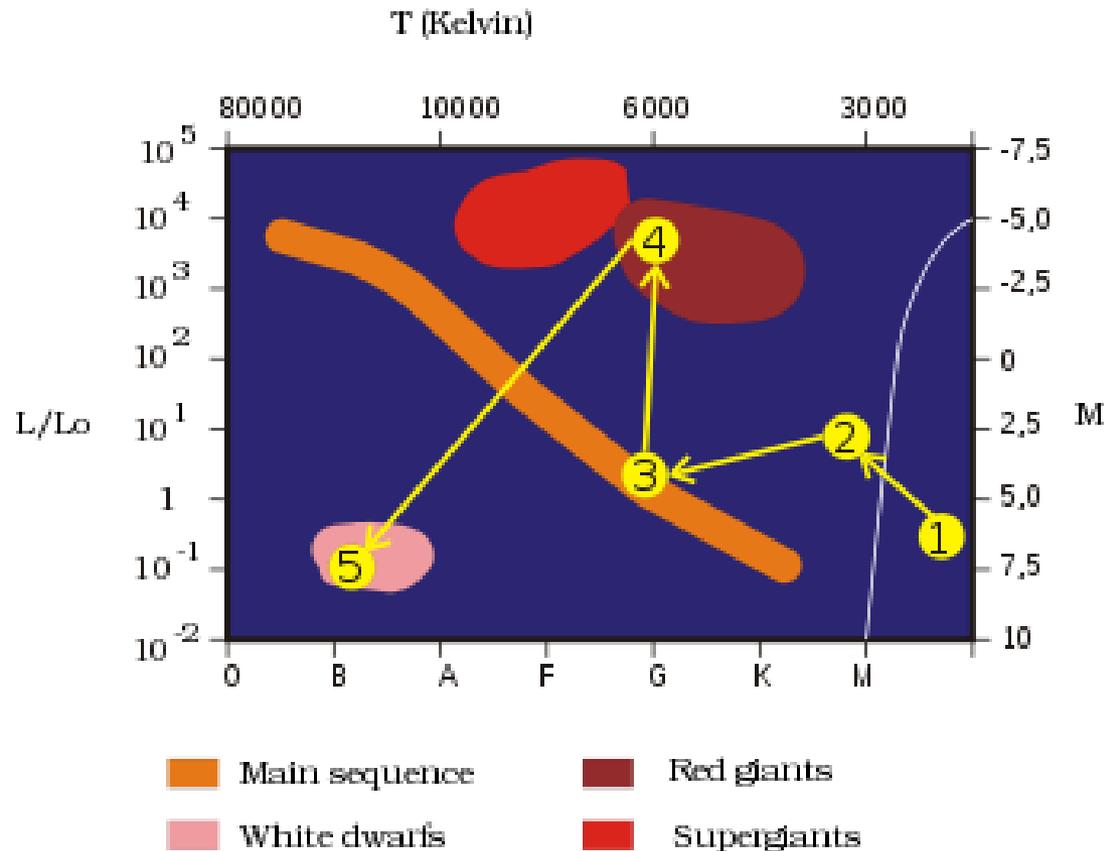
[About Google](#)

Clustering leads to discoveries: Galaxies in 2 dimensions



The Hertzsprung-Russel diagram clusters stars by temperature and luminosity

Clustering leads to discoveries:



Galaxies evolution

Main sequence stars generate energy by fusing Hydrogen to Helium

When the hydrogen is used up, Helium fusion occurs, the star expands -> red giant

The outer layer of gases is stripped away, the star cools -> white dwarf

Automatic clustering

- Discovering groups (classes) of objects from **unlabeled** data
- *Unsupervised learning*

Formulation for computer program

- Computer needs to know:
 - What is similarity (dissimilarity=distance)?
 - What exactly to look for?
- We need to:
 - define similarity as a numeric value
 - define the notion of a cluster
 - prescribe the precise algorithm for finding defined clusters

Part 1

DEFINE SIMILARITY/DISTANCE

Numeric *proximity* (similarity or distance) between data records

- Combination of proximity measures for each attribute
- Each attribute is a separate and independent (in this approach) *dimension* of the data
- First step: translate all fields into numeric variables, to make similarities (distances) numeric

Types of attributes

1. True measures (continuous)
2. Ranks (ordinal)
3. Categorical (nominal)

The distances are increasingly harder to convert into a numeric scale



How do we define the proximity measure for a single attribute of each type?

1. True measures

- True measures measure the value from a meaningful “0” point. The ratio between values is meaningful, and the distance is just an **absolute difference of values**.
- Examples: age, weight, length

2. Ordinal (Ranks)

- These values have an order, but the distance between different ranks is not defined

2. Ordinal (Ranks)

Example 1:

quality attribute of a product : {poor, fair, OK, good, wonderful}

Order is important, but exact difference between values is undefined

Solution: map the values of the attribute to successive integers

{poor=0, fair=1, OK=2, good=3, wonderful=4}

Dissimilarity

$$d(p,q) = |p - q| / (\max_d - \min_d)$$

e.g. $d(\text{wonderful}, \text{fair}) = |4-1| / (4-0) = .75$

Not always meaningful, but the best we can do

Similarity

$s(p,q) = 1 - d(p,q)$ e.g. $d(\text{wonderful}, \text{fair}) = .25$

2. Ordinal (Ranks)

Example 2:

Top 10 swimmers - 50m Fly				
1	KONOVALOV, Nikita	88	RUS	22.70
2	GOVOROV, Andriy	92	UKR	22.70
3	LEVEAUX, Amaury	85	FRA	22.74
4	CZERNIAK, Konrad	89	POL	22.77
5	KOROTYSHKIN, Evgeny	83	RUS	22.88
6	EIBLER, Steffen	87	GER	22.89
7	FESIKOV, Sergey	89	RUS	22.96
8	HEERSBRANDT, Francois	89	BEL	22.98
9	MUNOZ PEREZ, Rafael	88	ESP	23.07
10	JAMES, Antony	89	GBR	23.14

Distance between 3 and 1 (0.04 sec) is not the same as distance between 10 and 8 (0.16). It is better to use the attributes which contributed to this ranking

3. Categorical (nominal) attributes

- Each value is one of a set of unordered categories. We can only tell that $X \neq Y$, but not how much X is greater than Y .
- Example: ice cream pistachio is not equal to butter pecan, but we cannot tell which one is greater and which one is closer to black cherry ice cream
- The general approach: if equal then similarity = 1, if not equal then similarity = 0

Summary on proximity measures for a single attribute

Attribute type	Distance (dissimilarity)	Similarity
True measures	$d = x - y $	$s = -d$, $s = 1/(1+d)$, $s = 1 - (d - \min_d) / (\max_d - \min_d)$
Ordinal	$d = x - y / (n - 1)$ (values mapped to integers 0 to n-1 where n is the number of values)	$s = 1 - d$
Nominal	$d = 0$ if $x = y$ $d = 1$ if $x \neq y$	$s = 1$ if $x = y$ $s = 0$ if $x \neq y$

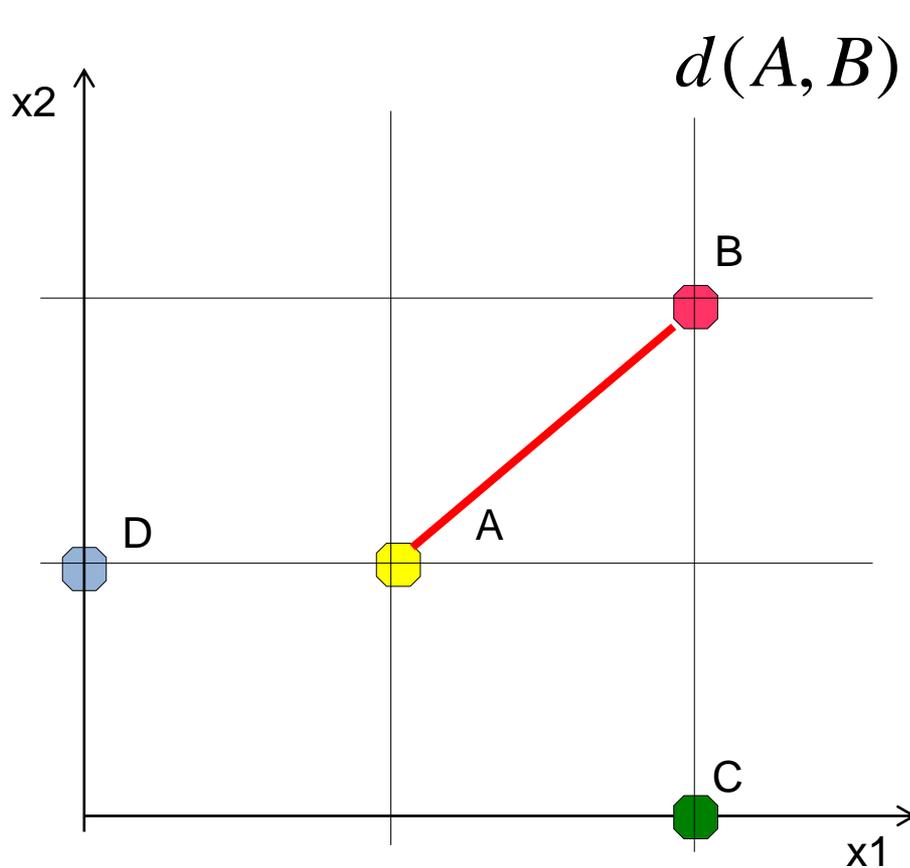
Combining measures for separate attributes into proximity measure for data records

- Hundreds of similarity measures were proposed
- We will look at:
 - Euclidean distance
 - Jaccard index
 - Tanimoto coefficient
 - Cosine similarity
 - Pearson similarity

Proximity measures for data records

1. Euclidean distance (all attributes are numeric)
2. Matching coefficients (all attributes are binary– from categorical attributes transformed into binary)
3. Cosine similarity (true values as vectors)

Distance 1. All attributes are numeric: Euclidean distance



$$d(A, B) = \sqrt{|A_X - B_X|^2 + |A_Y - B_Y|^2}$$

For N dimensions:

$$d(A, B) = \sqrt{\sum_{i=1}^N |A_i - B_i|^2}$$

Similarity:

$$s(A, B) = 1 / (1 + d(A, B))$$

It is hard to visualize points in more than 3 dimensions, but for computer it is not a problem.

Distance 2. Matching coefficients (all attributes are binary)

	Y	Y
X	M_{11}	M_{10}
X	M_{01}	M_{00}

M_{11} : number of attributes with value 1 in both X and Y

M_{10} : number of attributes with value 1 in X and 0 in Y

M_{01} : number of attributes with value 0 in X but 1 in Y

M_{00} : number of attributes with value 0 in both X and Y

Matching coefficients and Jaccard index

	Y	Y
X	M_{11}	M_{10}
X	M_{01}	M_{00}

Jaccard index is used for **asymmetric binary attributes**, where only value 1 is important

Simple Matching Coefficient

$$\begin{aligned} \text{SMC} &= \text{number of matches} / \text{number of attributes} \\ &= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) \end{aligned}$$

Jaccard Index

$$\begin{aligned} J &= \text{number of } M_{11} \text{ matches} / \text{number of not-both-zero attributes values} \\ &= (M_{11}) / (M_{01} + M_{10} + M_{11}) \end{aligned}$$

SMC and Jaccard example

$$x = (1 \quad 0)$$

$$y = (0 \quad 1 \quad 0 \quad 0 \quad 1)$$

$$SMC = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0 + 7) / 10 = 0.7$$

$$J = M_{11} / (M_{01} + M_{10} + M_{11}) = (0) / 3 = 0.0$$

The choice is application-dependent.

SMC and Jaccard example

$$x = (1 \quad 0)$$

$$y = (0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 1)$$

$$\text{SMC} = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0+7)/10=0.7$$

$$J = M_{11} / (M_{01} + M_{10} + M_{11}) = (0)/3=0.0$$

The choice is application-dependent.

Which measure to choose for:

Comparing documents by common words?

Comparing transactions by common items?

Comparing students by knowledge of 10 topics?

Tanimoto similarity coefficient

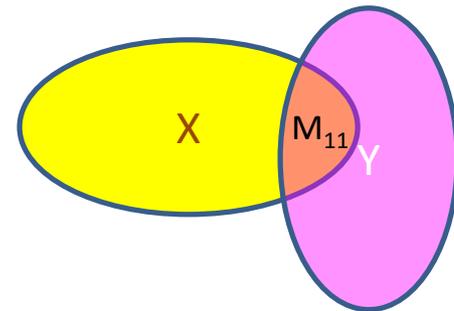
- **Jaccard** index is defined as the number of attributes with value 1 in both records, divided by the total number of records for which there is at least one 1 value: $J = M_{11} / (M_{01} + M_{10} + M_{11})$

- **Tanimoto** coefficient is similar but is defined in terms of set operations: it is an intersection over union of all attribute values without attributes for which both binary values are False(0):

$$T = M_{11} / (M_{-1} + M_{1-} - M_{11})$$

The formulas show that Jaccard and Tanimoto are **exactly the same!**

	Y	Y
X	M_{11}	M_{10}
X	M_{01}	M_{00}



$M_{11} \leftarrow$ intersection

$$(M_{01} + M_{10} + M_{11}) = (M_{-1} + M_{1-} - M_{11}) \leftarrow \text{union}$$

Distance 3. Cosine similarity

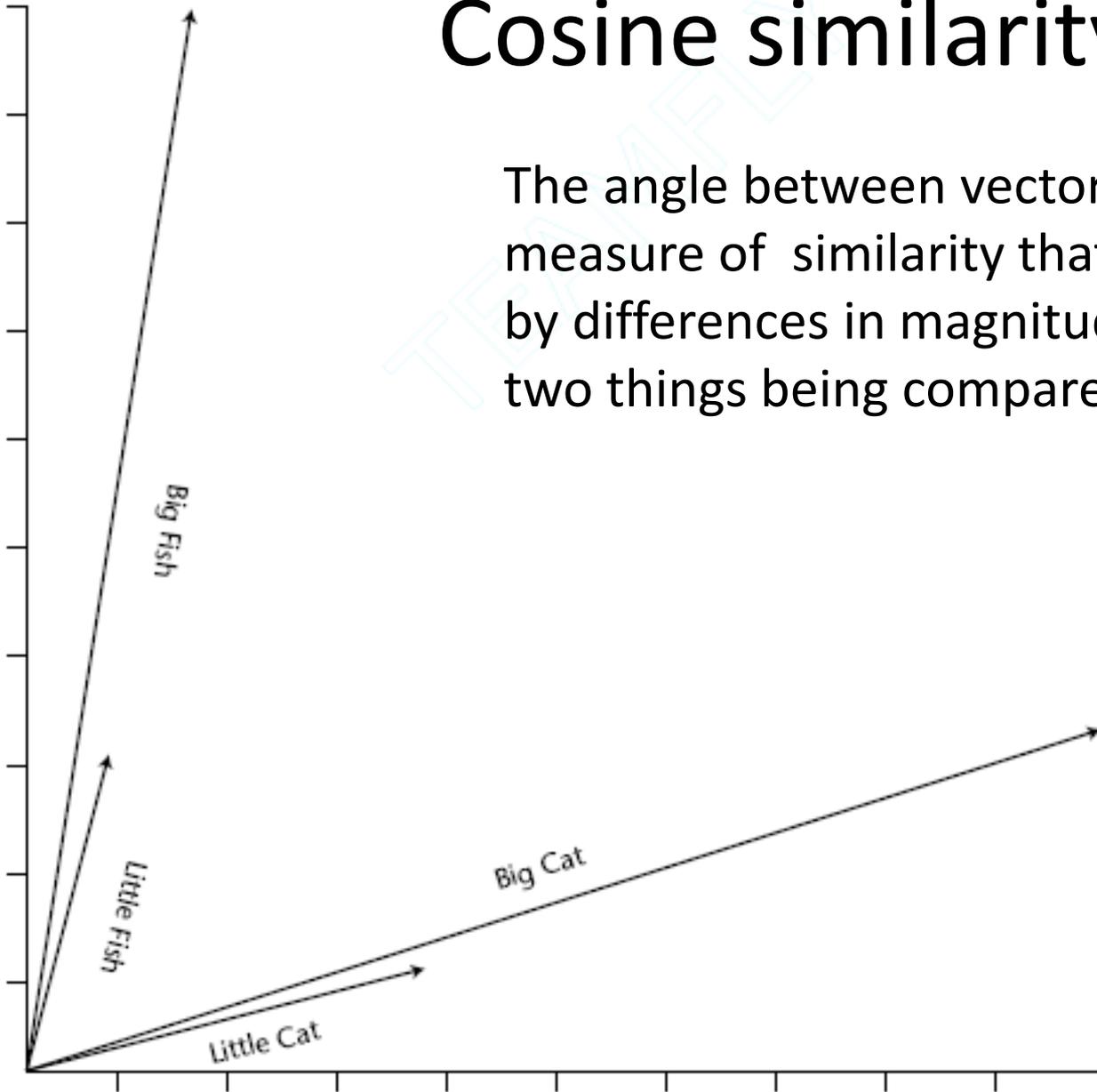
- Sometimes it makes more sense to consider two records closely associated because of similarities in the way the fields *within each record are related*
- Example: sardines should cluster with cod and tuna, while kittens cluster with cougars and lions, but if we use the Euclidean distance of body-part lengths, the sardine is closer to a kitten than it is to a catfish.

Cosine similarity

- Sometimes it makes more sense to consider two records closely associated because of similarities in the way the fields *within each record are related*
- Solution: use a different geometric interpretation. Instead of thinking of X and Y as *points in space*, think of them as *vectors and measure the angle between them*.
- In this context, a vector is the line segment connecting the origin of a coordinate system to the point described by the vector values.

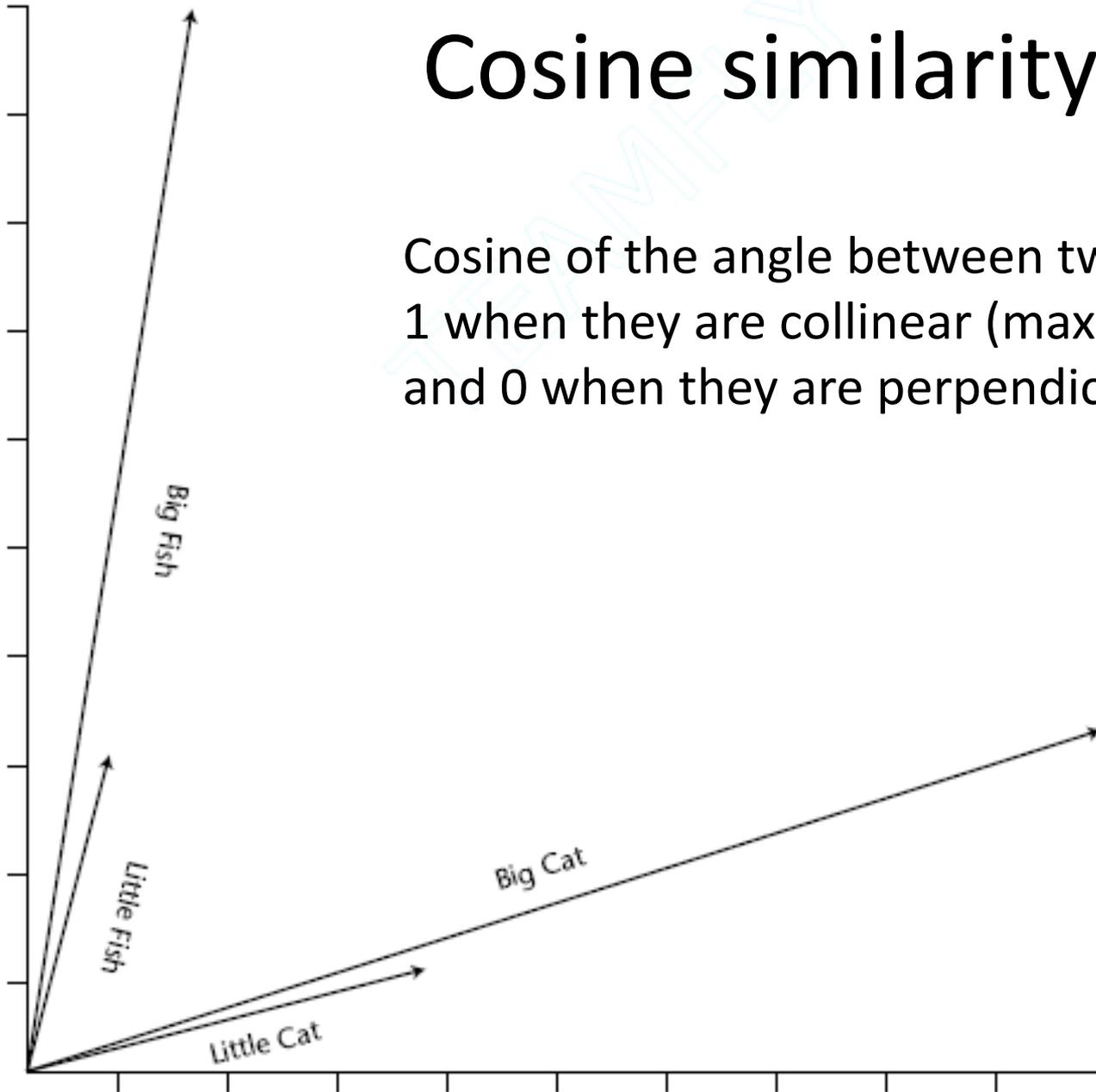
Cosine similarity

The angle between vectors provides a measure of similarity that is not influenced by differences in magnitude between the two things being compared



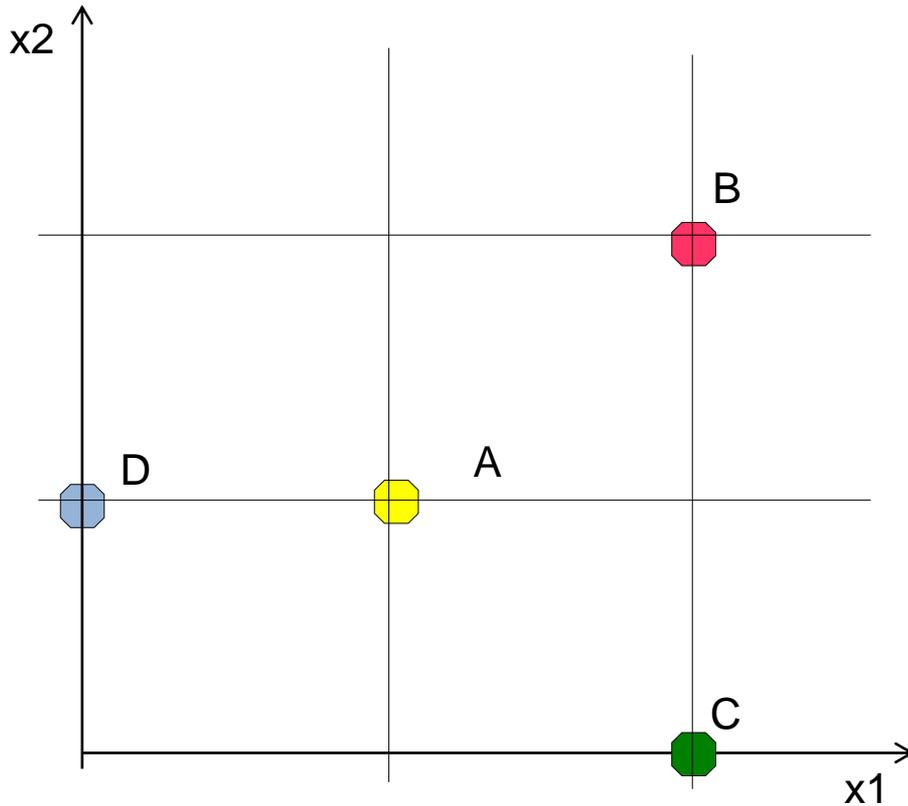
Cosine similarity

Cosine of the angle between two vectors is 1 when they are collinear (maximum similarity) and 0 when they are perpendicular



Cosine similarity

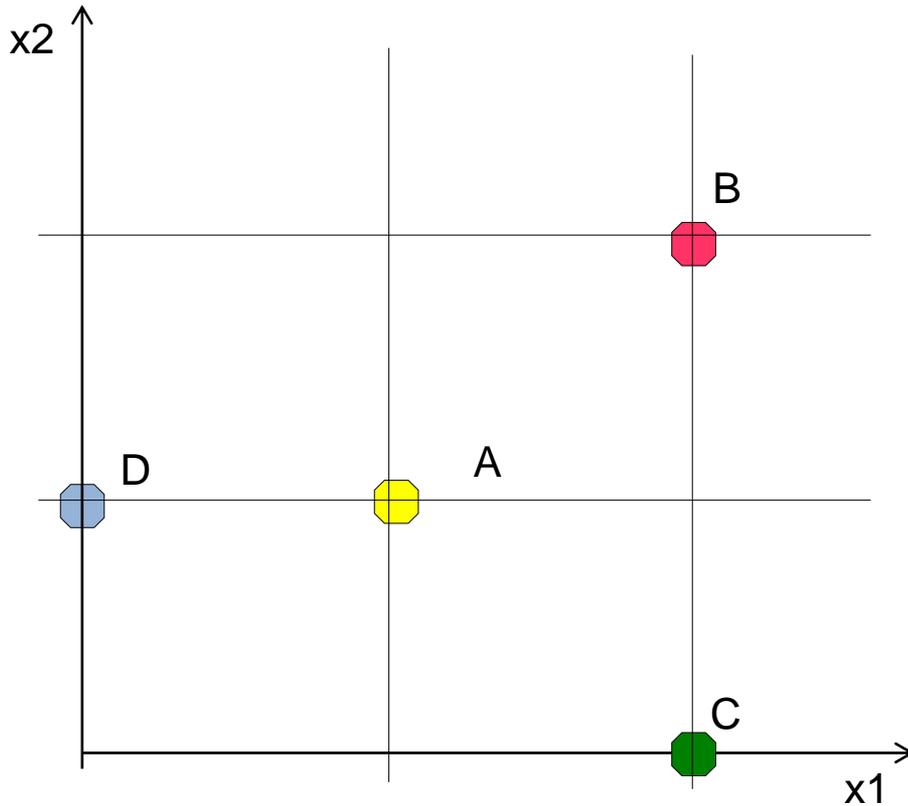
$$s(A,B) = \cos(\mathbf{A}, \mathbf{B}) = (\mathbf{A} \bullet \mathbf{B}) / \|\mathbf{A}\| \cdot \|\mathbf{B}\|$$



Dot-product of
vectors

Cosine similarity

$$s(A,B) = \cos(\mathbf{A}, \mathbf{B}) = (\mathbf{A} \cdot \mathbf{B}) / \|\mathbf{A}\| \cdot \|\mathbf{B}\|$$



Absolute length of
vector A

Cosine similarity

$$s(\mathbf{A}, \mathbf{B}) = \cos(\mathbf{A}, \mathbf{B}) = (\mathbf{A} \cdot \mathbf{B}) / \|\mathbf{A}\| \cdot \|\mathbf{B}\|$$

$$\mathbf{A} = (1, 1)$$

$$\mathbf{B} = (2, 2)$$

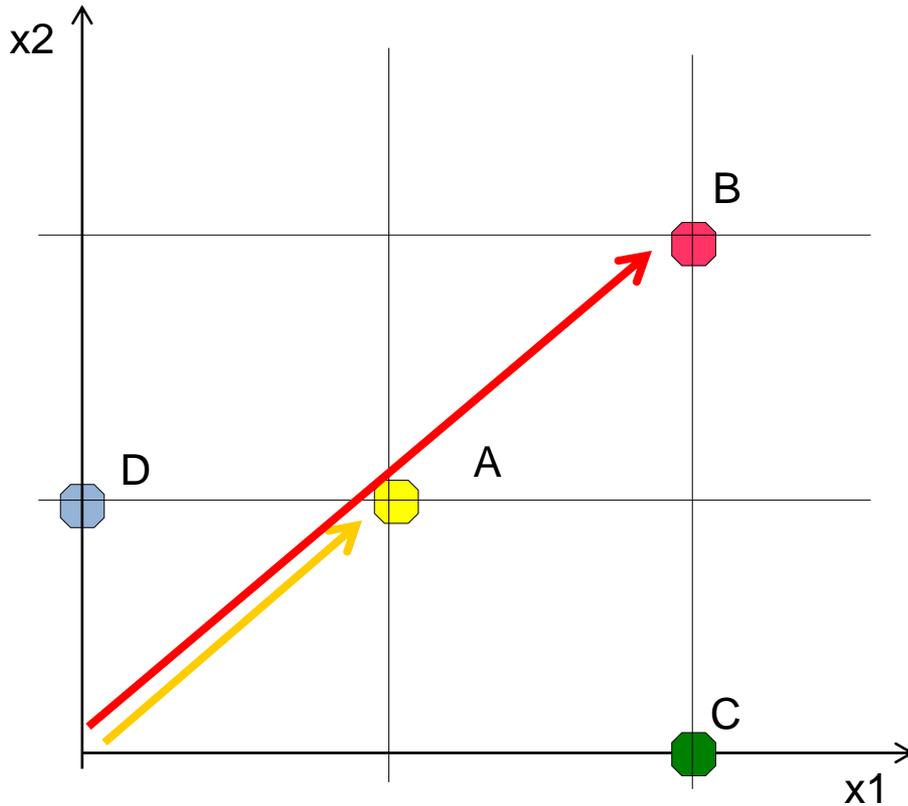
$$\mathbf{A} \cdot \mathbf{B} = 1 * 2 + 1 * 2 = 4$$

$$\|\mathbf{A}\| = \sqrt{1+1}$$

$$\|\mathbf{B}\| = \sqrt{4+4}$$

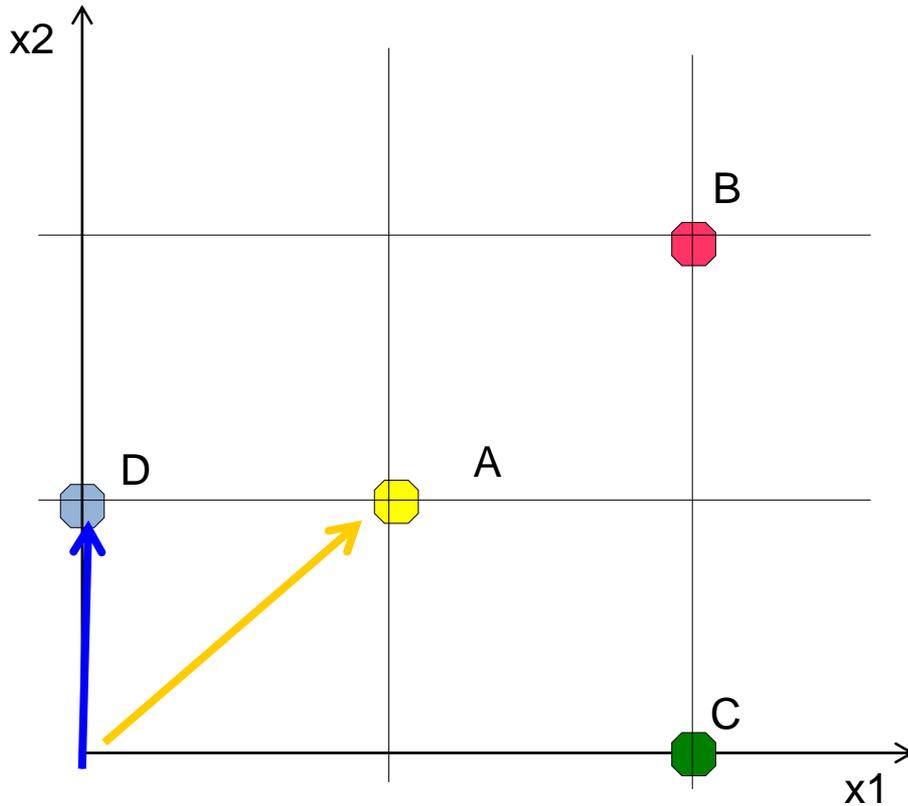
$$\|\mathbf{A}\| \cdot \|\mathbf{B}\| = \sqrt{16}$$

$$s(\mathbf{A}, \mathbf{B}) = \cos(\mathbf{A}, \mathbf{B}) = 1$$



Cosine similarity

$$s(\mathbf{A}, \mathbf{D}) = \cos(\mathbf{A}, \mathbf{D}) = (\mathbf{A} \bullet \mathbf{D}) / \|\mathbf{A}\| \cdot \|\mathbf{D}\|$$



$$\mathbf{A} = (1, 1)$$

$$\mathbf{D} = (0, 1)$$

$$\mathbf{A} \bullet \mathbf{D} = 0 + 1 = 1$$

$$\|\mathbf{A}\| = \sqrt{2}$$

$$\|\mathbf{D}\| = 1$$

$$\|\mathbf{A}\| \cdot \|\mathbf{D}\| = \sqrt{2}$$

$$s(\mathbf{A}, \mathbf{D}) = \cos(\mathbf{A}, \mathbf{D})$$

$$= \sqrt{1/2} \approx 0.7$$

Cosine similarity

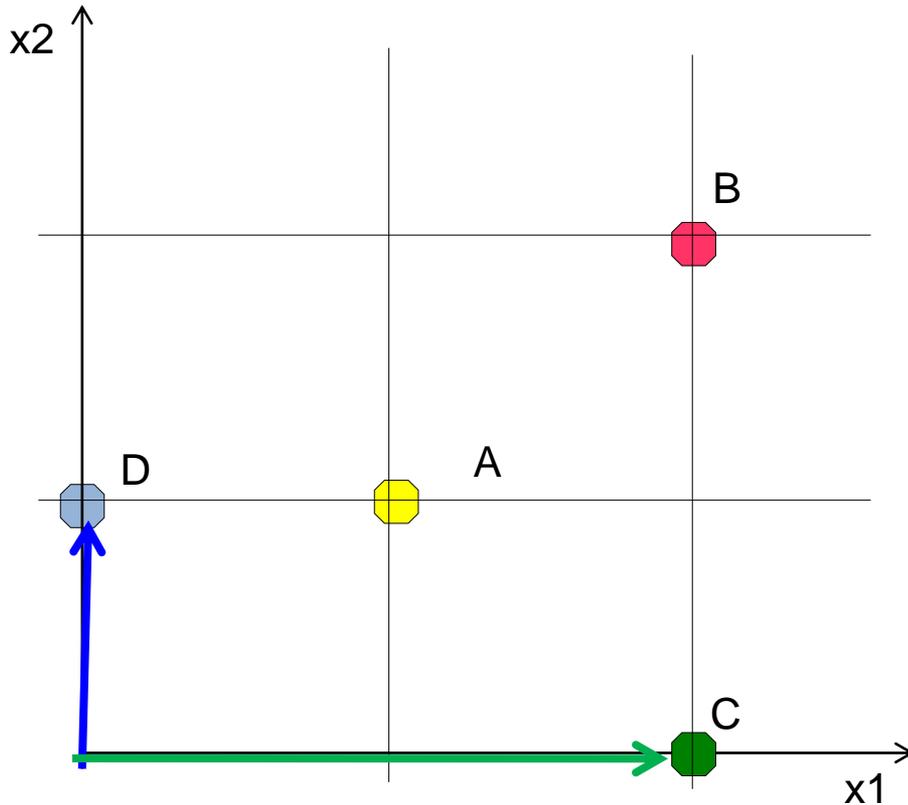
$$s(\mathbf{C}, \mathbf{D}) = \cos(\mathbf{C}, \mathbf{D}) = (\mathbf{C} \bullet \mathbf{D}) / \|\mathbf{C}\| \cdot \|\mathbf{D}\|$$

$$\mathbf{C} = (2, 0)$$

$$\mathbf{D} = (0, 1)$$

$$\mathbf{C} \bullet \mathbf{D} = 0$$

$$s(\mathbf{C}, \mathbf{D}) = \cos(\mathbf{C}, \mathbf{D}) = 0$$



Cosine Similarity for document vectors

	w1	w2	w3	w4	w5	w6
$x=($	1	0	0	0	0	0)
$y=($	0	0	0	1	2	0)
$z=($	0	0	0	4	8	0)

Cosine between \mathbf{x} and \mathbf{y} is 0 (dot-product is 0). These documents are not similar.

Cosine between \mathbf{y} and \mathbf{z} is 1: though the number of times each word occurs in y and z is different, these documents are about the same topic

Pearson correlation

- A ***correlation*** is a number between -1 and +1 that measures the degree of association between two variables (in our case – between 2 data objects for which we recorded n observations)
- A **positive** value for the correlation implies a **positive association** (the values across all observations vary in the same direction)
- A **negative** value for the correlation implies a negative or **inverse association** – which makes 2 data objects dissimilar
- A values close to **0** implies that there is **no correlation** between two data objects

Pearson correlation formula

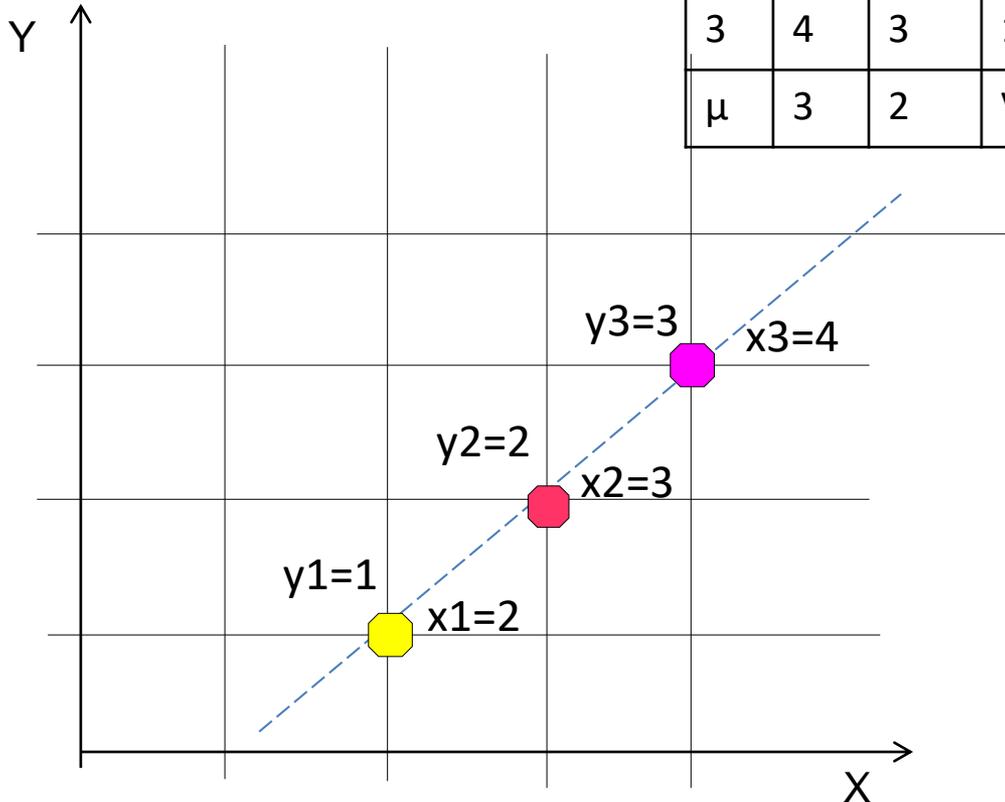
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- In numerator we see a *covariance* - a measure of the joint variability of 2 data objects across n observations
- We normalize it by dividing by a *variance* inside each separate data object

Pearson correlation example 1

Are all (3) observations for objects X and Y change in the same direction?

	x	y	$x_i - \mu_x$	$y_i - \mu_y$	Cov(x,y)
1	2	1	-1	-1	1
2	3	2	0	0	0
3	4	3	1	1	1
μ	3	2	Var(x) = $\sqrt{2}$	Var(y) = $\sqrt{2}$	2



$$r = 2/(\sqrt{2} * \sqrt{2}) = 1$$

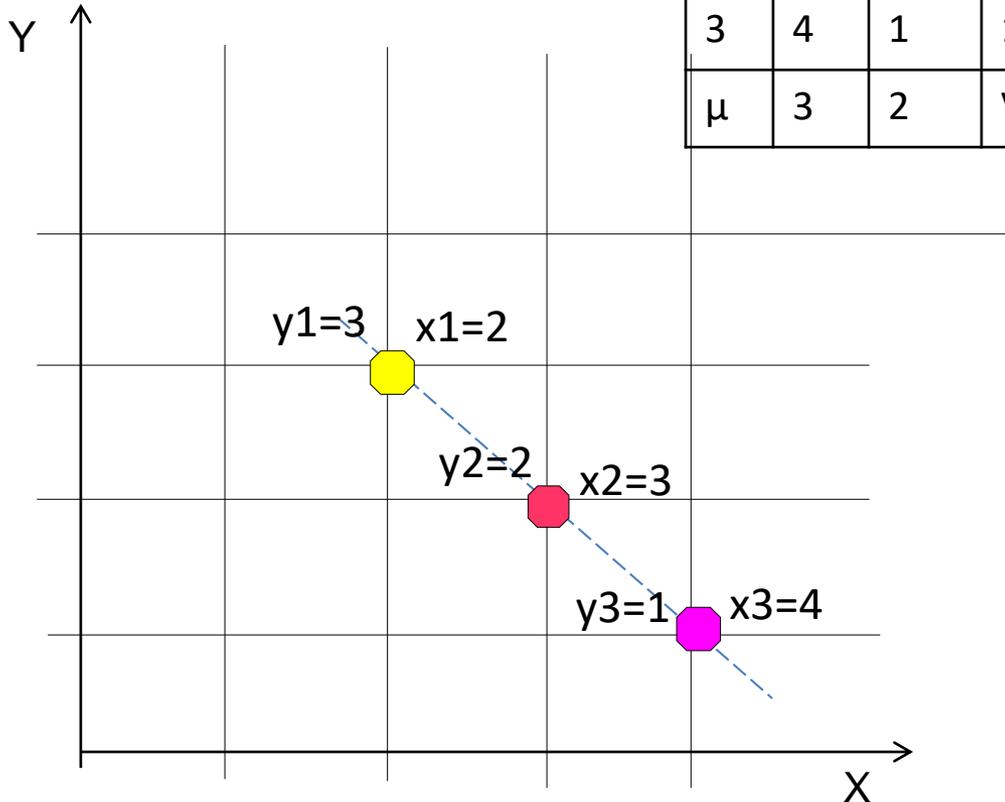
Full positive correlation
(X and Y are very similar)

Note that absolute values across each dimension do not play any role – only direction is important

Pearson correlation example 2

Are all (3) observations for objects X and Y change in the same direction?

	x	y	$x_i - \mu_x$	$y_i - \mu_y$	Cov(x,y)
1	2	3	-1	1	-1
2	3	2	0	0	0
3	4	1	1	-1	-1
μ	3	2	$\text{Var}(x) = \sqrt{2}$	$\text{Var}(y) = \sqrt{2}$	-2



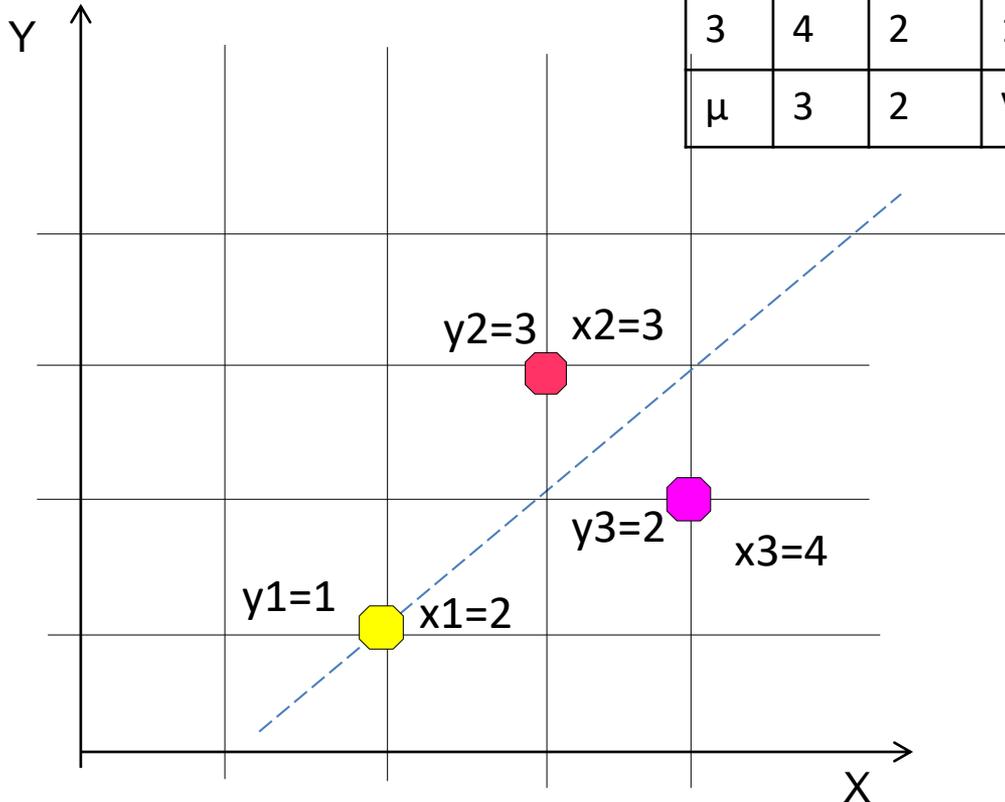
$$r = -2/(\sqrt{2} * \sqrt{2}) = -1$$

Full negative correlation
(X and Y are least similar – quite opposite)

Pearson correlation example 3

Are all (3) observations for objects X and Y change in the same direction?

	x	y	$x_i - \mu_x$	$y_i - \mu_y$	Cov(x,y)
1	2	1	-1	-1	1
2	3	3	0	1	0
3	4	2	1	0	0
μ	3	2	$\text{Var}(x) = \sqrt{2}$	$\text{Var}(y) = \sqrt{2}$	1



$$r = 1/(\sqrt{2} * \sqrt{2}) = 1/2$$

Partly correlated objects (less similar)

Relationship between Pearson correlation and cosine similarity

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}},$$

- By subtracting mean from each value, we effectively just transposing vectors to center them around mean
- Pearson correlation is nothing else as a cosine between two vectors after they are centered around the mean for each dimension
- Pearson correlation is a cosine of centered vectors

Combining Similarities

- Sometimes attributes are of many different types, but an overall similarity/dissimilarity is needed.
- For each attribute k , compute a similarity s_k
- Then average,

$$\text{similarity}(p, q) = \frac{\sum_{k=1}^n s_k}{n}$$

- Similar formula for dissimilarity (distance)

Scaling attributes for consistency

- X- in yards, Y in cm
- X- number of children, Y – income

Difference in 1 dollar = difference in 1 child?

Scaling: map all variables to a common range 0-1

Scaling for consistency

$$a_i = \frac{v_i - \min(\text{all } v)}{\max(\text{all } v) - \min(\text{all } v)}$$

For numeric attributes: convert all values into the range between 0 and 1

Scaling vectors

- Vector normalization – changes the vector values so that the length of the vector is 1, only the direction is compared
- $X = \{\text{Debt} = 200,000 \text{ equity} = 100,000\}$
- $Y = \{\text{Debt} = 2,000 \text{ equity} = 1,000\}$

Emphasizes internal **relation** between different attributes of each record

Encode expert knowledge with weights

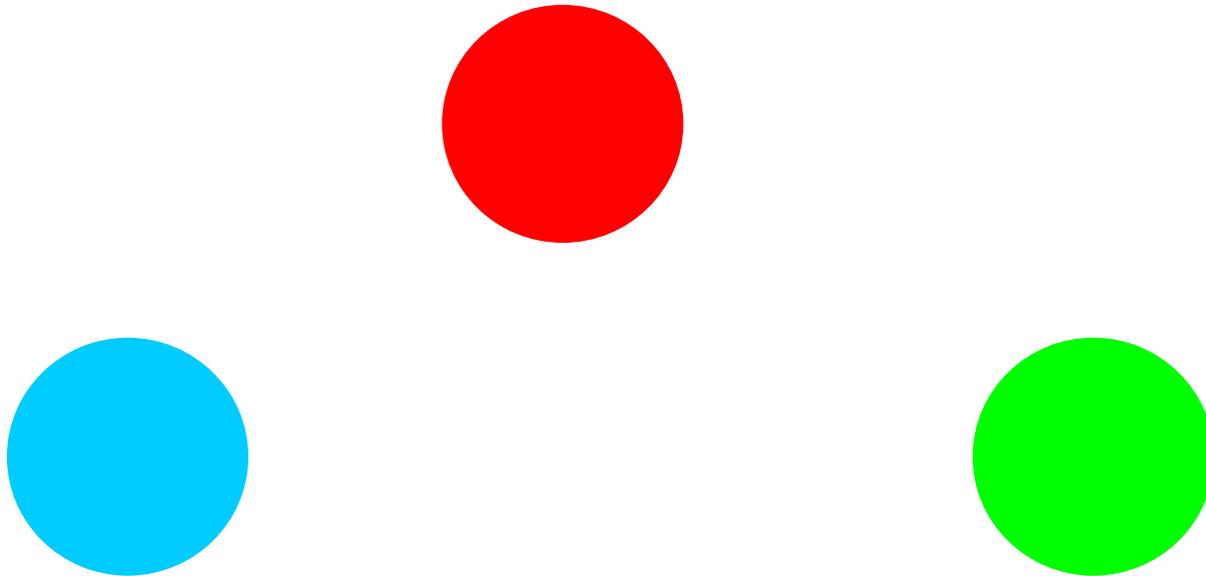
- Changes in one variable should not be more significant only because of differences in magnitudes of values
- After scaling to **get rid of bias due to the units**, use weights to **introduce bias** based on expert knowledge of context:
 - 2 families with the same income and number of children are more similar than 2 families living in the same neighborhood
 - Number of children is more important than the number of credit cards

Part 2

HOW DO WE DEFINE A CLUSTER

Types of Clusters 1/4: Well-Separated

- Any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.



3 well-separated clusters

Types of Clusters 2/4: Center-Based

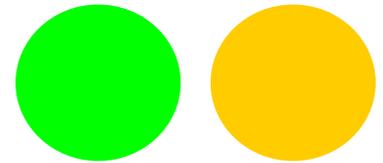
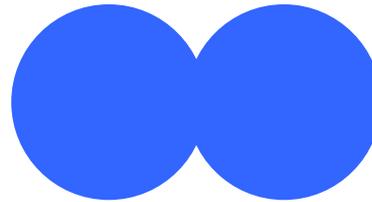
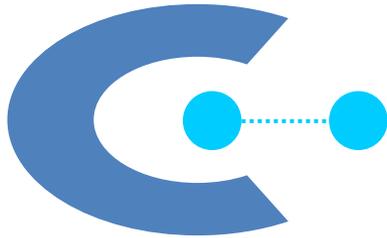
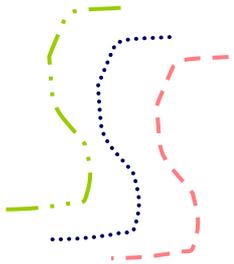
- An object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often called a **centroid**, the average of all the points in the cluster, or a **medoid**, the most “representative” point of a cluster



4 center-based clusters

Types of Clusters 3/4: Contiguity-Based

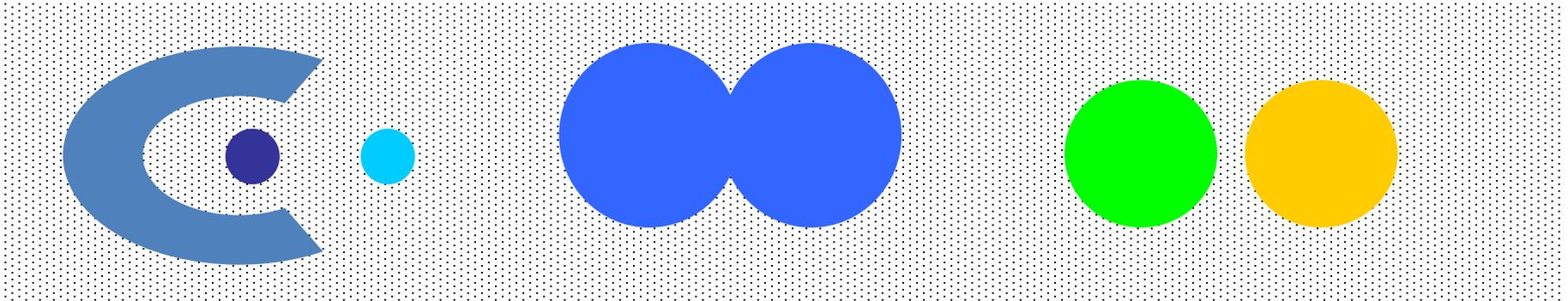
- Contiguous Cluster (Nearest neighbor or Transitive)
- A point in a cluster is closer to at least one point in the cluster than to any point not in the cluster. The group of objects that are connected to one another.



8 contiguous clusters

Types of Clusters 4/4: Density-Based

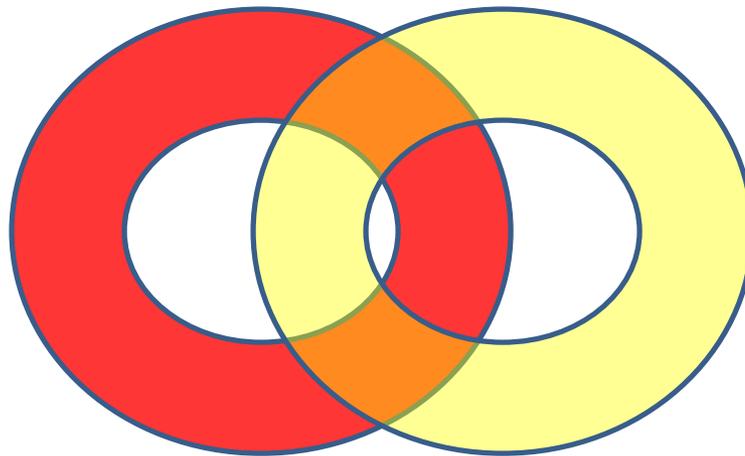
- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.



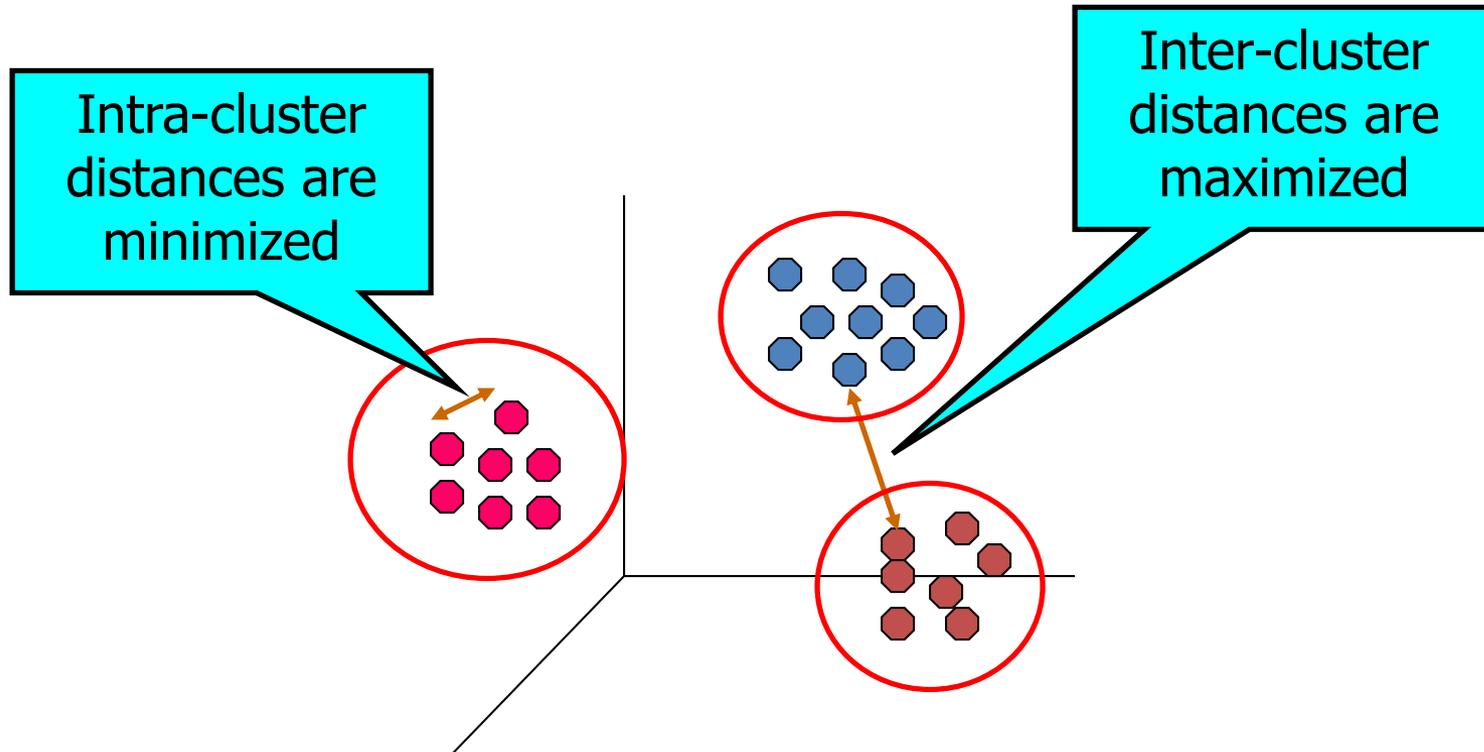
6 density-based clusters

General definition: conceptual clusters

- A set of objects that **share some property**. This includes all previous cluster types
- In addition it includes clusters defined by a *concept*. Such clusters are used in pattern recognition. To discover such clusters automatically, the concept should be defined first.



Clustering algorithm: goal



Clustering algorithms

- ▶ • *K*-means clustering
- Agglomerative hierarchical clustering
- Density-based clustering

Iterative solution: K-means clustering algorithm

Select K random **seeds**

Do

Assign each record to the closest **seed**

Calculate **centroid** of each cluster

(take average value for each dimension
of all records in the cluster)

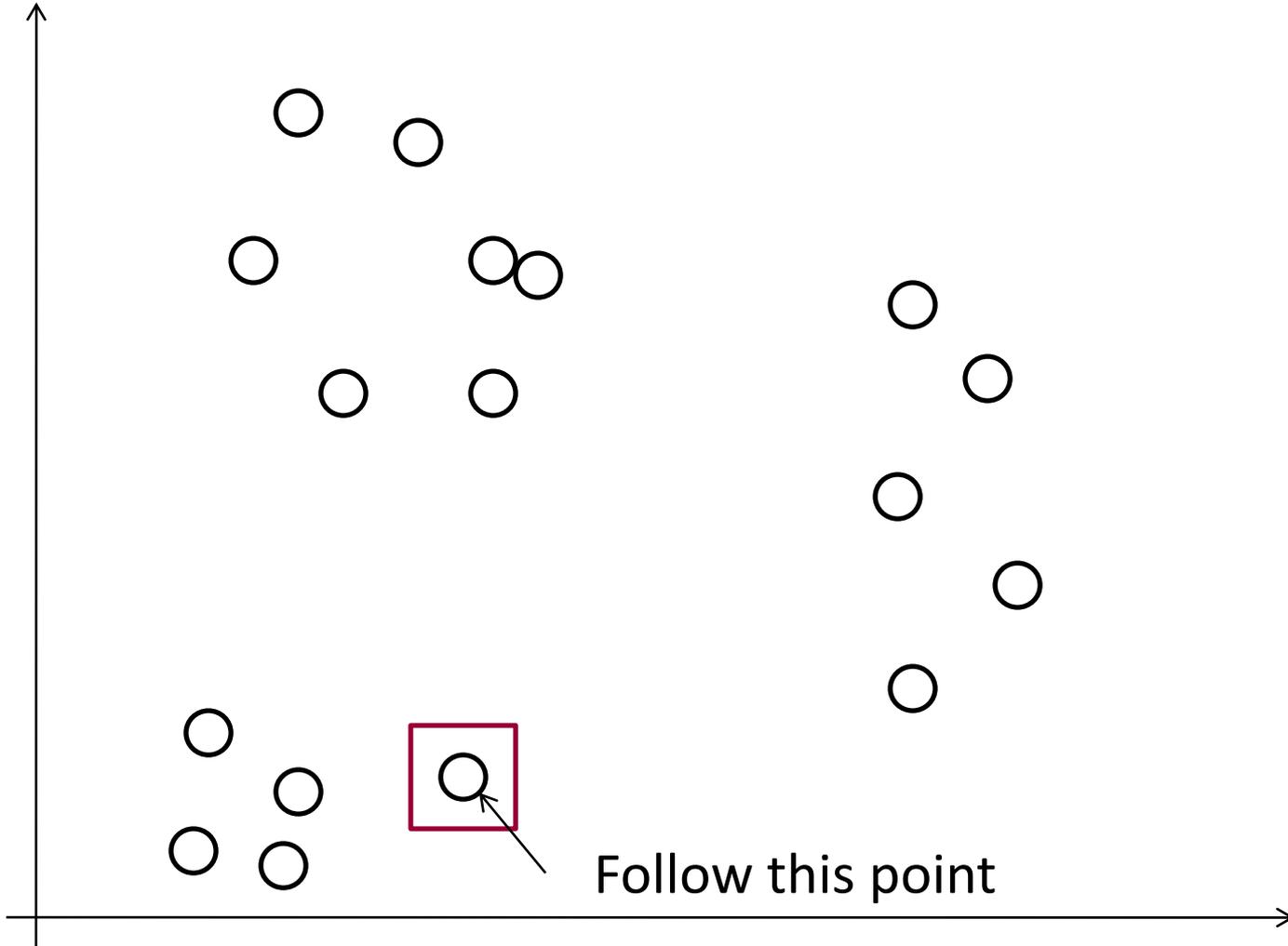
Set these **centroids** as new **seeds**

Until coordinates of **seeds** *do not change*

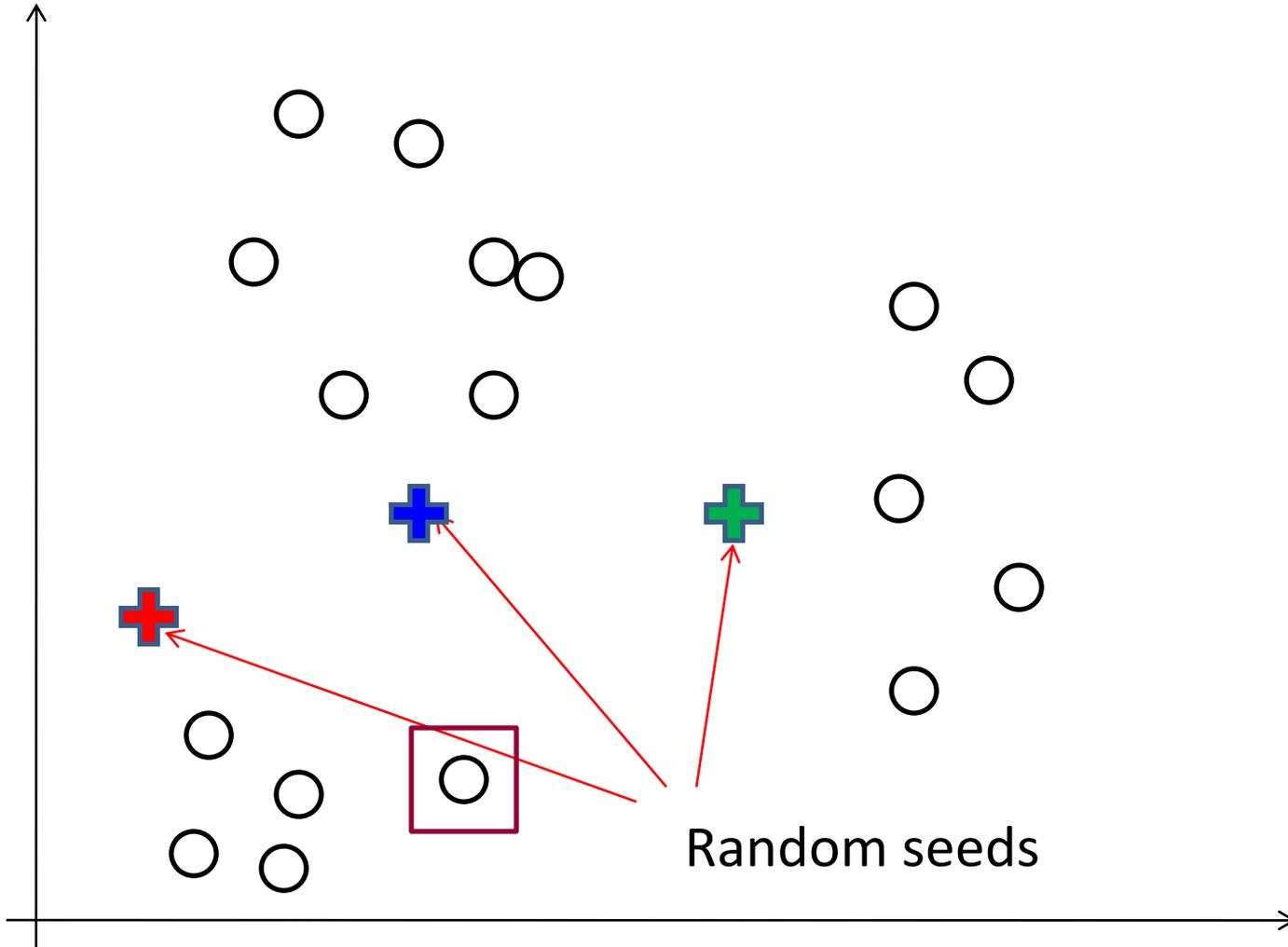
This algorithm in each iteration makes assignment of points such that intra-cluster distances are decreasing.

Local optimization technique – moves into the direction of local minimum, might miss the best solution

Example 1: $K=3$

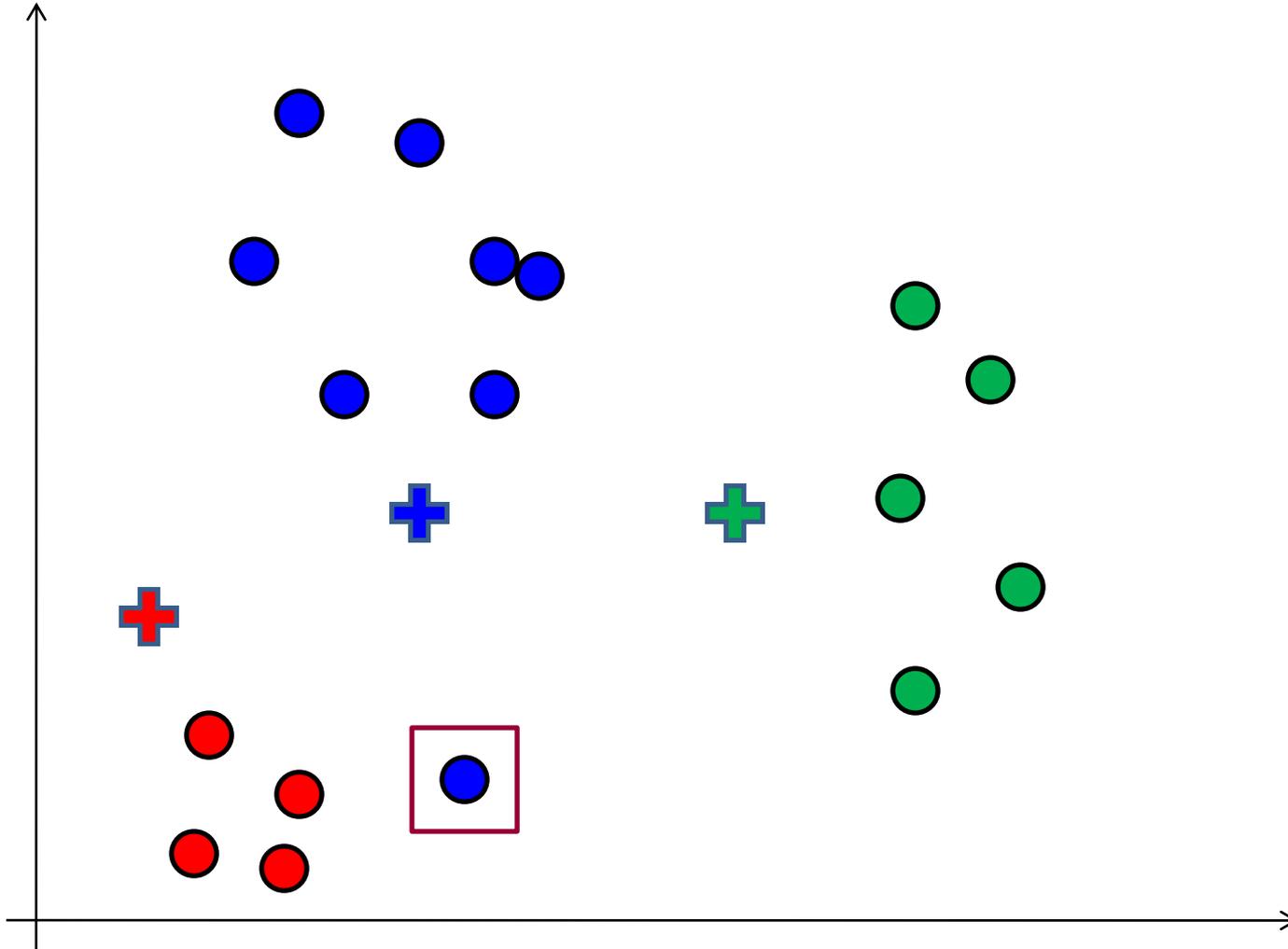


Example 1: initialization



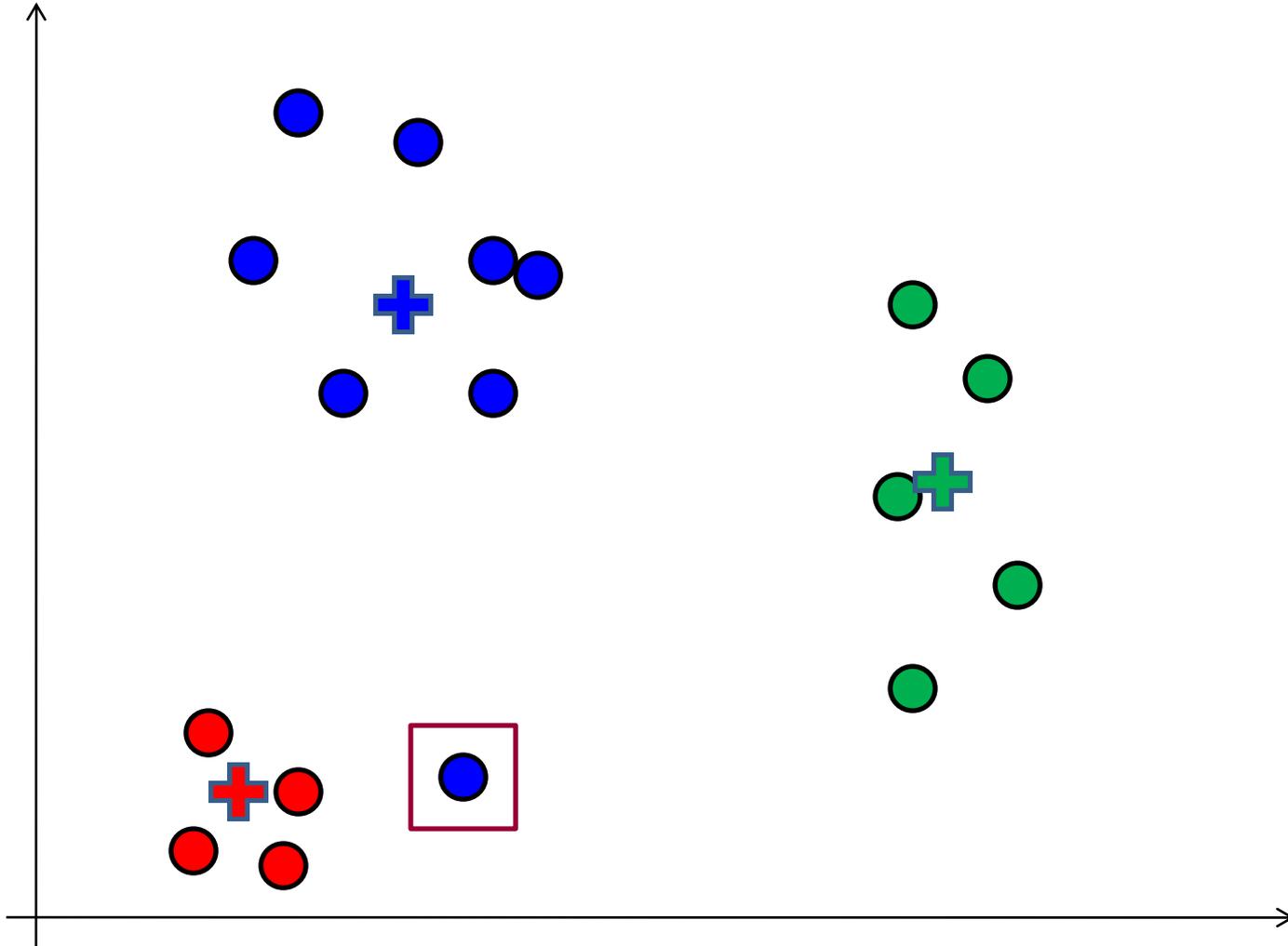
Example 1: iteration 1.

Assign each point to the closest seed



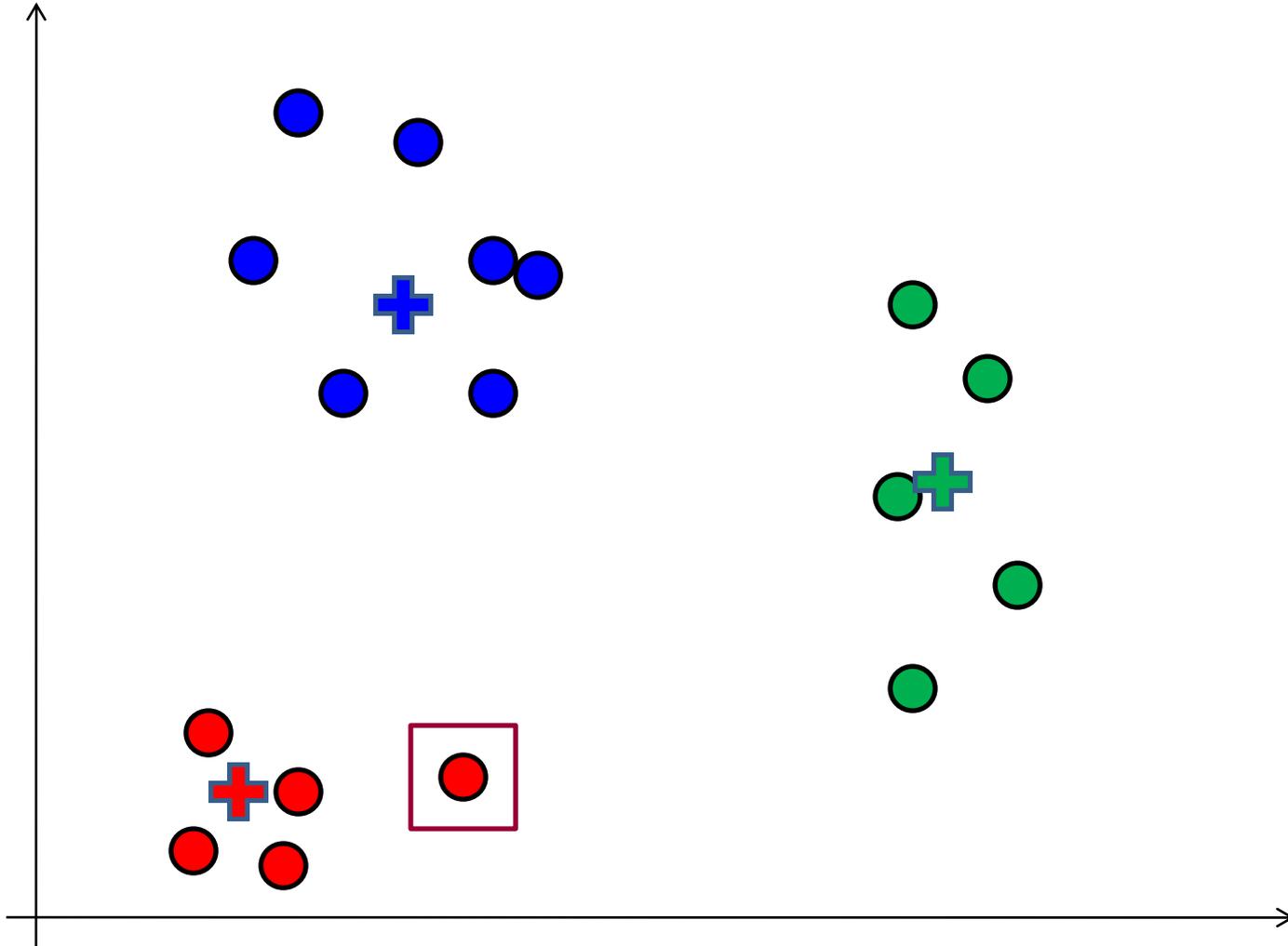
Example 1: iteration 1.

Recalculate centroids



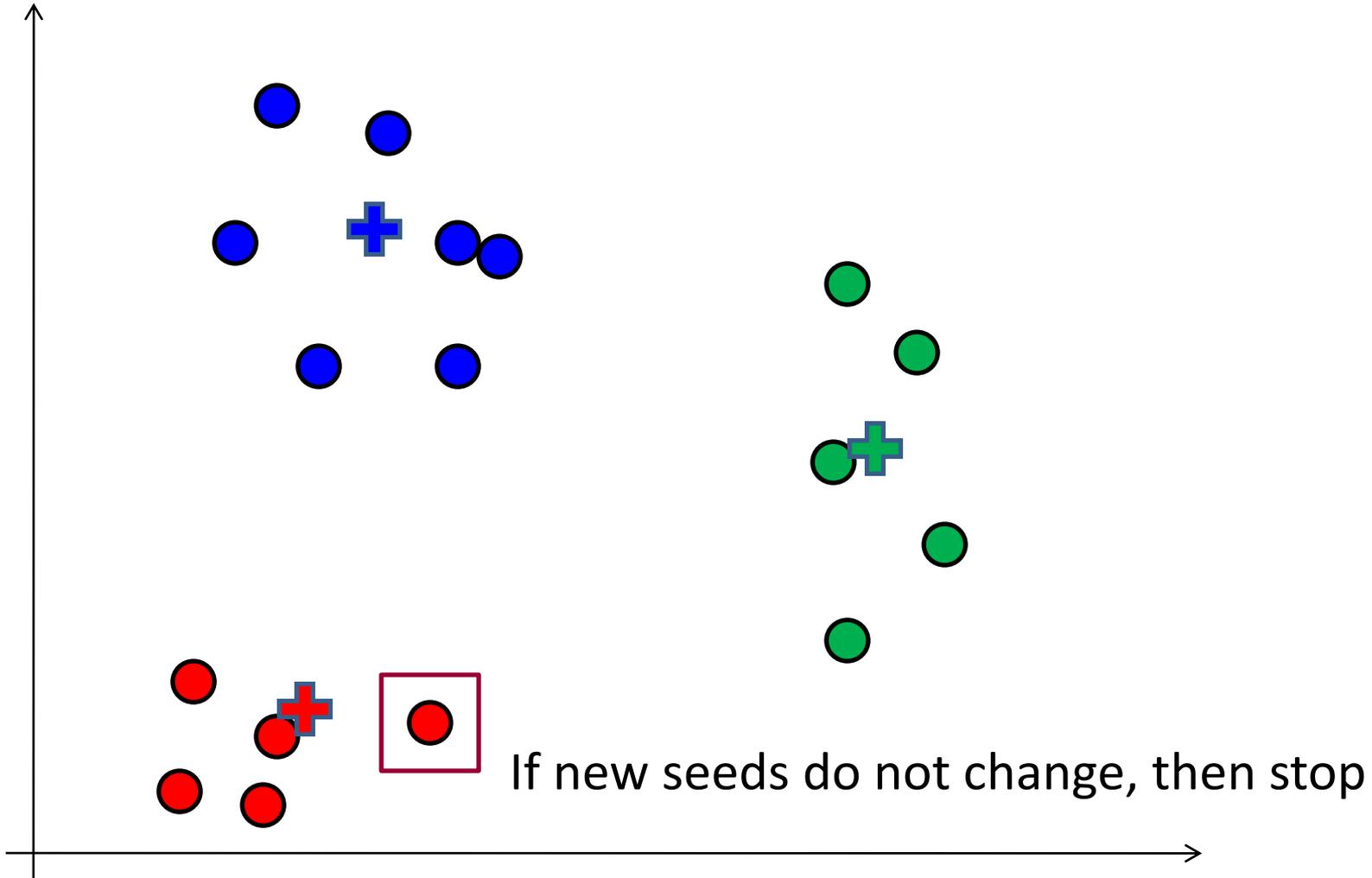
Example 1: iteration 2.

Assign each point to the closest seed

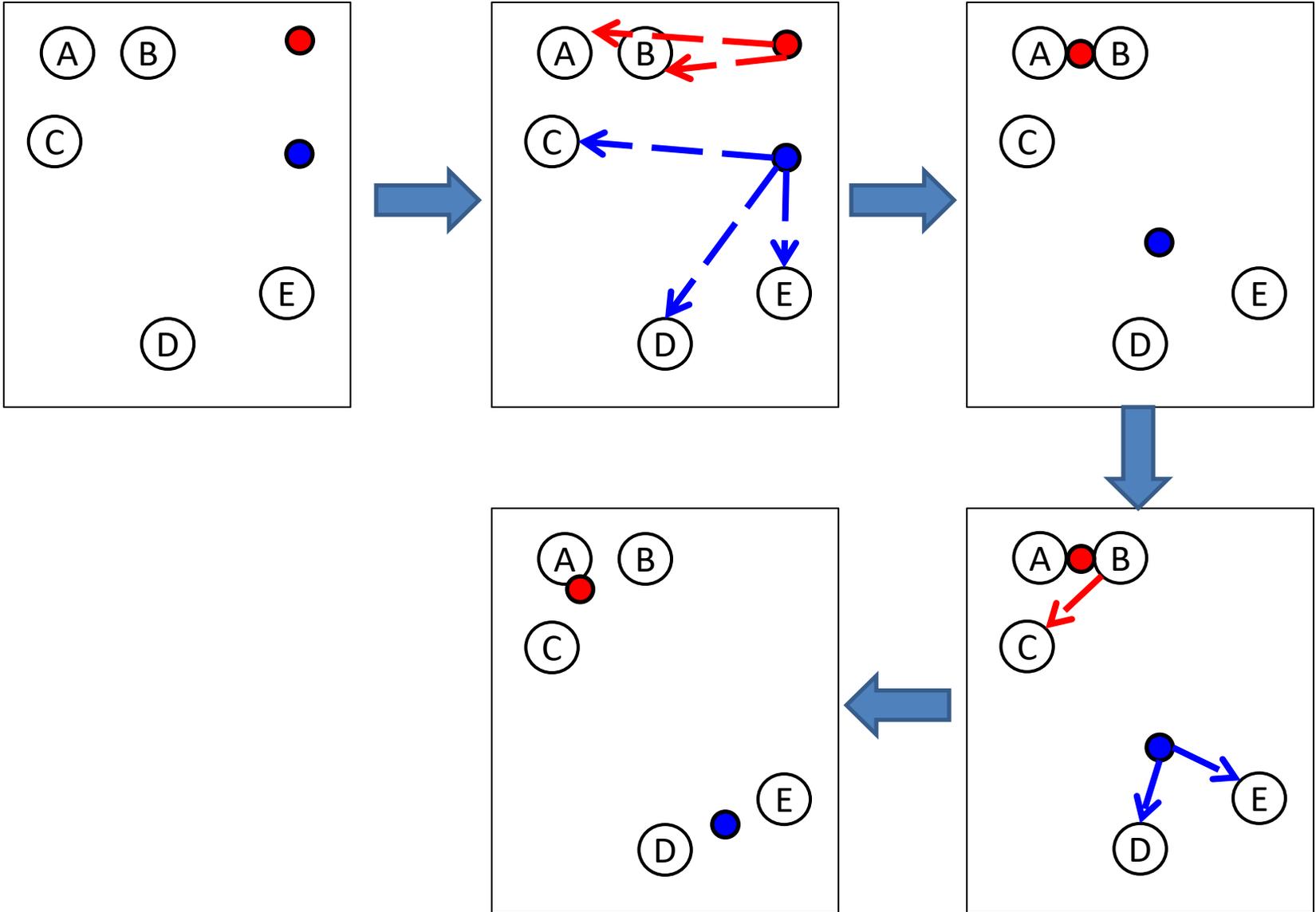


Example 1: iteration 2.

recalculate centroids – new seeds



Example 2: K=2



Evaluating K-means Clusters

- Most common measure is **Sum of Squared Error (SSE)**
 - For each point, the error is the distance to the nearest cluster centroid
 - To get **SSE**, we square these errors and sum them up.

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} [dist(m_i, x)]^2$$

x is a data point in cluster C_i and

m_i is the representative point for cluster C_i (in our case, centroid)

Centroid that minimizes SSE of each cluster is a mean

At each iteration, we decrease total SSE, but with respect to a given set of centroids and point assignments

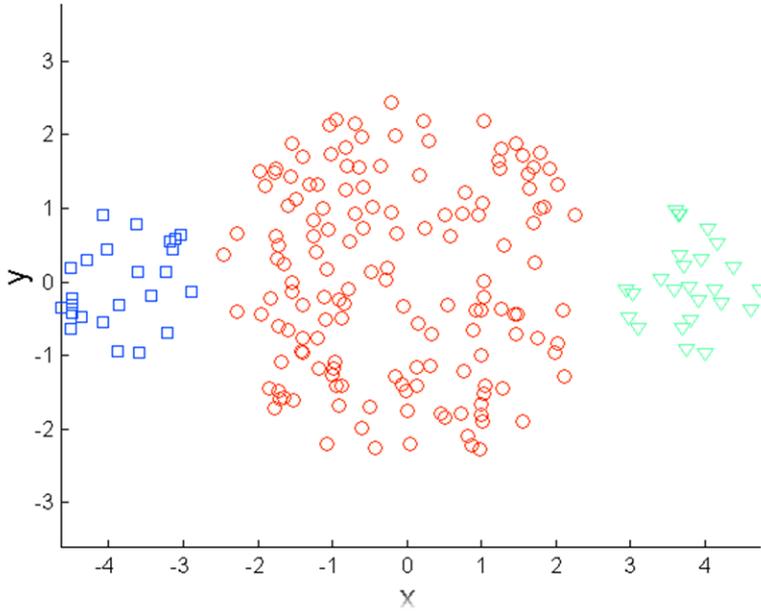
K-means Clustering – Details

- Initial centroids may be chosen randomly.
 - Clusters produced vary from one run to another.
- Most of the convergence happens in the first few iterations.
 - Often the stopping condition is changed to ‘Until relatively few points change clusters’
- Complexity is $O(l * K * n * d)$
 - n = number of points, K = number of clusters,
 l = number of iterations, d = number of attributes

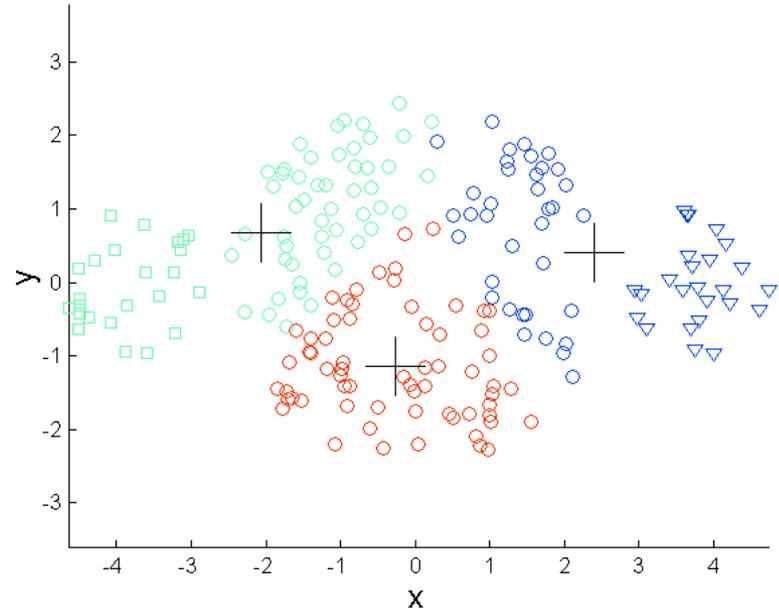
Limitations of K-means

- **K-means** has problems when clusters are of
 - Differing **Sizes**
 - Differing **Densities**
 - **Non-globular shapes**

Limitations of K-means: Differing Sizes

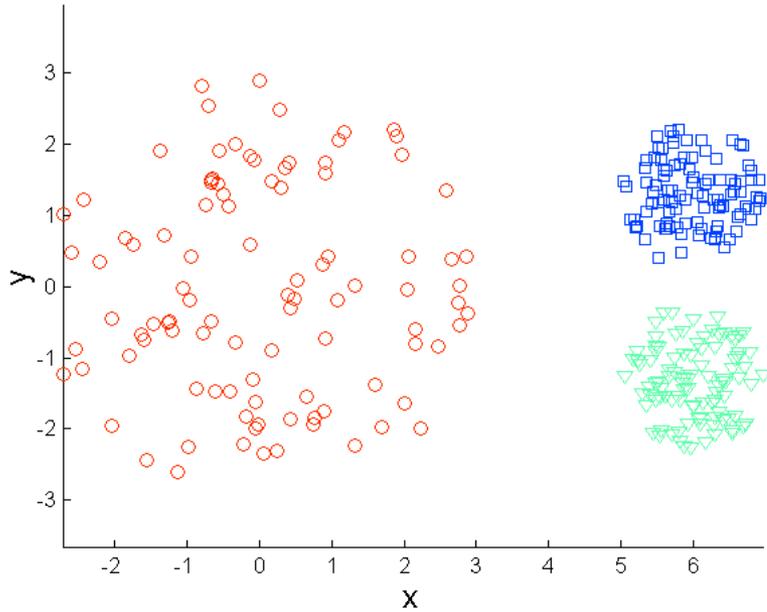


Original Points

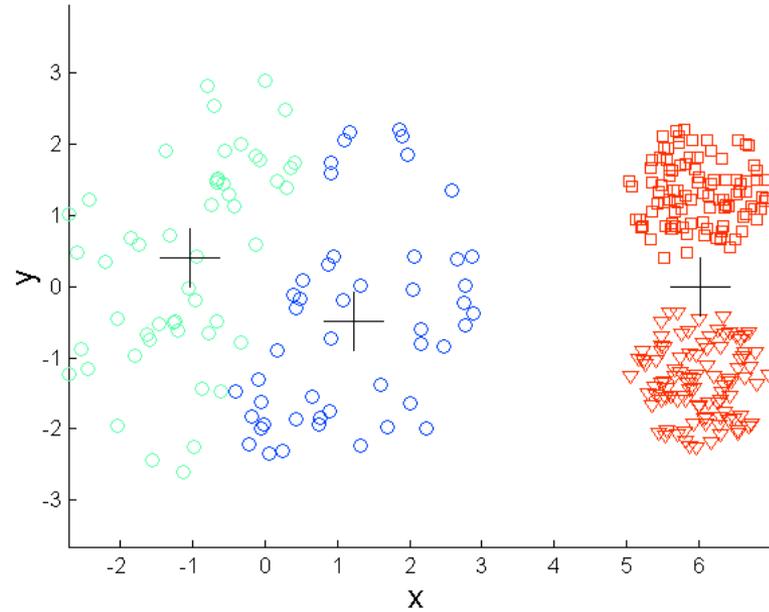


K-means (3 Clusters)

Limitations of K-means: Differing Density

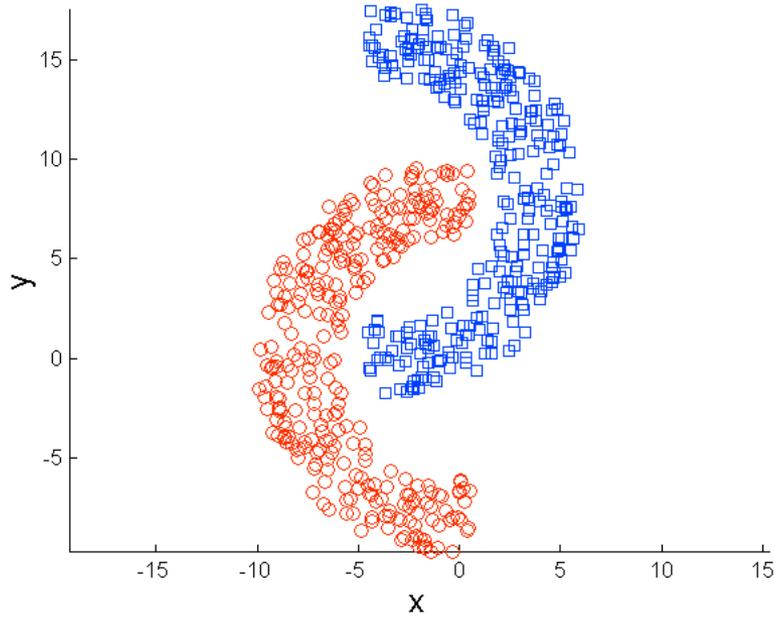


Original Points

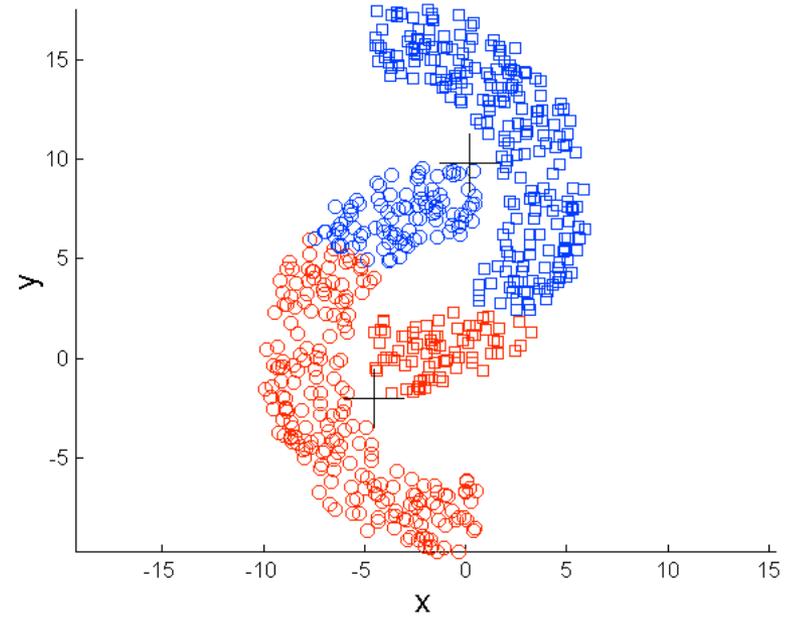


K-means (3 Clusters)

Limitations of K-means: Non-globular Shapes



Original Points

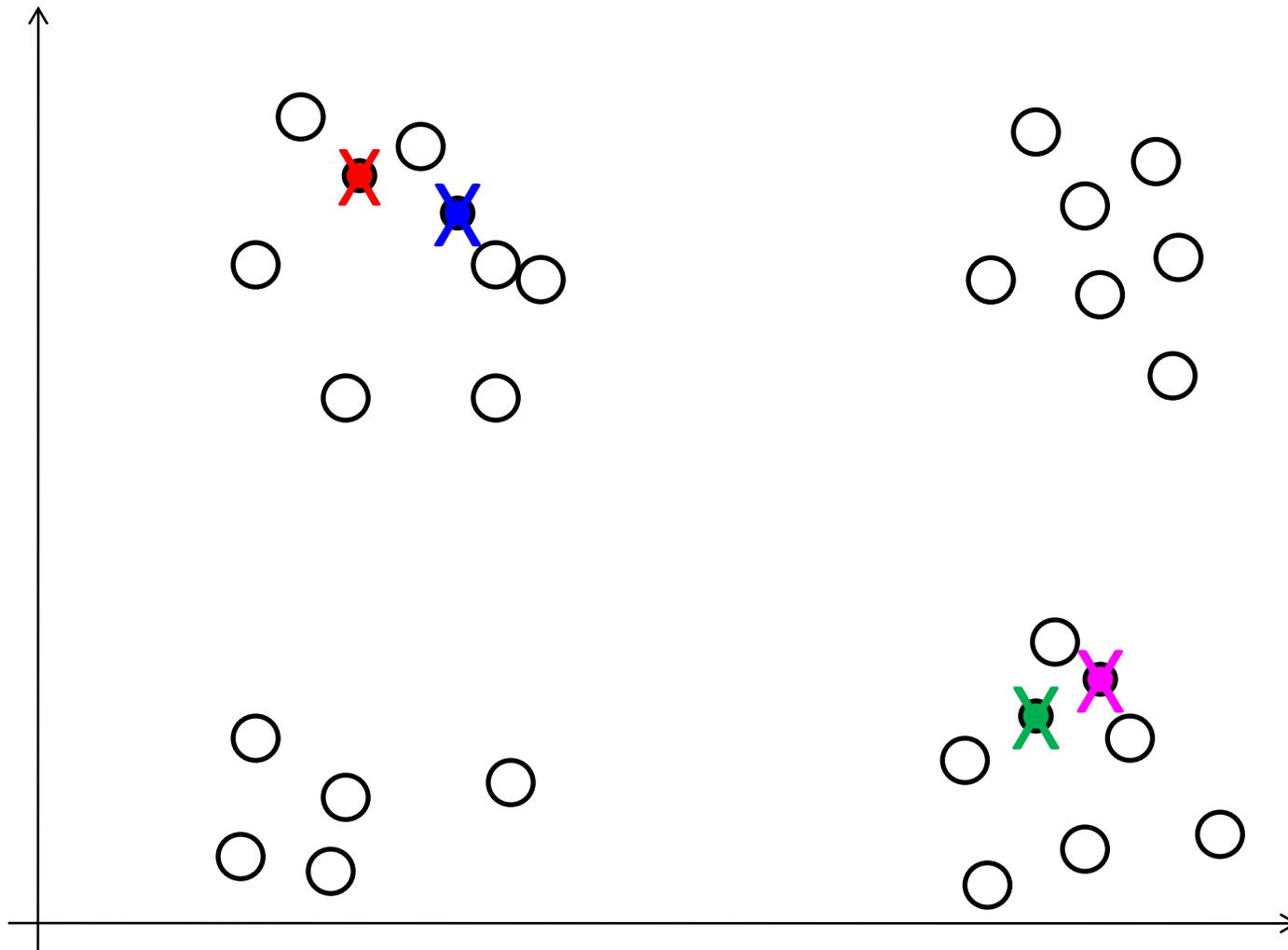


K-means (2 Clusters)

Limitations of K-means

- **K-means** has problems when clusters are of
 - Differing **Sizes**
 - Differing **Densities**
 - **Non-globular shapes**
- But even for globular clusters, the choice of initial centroids influences the quality of clustering

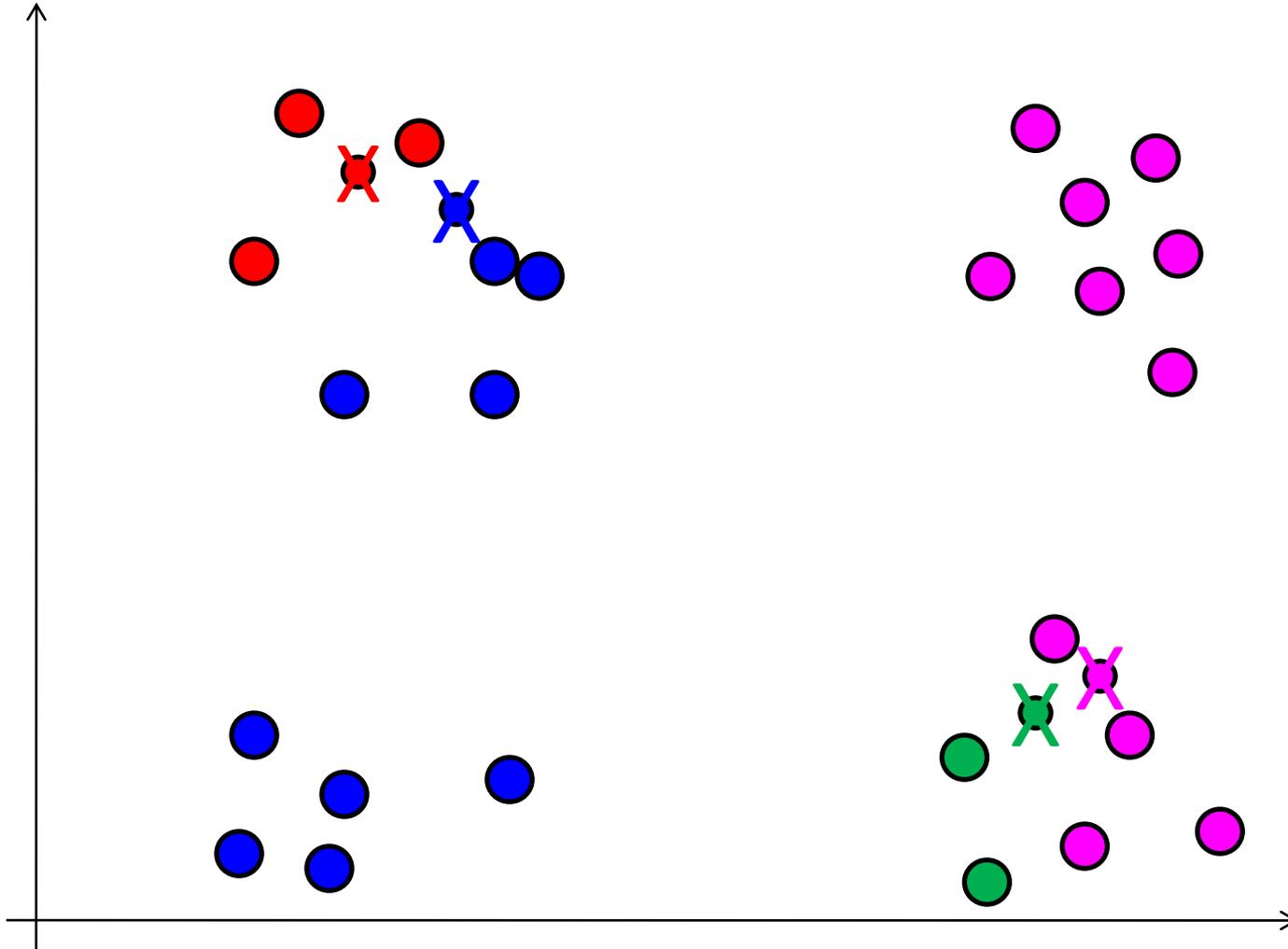
1. Importance of choosing initial centroids: $K=4$



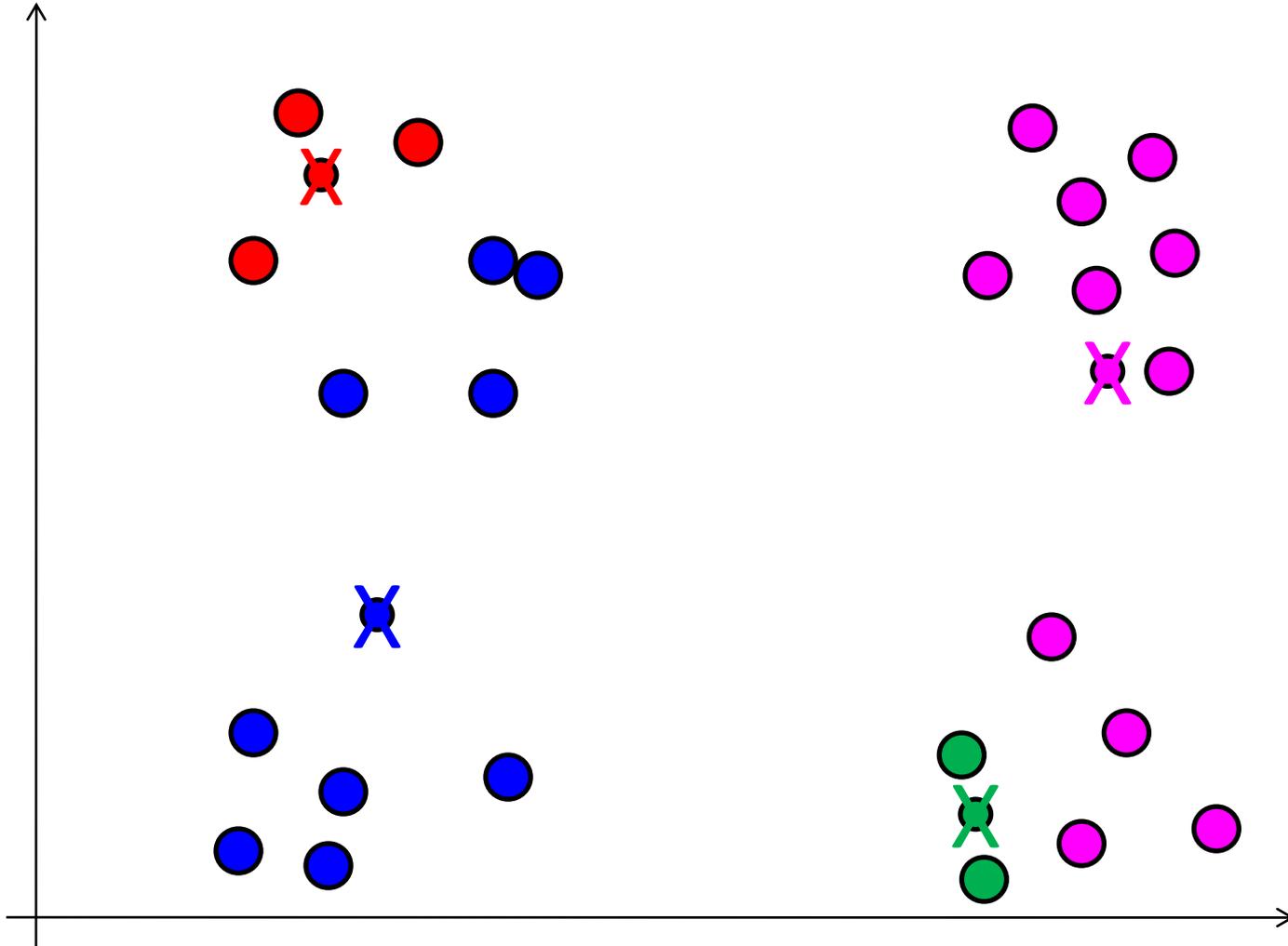
2 pairs of clusters

Initial seeds are chosen: 2 seeds per each pair

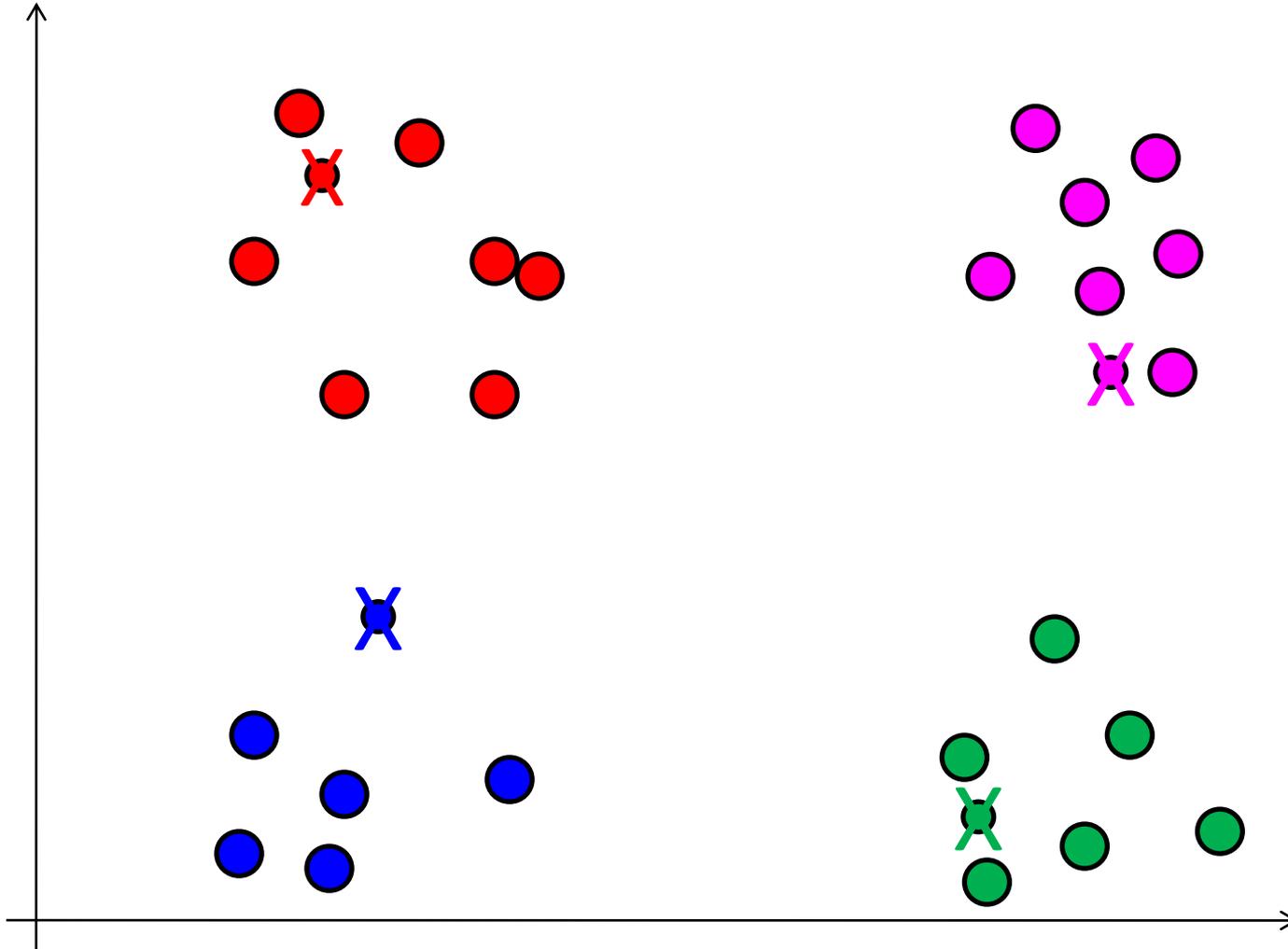
1. Importance of choosing initial centroids: point assignments



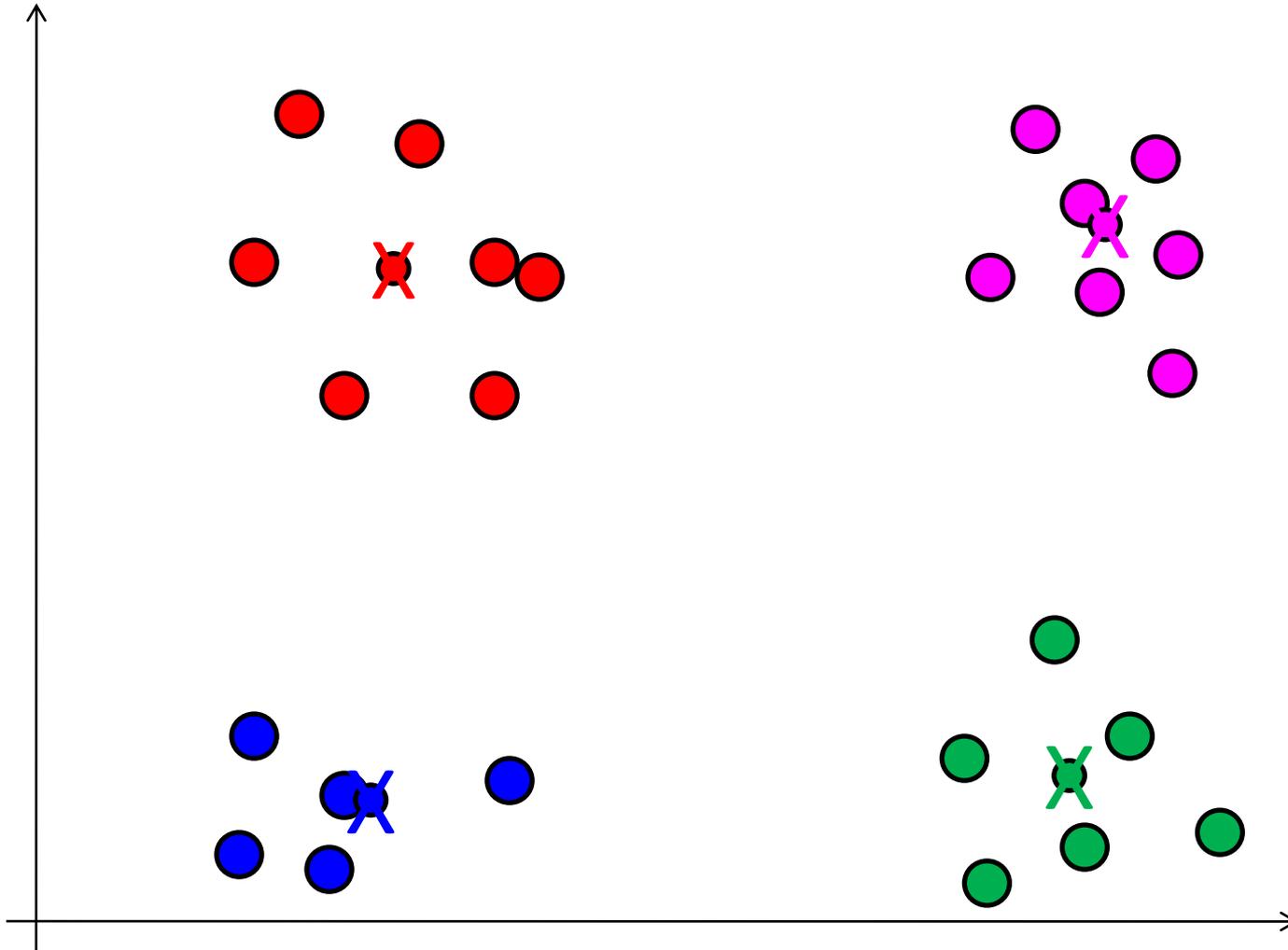
1. Importance of choosing initial centroids: recalculate centroids



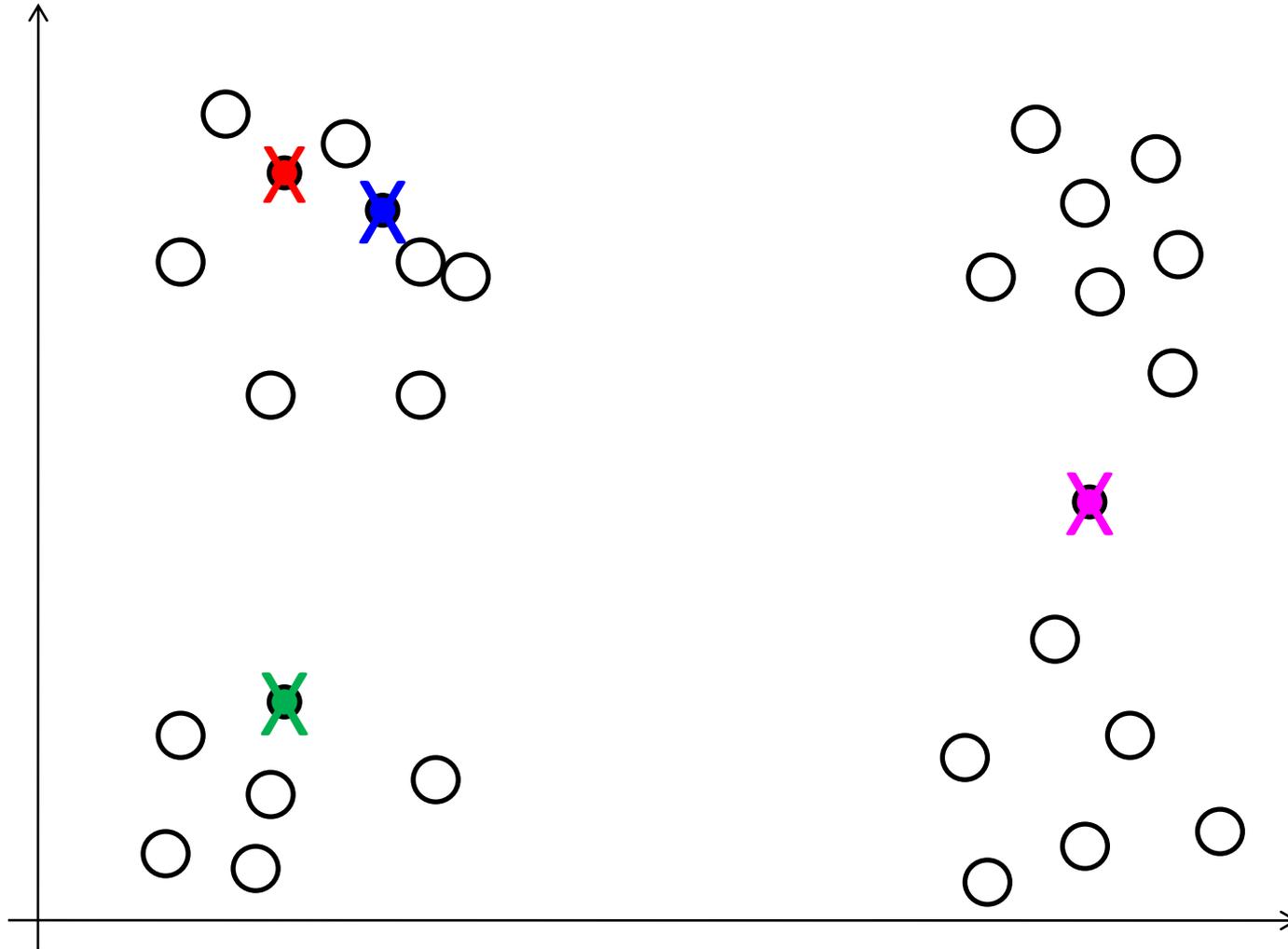
1. Importance of choosing initial centroids: points re-assignments



1. Importance of choosing initial centroids: success – correct clusters



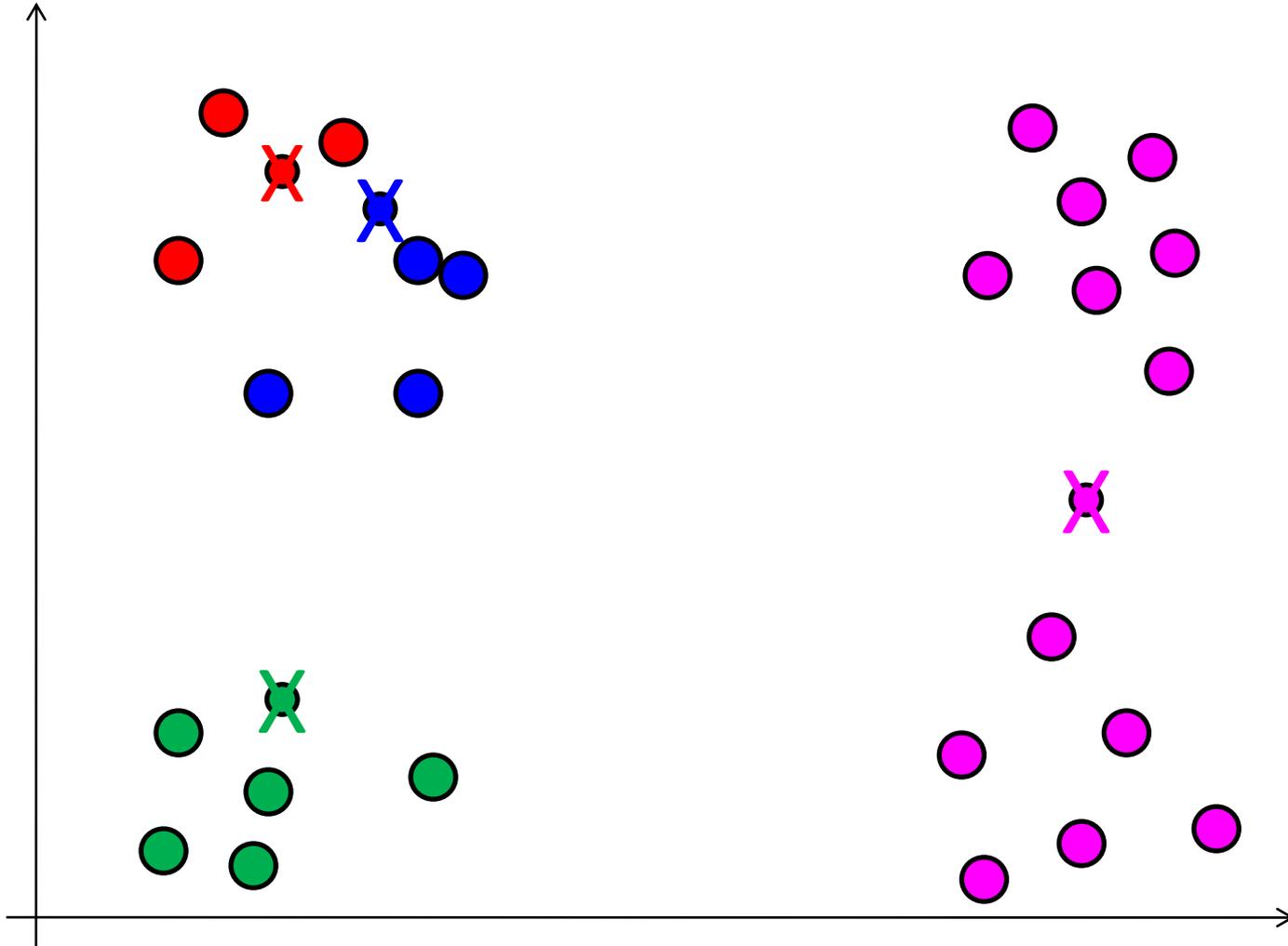
2. Importance of choosing initial centroids: $K=4$



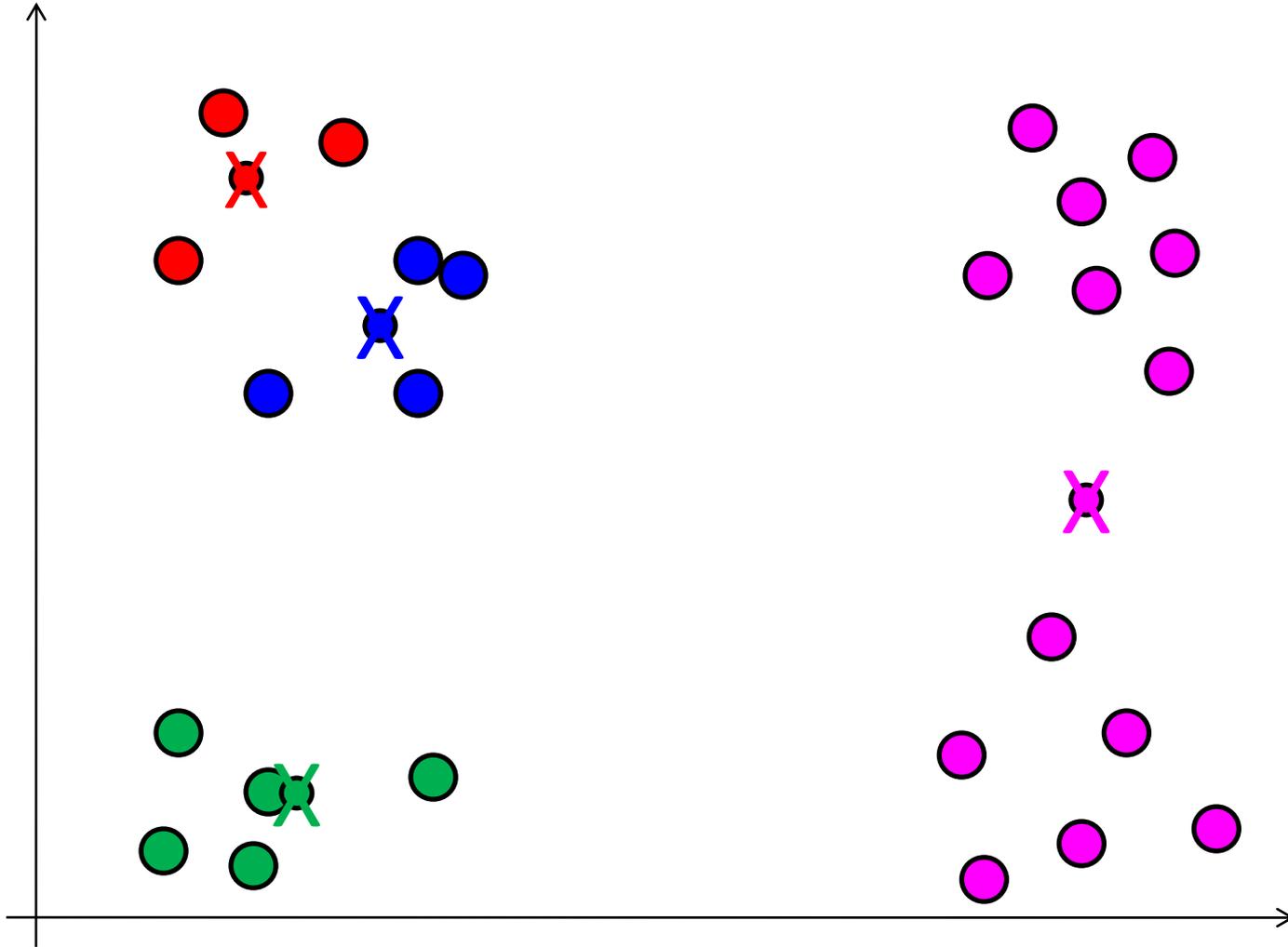
2 pairs of clusters

Initial seeds are chosen: 3 seeds in one pair

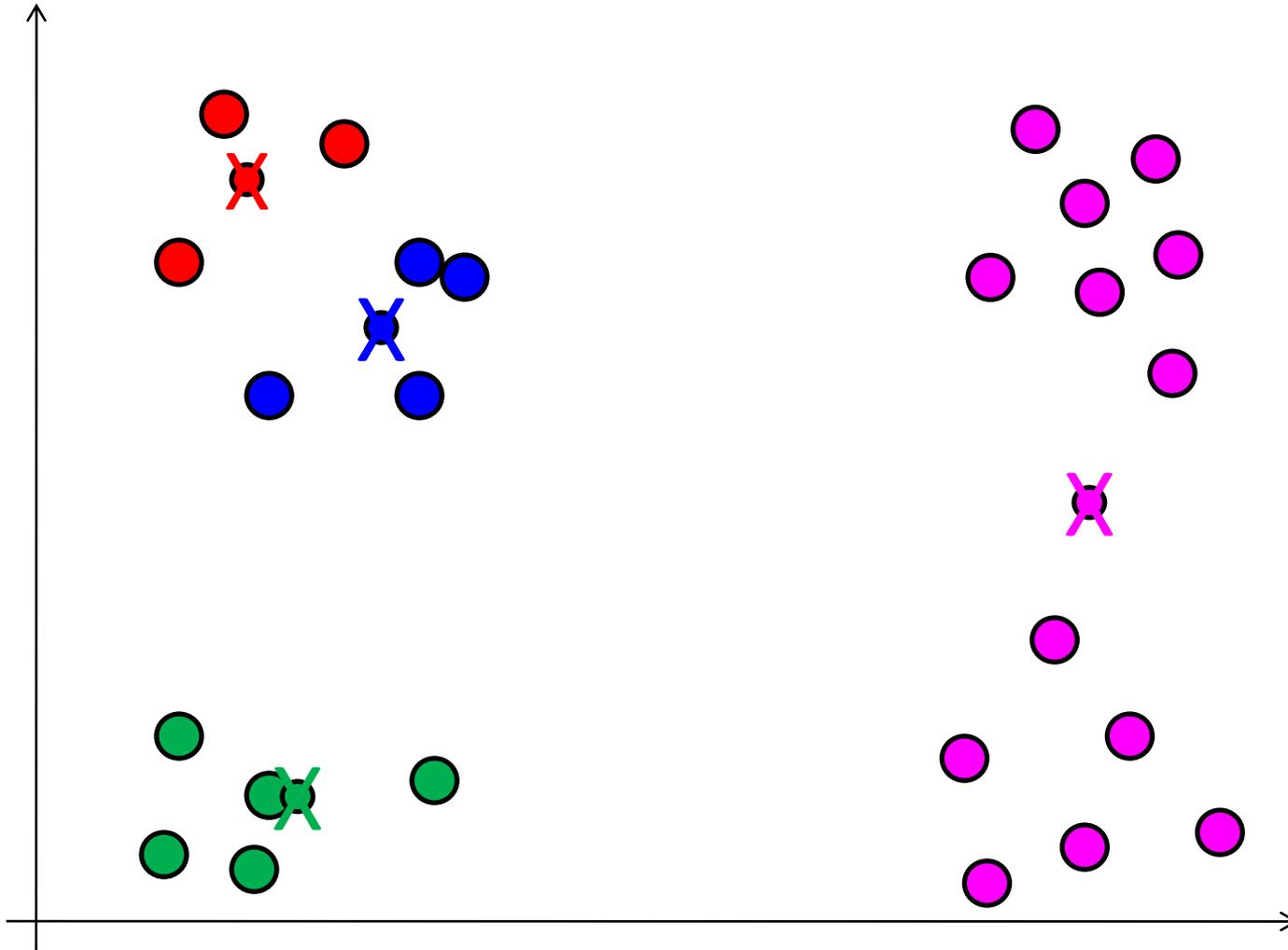
2. Importance of choosing initial centroids: assign points



2. Importance of choosing initial centroids: re-compute centroids



2. Importance of choosing initial centroids: found 4 clusters - **incorrect**



Problems with Selecting Initial Centroids

- Of course, the ideal would be to choose initial centroids, one from each true cluster.
- If there are K 'real' clusters then the chance of selecting one centroid from each cluster is small.
 - Chance is relatively small when K is large
 - If clusters are the same size, n , then

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{K^K}$$

- For example, if $K = 10$, then *probability* = $10!/10^{10} = 0.00036$
- Sometimes the initial centroids readjust themselves in the 'right' way, and sometimes they don't.

Solutions to Initial Centroids Problem

- Multiple runs
 - Helps, but probability is not on your side
- Bisecting K-means
 - Not as susceptible to initialization issues

Bisecting *K*means

- Straightforward extension of the basic *K*means algorithm.

Simple idea:

To obtain *K* clusters, split the set of points into two clusters, select one of these clusters to split, and so on, until *K* clusters have been produced.

Bisecting Kmeans

Initialize the list of clusters with the cluster consisting of all points.

Do

Remove a cluster with the highest SSE from the list of clusters.

//Perform several “trial” bisections of the chosen cluster.

for $i = 1$ **to** number of trials **do**

Bisect the selected cluster using basic K -means (i.e. 2-means).

end for

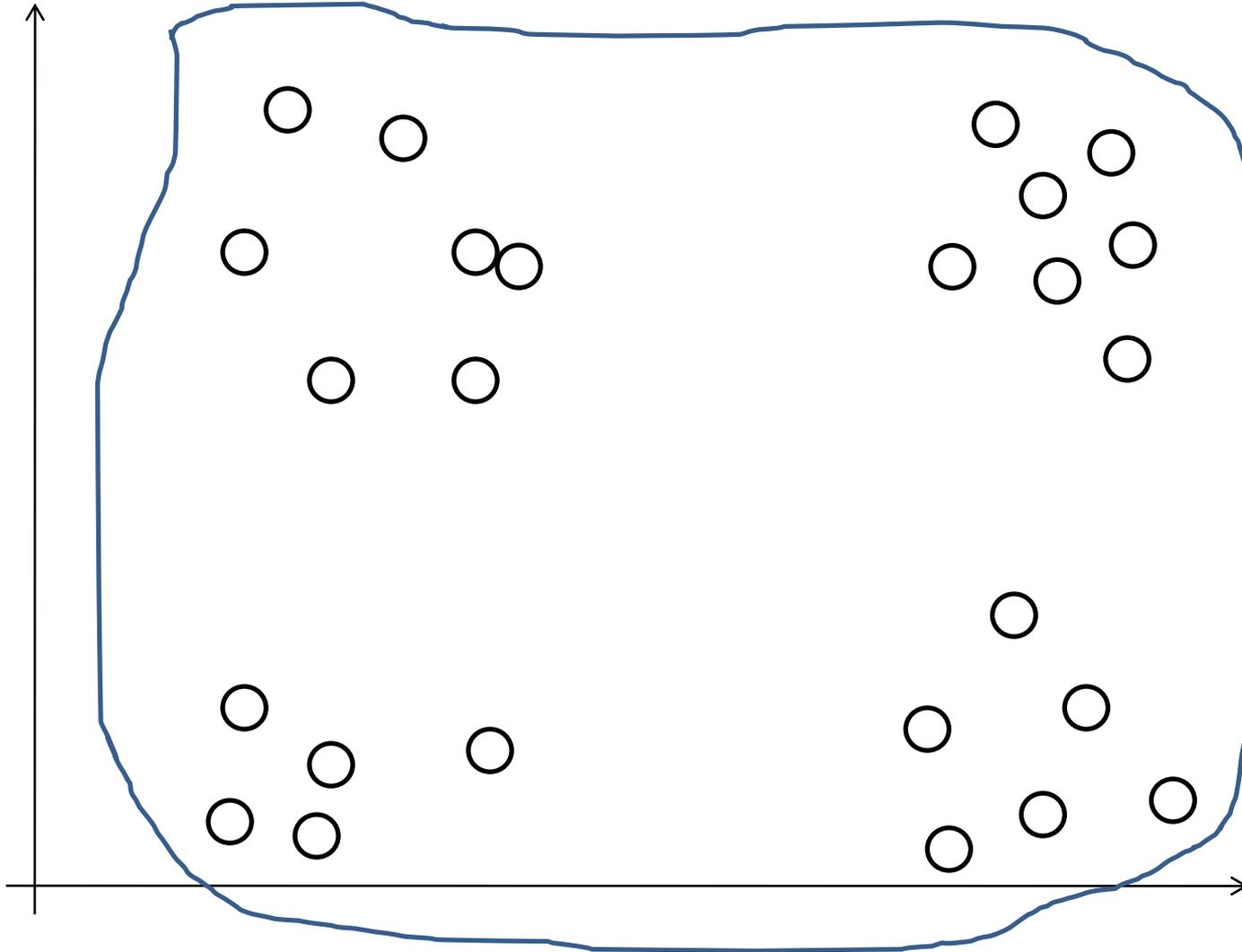
Select the two clusters from the bisection

with the lowest intra-cluster distances (SSE)

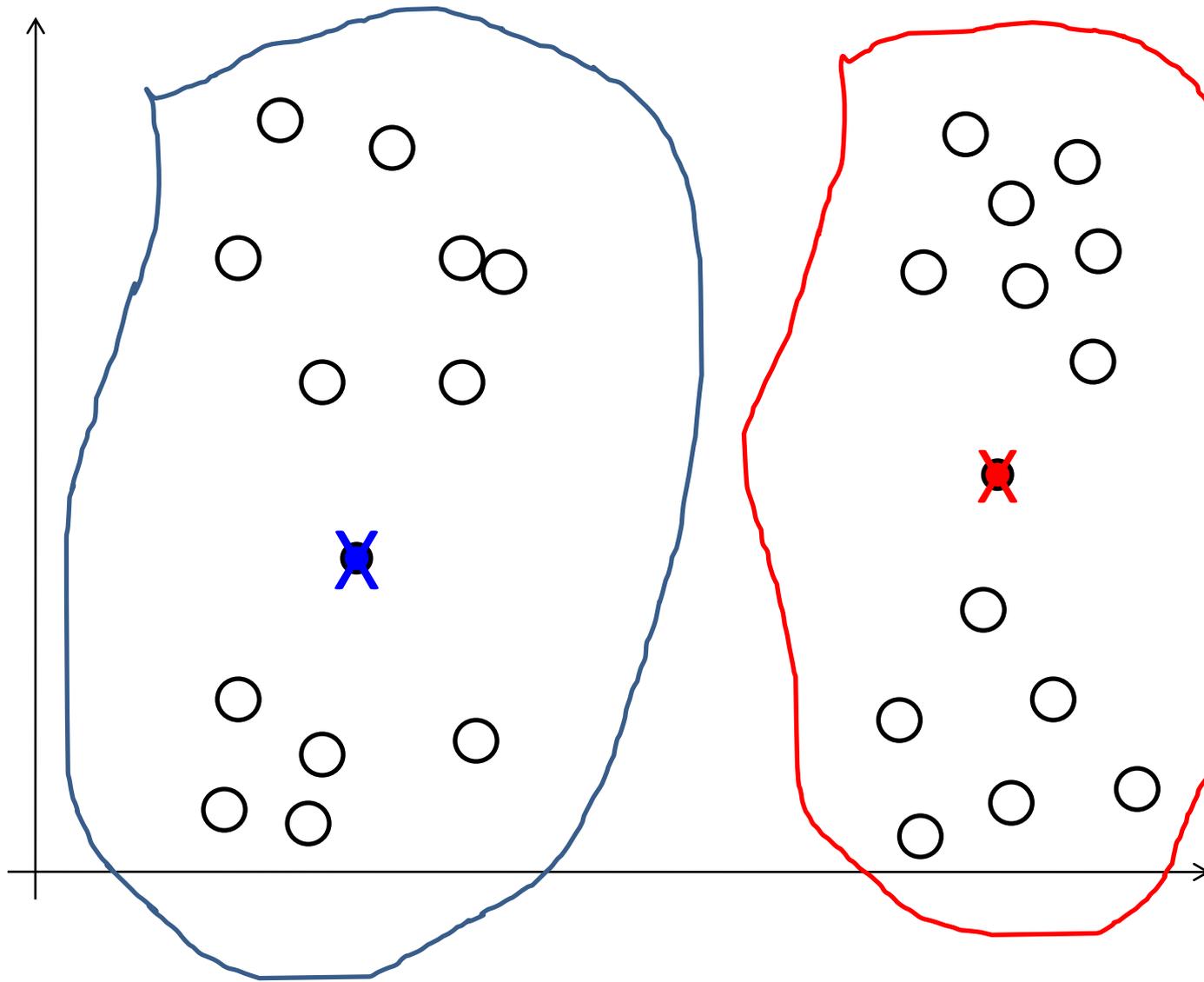
Add these two clusters to the list of clusters.

Until the list of clusters contains K clusters.

Bisecting K-means example: one initial cluster



Bisecting K-means example: bisecting initial cluster

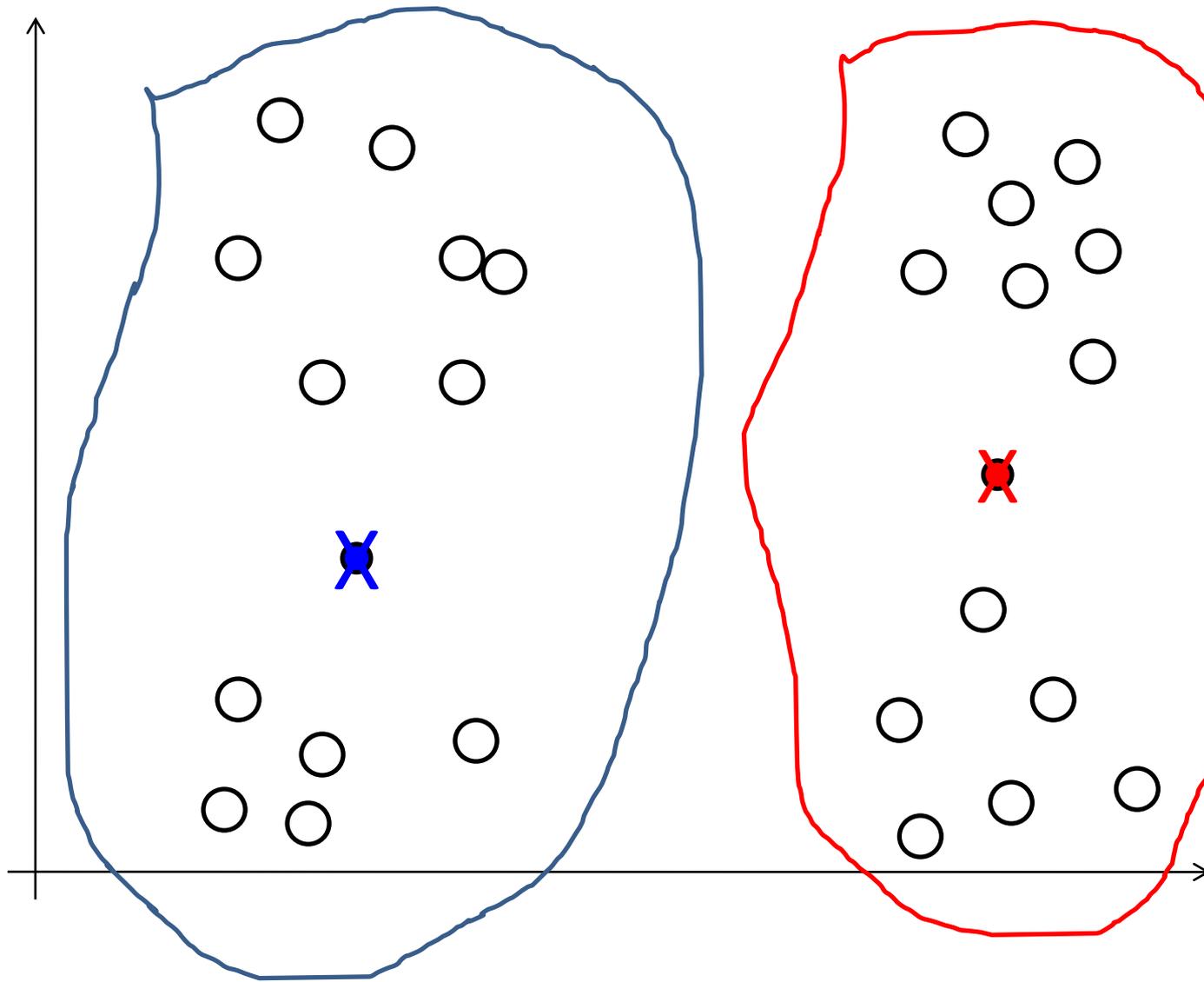


Perform K-means algorithm for $K=2$

Discovered 2 clusters

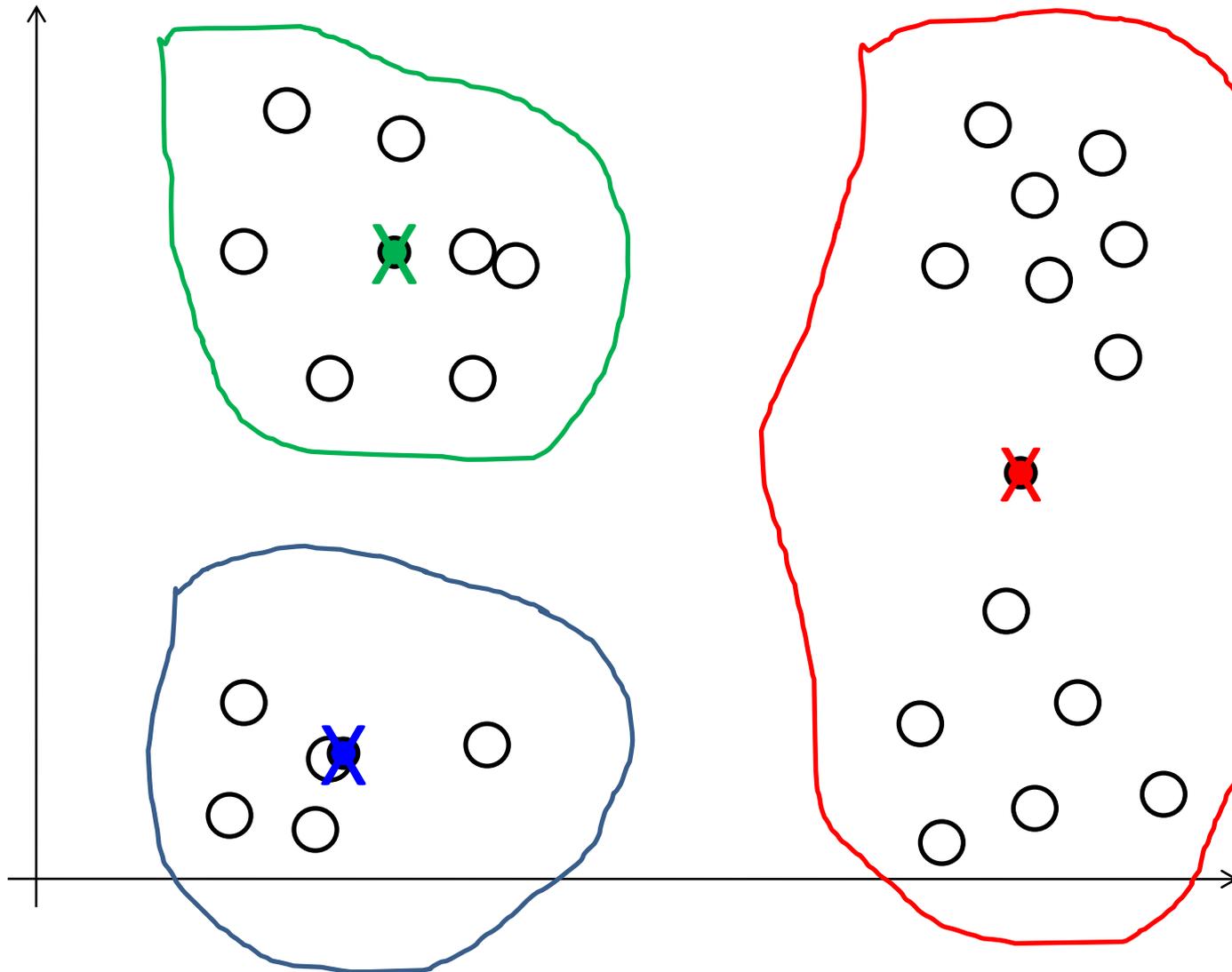
Blue cluster has larger SSE

Bisecting K-means example: bisecting blue cluster



Perform K-means algorithm for $K=2$ on a blue cluster

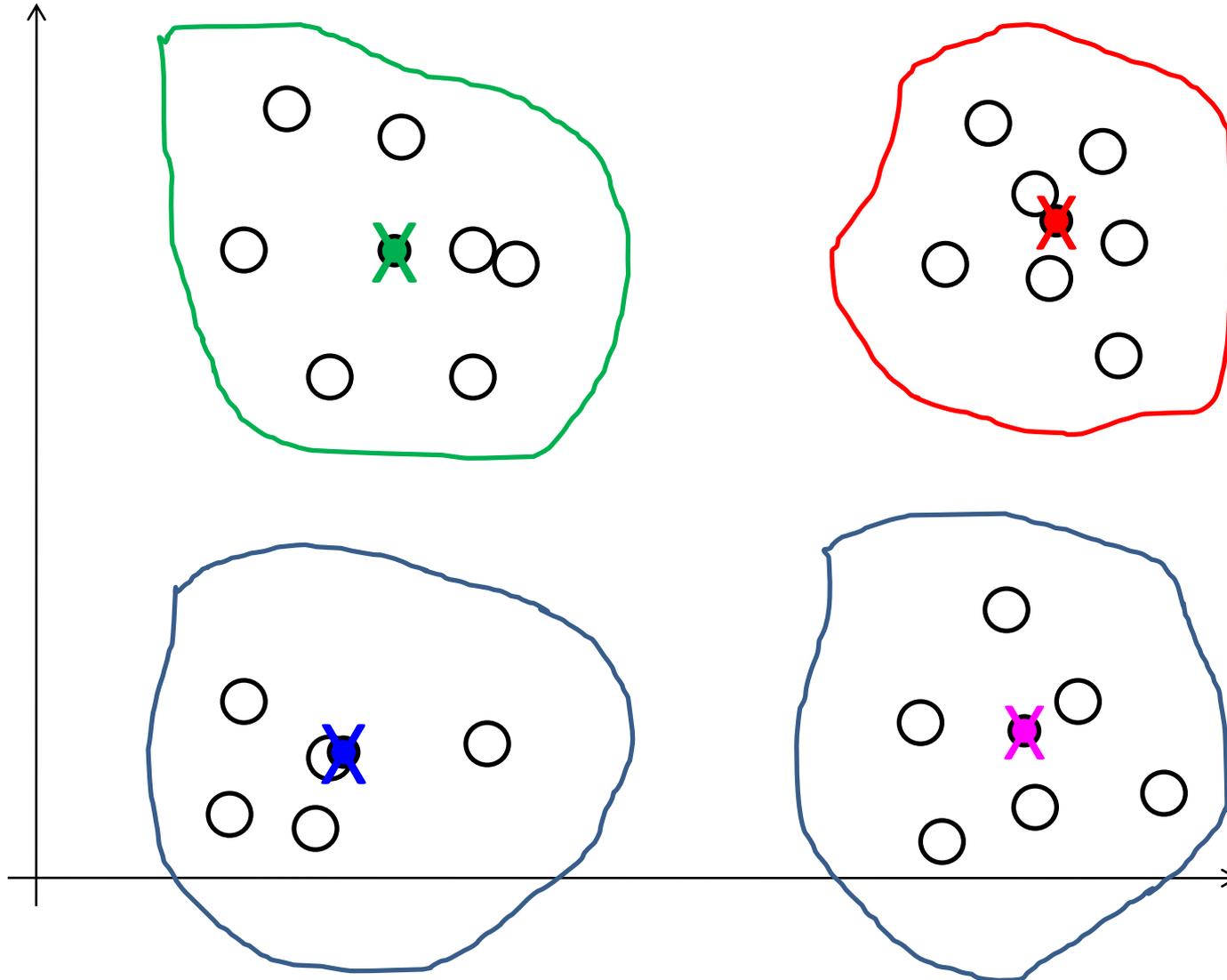
Bisecting K-means example: 3 clusters



Found 2
clusters.

Now red
cluster has
the largest
SSE.

Bisecting K-means example: bisecting red cluster



Process red
cluster.

Found 4
clusters.

Stop.

Bisecting K-means Example

Iteration 10

