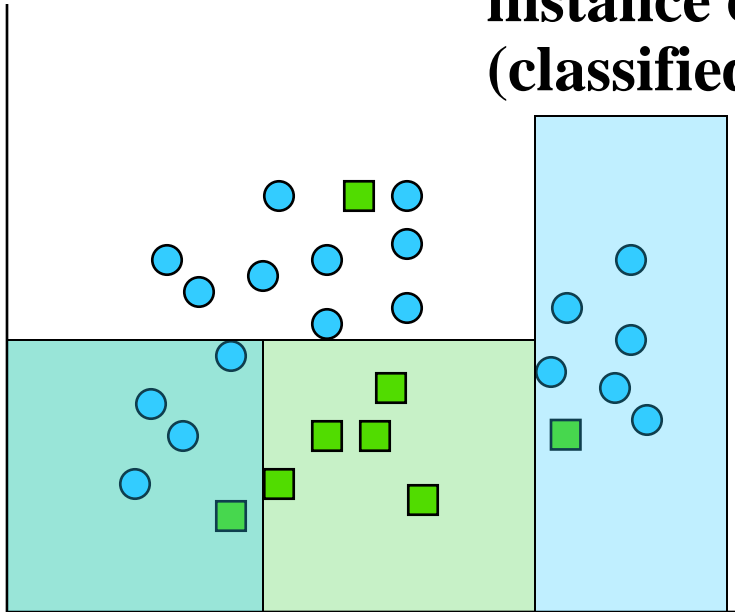


Clustering

Classification vs. Clustering

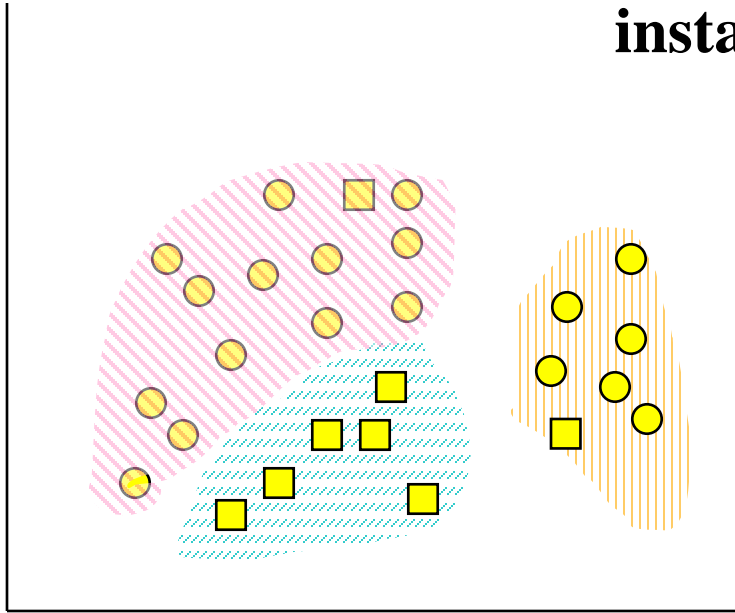
**Classification: Supervised learning:
Learns a method for predicting the
instance class from pre-labeled
(classified) instances**



Clustering

Unsupervised learning:

**Finds “natural” grouping of
instances given un-labeled data**

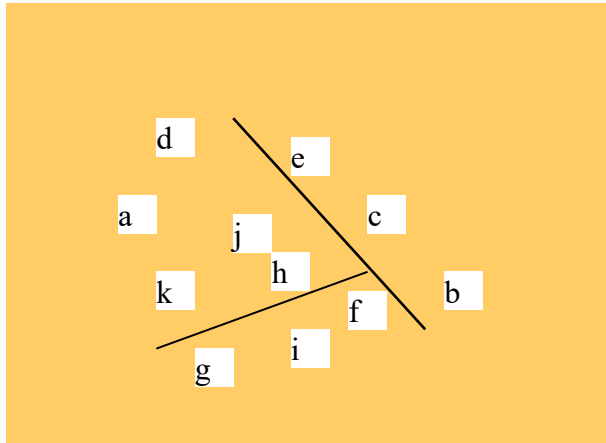


Clustering Methods

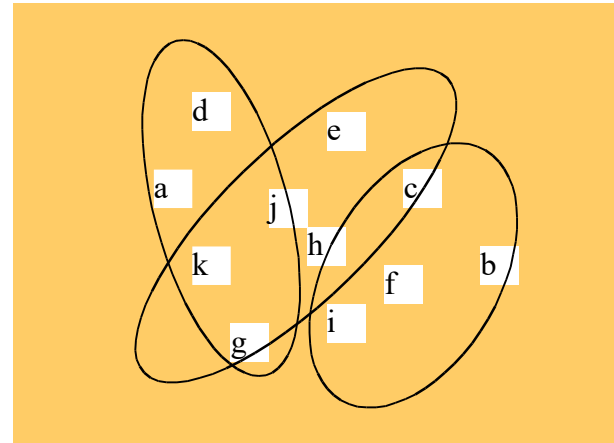
- Many different method and algorithms:
 - For numeric and/or symbolic data
 - Deterministic vs. probabilistic
 - Exclusive vs. overlapping
 - Hierarchical vs. flat
 - Top-down vs. bottom-up

Representing clusters

Non-overlapping



Overlapping



Representing clusters

Probabilistic assignment

	1	2	3
a	0.4	0.1	0.5
b	0.1	0.8	0.1
c	0.3	0.3	0.4
d	0.1	0.1	0.8
e	0.4	0.2	0.4
f	0.1	0.4	0.5
g	0.7	0.2	0.1
h	0.5	0.4	0.1

The distance function

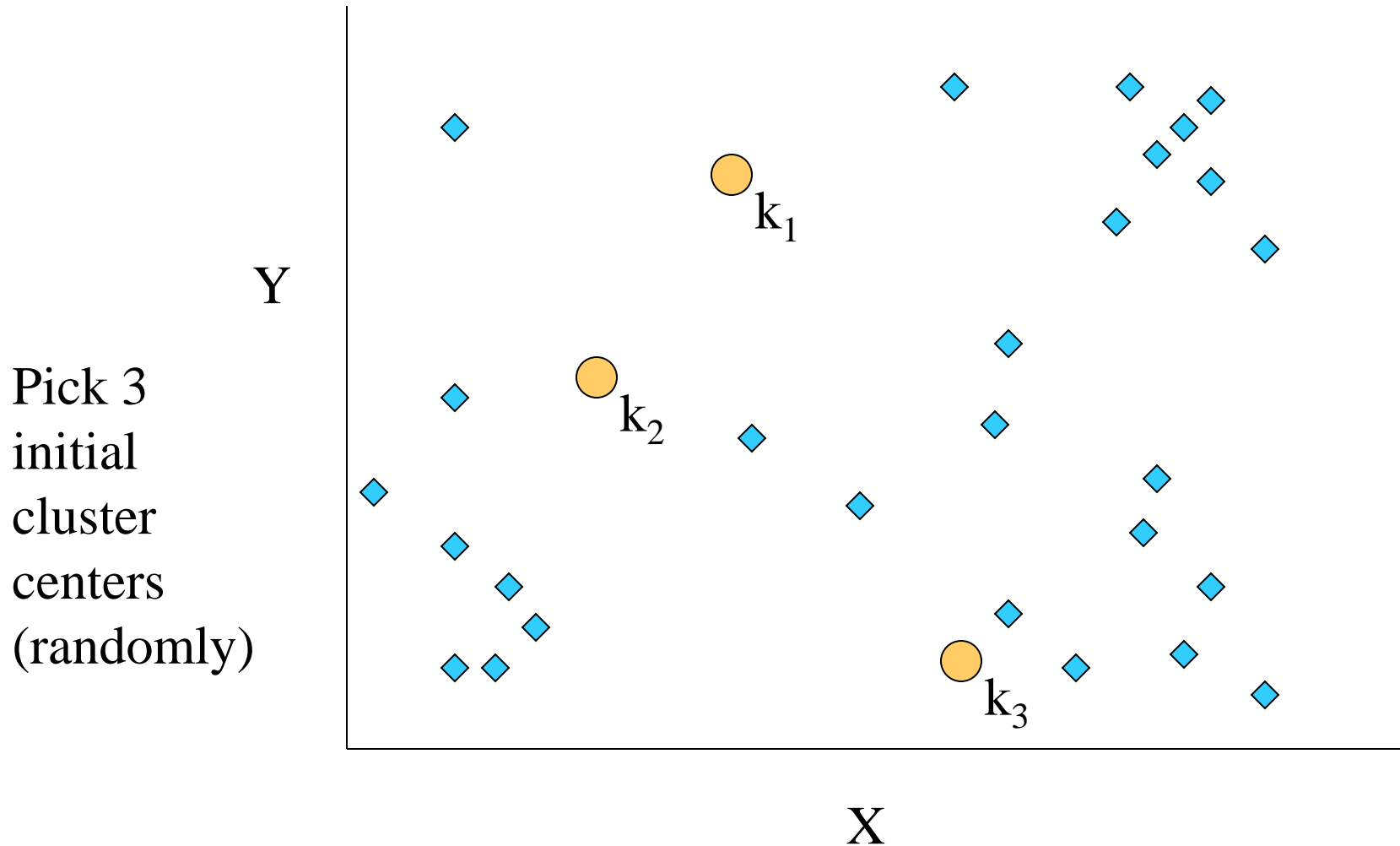
- Several numeric attributes:
 - $\text{Distance}(X,Y)$ = Euclidean distance between X,Y
- Nominal attributes: distance is set to 1 if values are different, 0 if they are equal
- Are all attributes equally important?
 - Weighting the attributes might be necessary

Simple Clustering: K-means

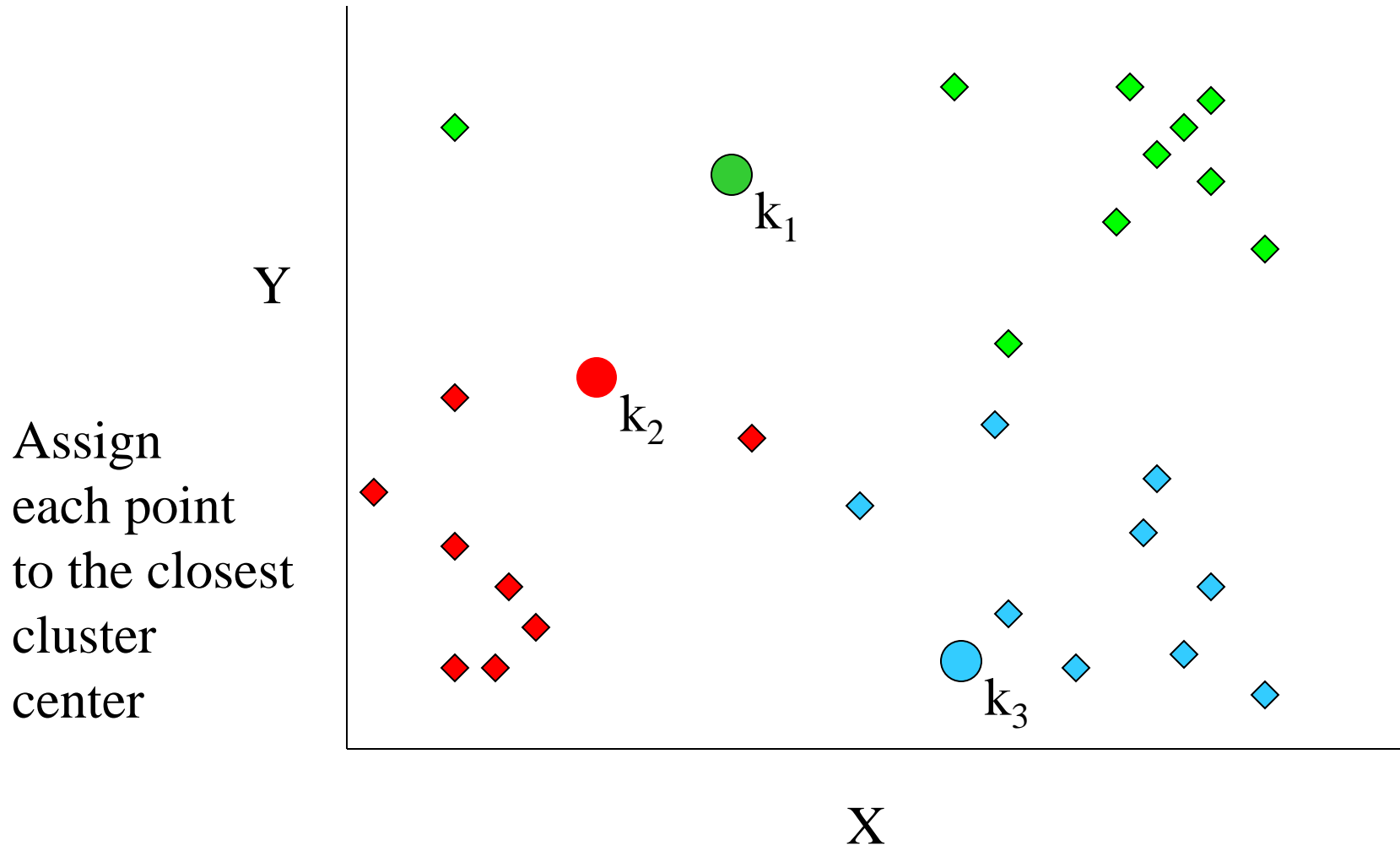
Works with numeric data only

- 1) Pick a number (K) of cluster centers (at random)
- 2) Assign every item to its nearest cluster center (e.g. using Euclidean distance)
- 3) Move each cluster center to the mean of its assigned items
- 4) Repeat steps 2,3 until convergence (change in cluster assignments less than a threshold)

K-means example, step 1

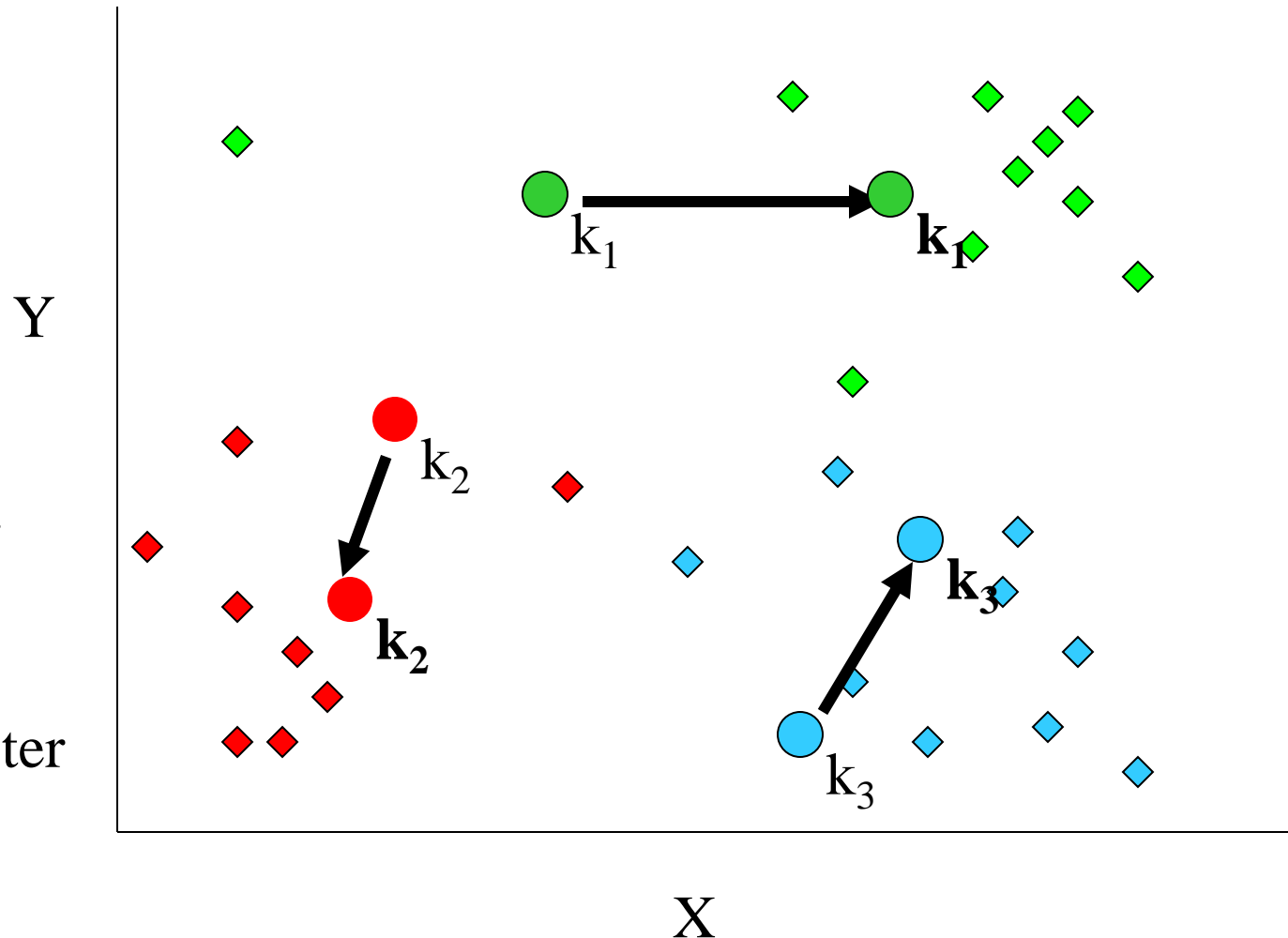


K-means example, step 2



K-means example, step 3

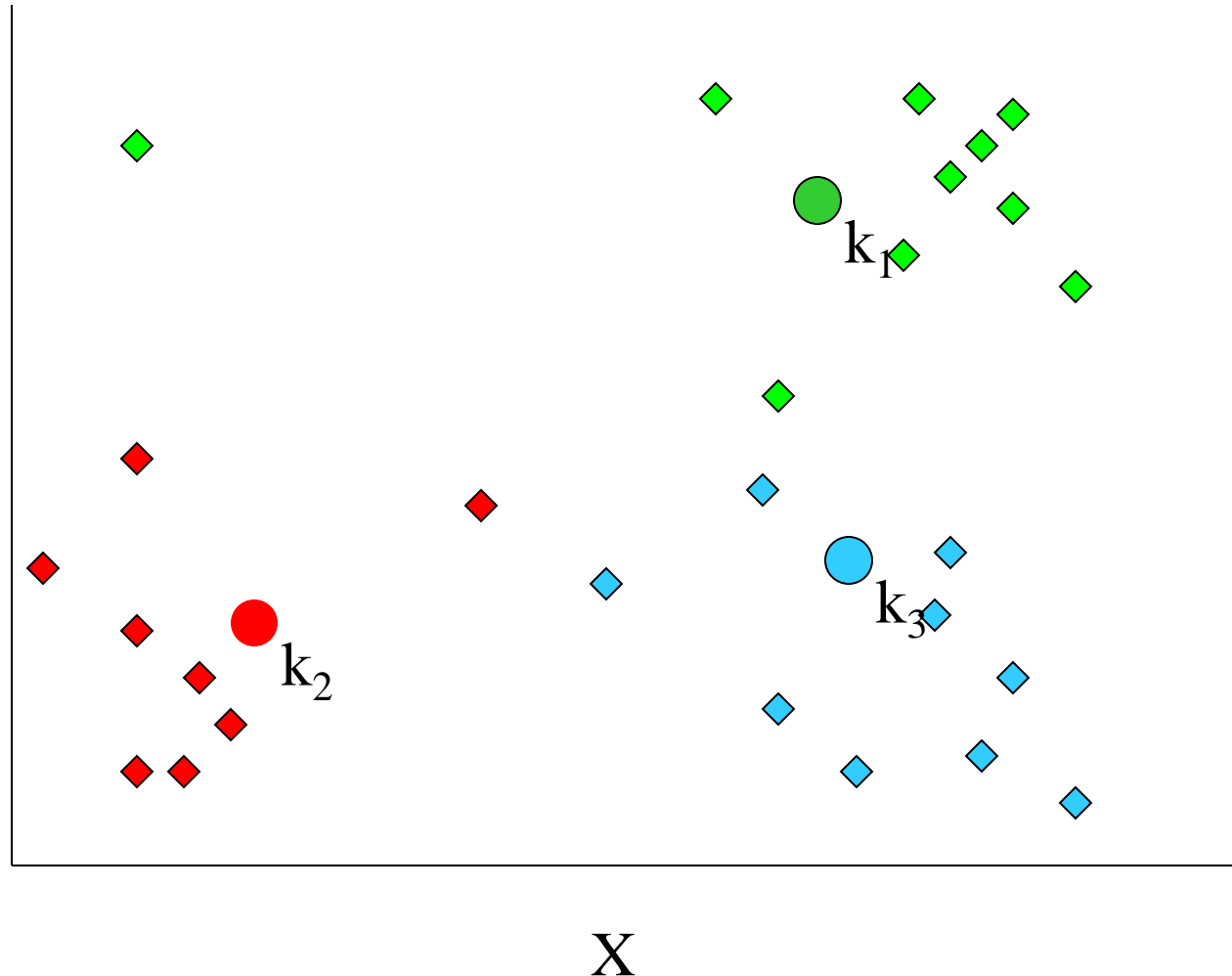
Move
each cluster
center
to the mean
of each cluster



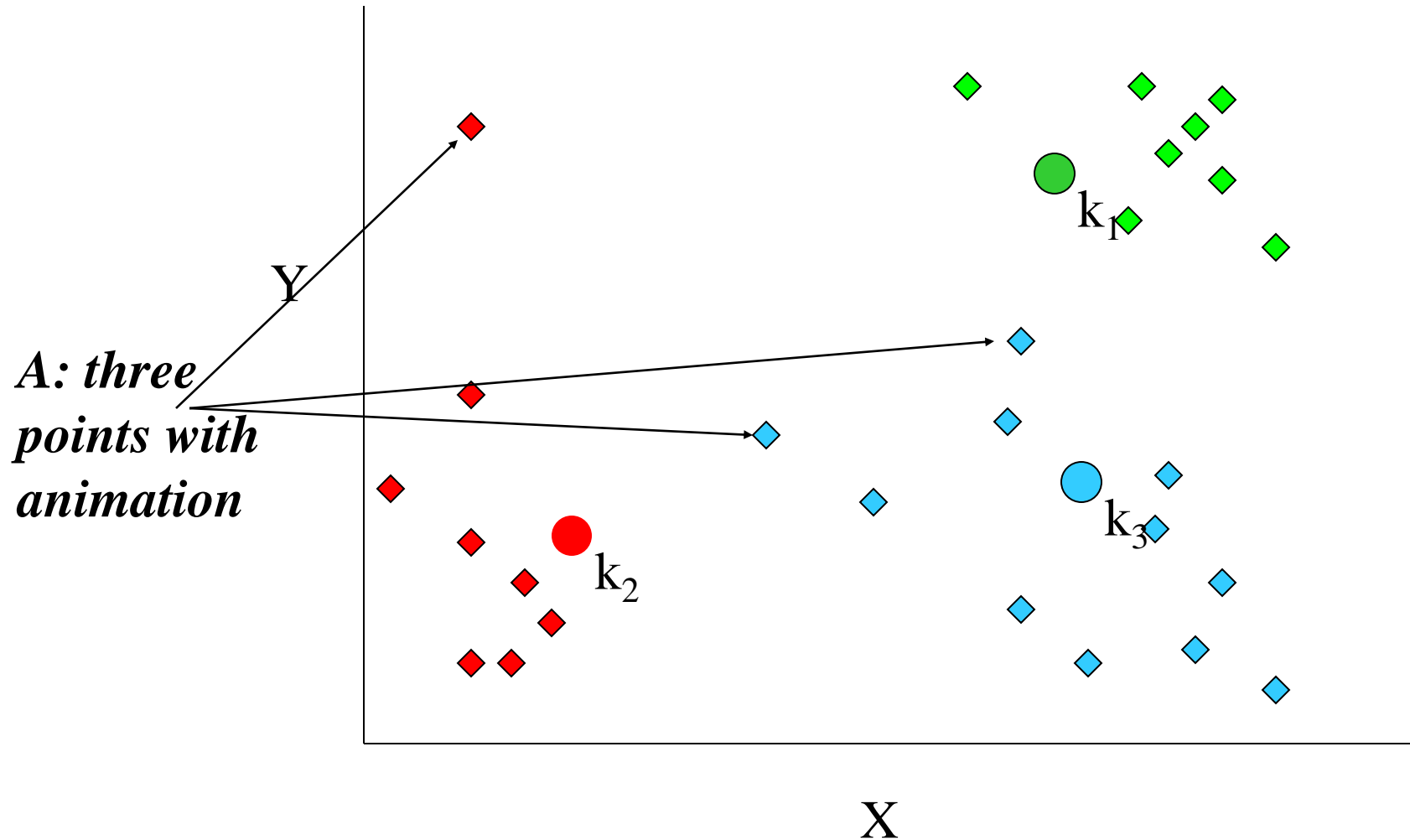
K-means example, step 4

Reassign
points
closest to a
different new
cluster center

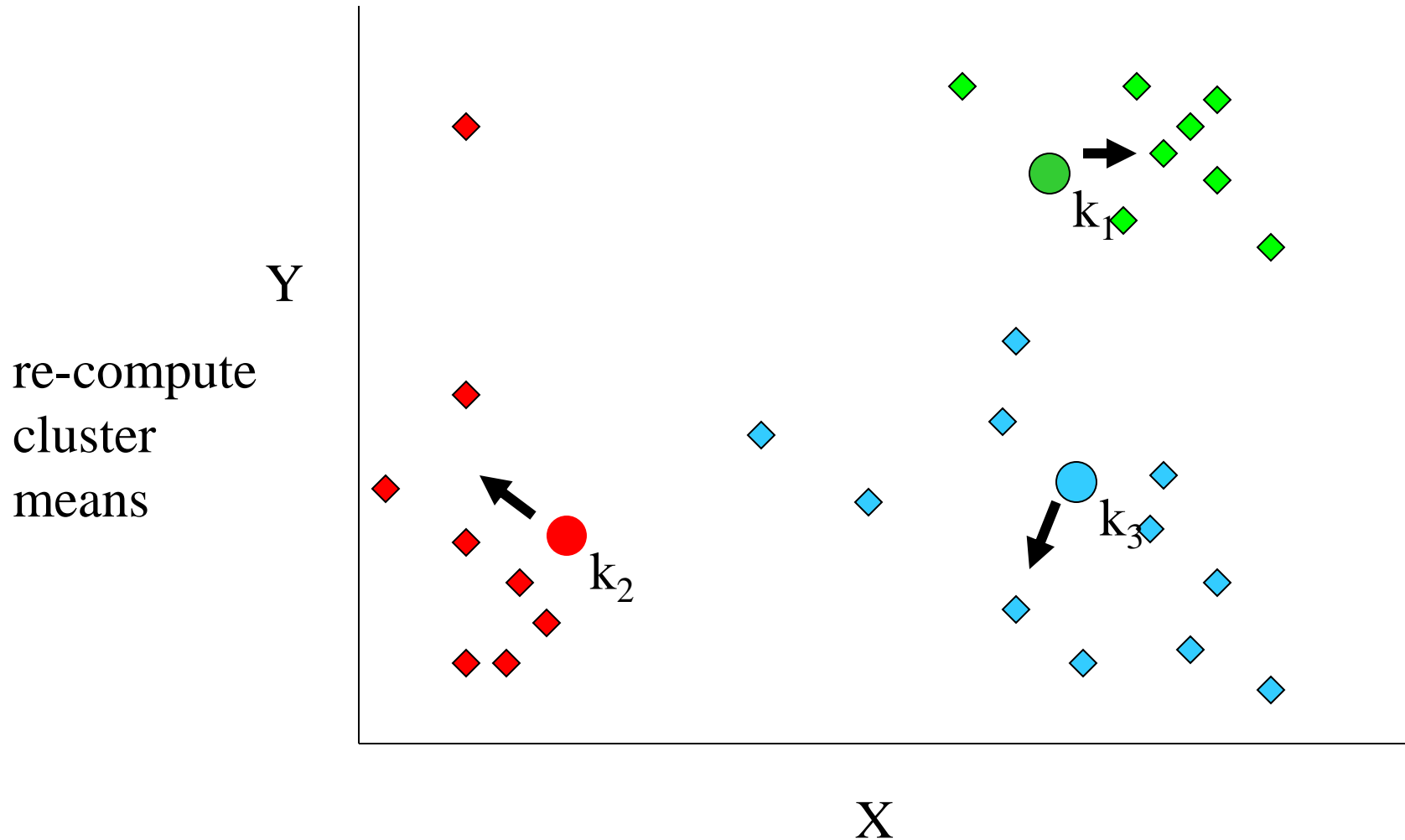
*Q: Which
points are
reassigned?*



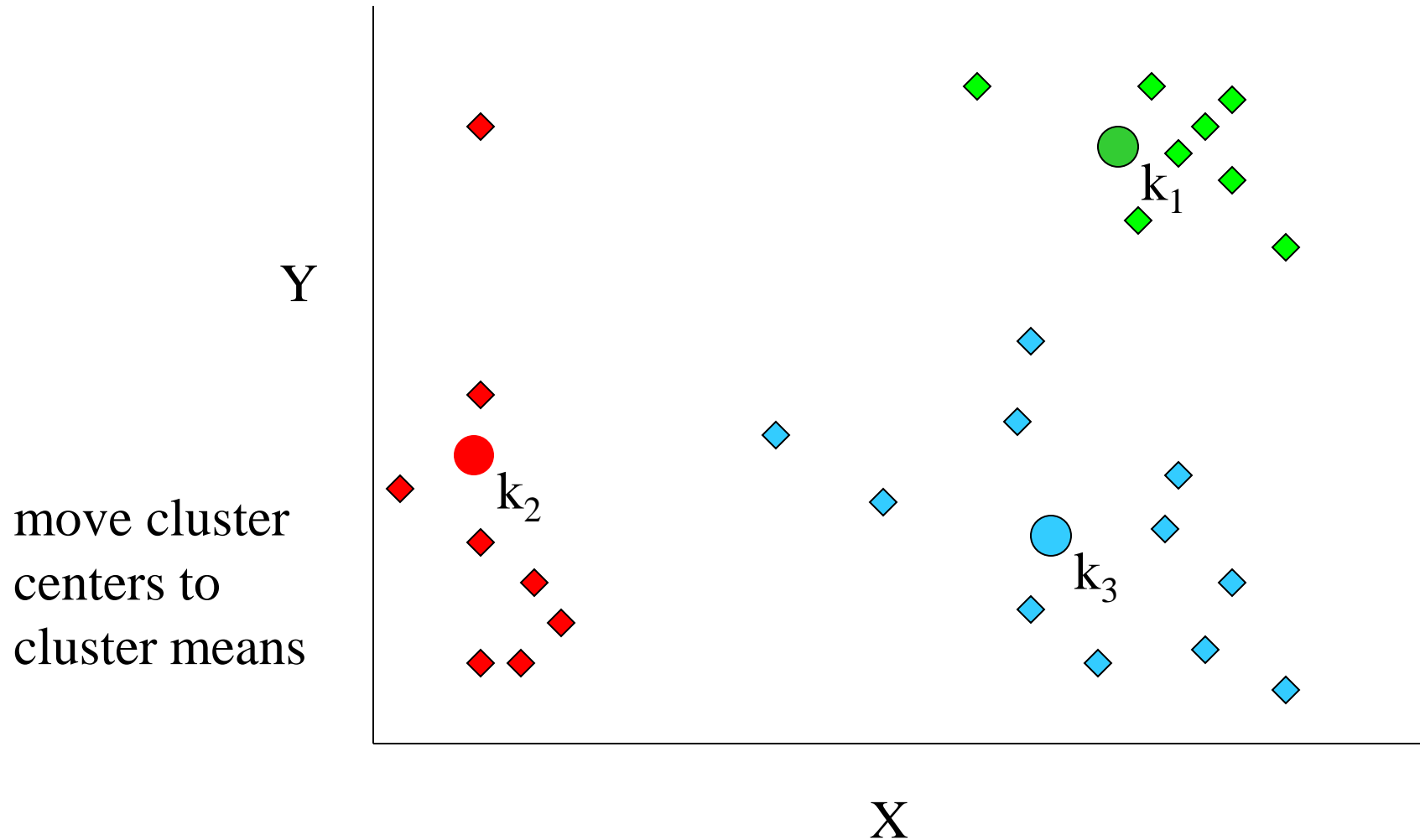
K-means example, step 4 ...



K-means example, step 4b



K-means example, step 5

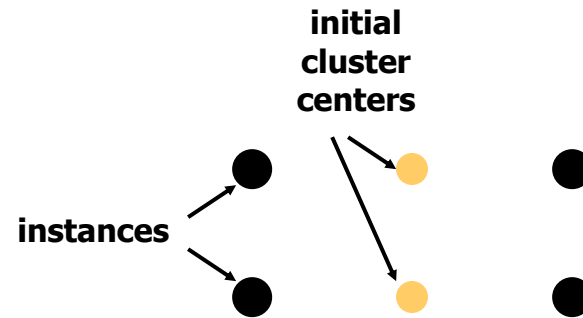


Discussion, 1

What can be the problems with K-means clustering?

Discussion, 2

- Result can vary significantly depending on initial choice of seeds (number and position)
- Can get trapped in local minimum
 - Example:



- Q: What can be done?

Discussion, 3

A: To increase chance of finding global optimum: restart with different random seeds.

K-means clustering summary

Advantages

- Simple, understandable
- items automatically assigned to clusters

Disadvantages

- Must pick number of clusters before hand
- All items forced into a cluster
- Too sensitive to outliers

K-means clustering - outliers ?

What can be done about outliers?

K-means variations

- **K-medoids** – instead of mean, use medians of each cluster
 - Mean of 1, 3, 5, 7, 9 is 5
 - Mean of 1, 3, 5, 7, 1009 is 205
 - Median of 1, 3, 5, 7, 1009 is 5
 - Median advantage: not affected by extreme values
- For large databases, use sampling

Other Clustering Approaches

- EM – probability based clustering
- Bayesian clustering
- SOM – self-organizing maps
- ...

Examples of Clustering Applications

- **Marketing:** discover customer groups and use them for targeted marketing and re-organization
- **Astronomy:** find groups of similar stars and galaxies
- **Genomics:** finding groups of gene with similar expressions
- ...

Clustering Summary

- unsupervised
- many approaches
 - K-means – simple, sometimes useful
 - K-medoids is less sensitive to outliers
 - Hierarchical clustering – works for symbolic attributes
- Evaluation is a problem