

# Clustering

kMeans,  
Expectation Maximization,  
Self-Organizing Maps

# OUTLINE

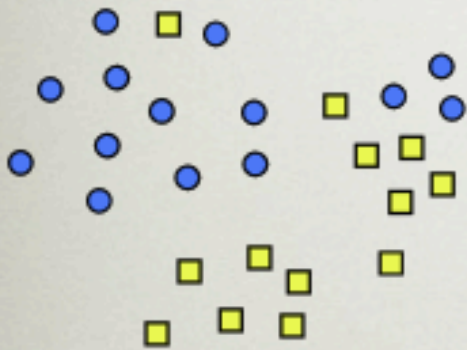
---

- K-means clustering
- Hierarchical clustering
- Incremental clustering
- Probability-based clustering
- Self-Organising Maps

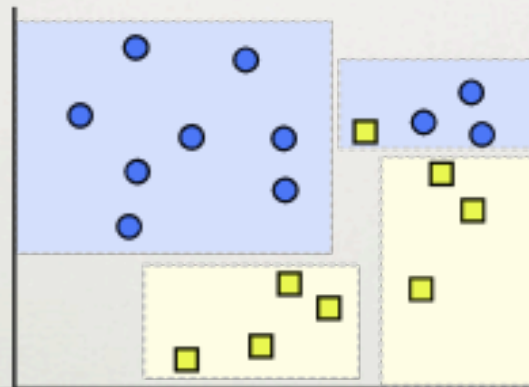
# CLASSIFICATION VS. CLUSTERING

---

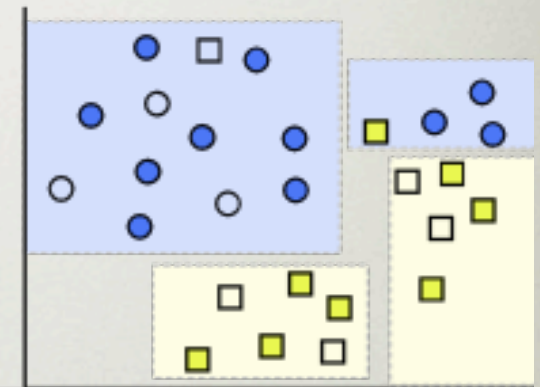
Database



Training & build Model



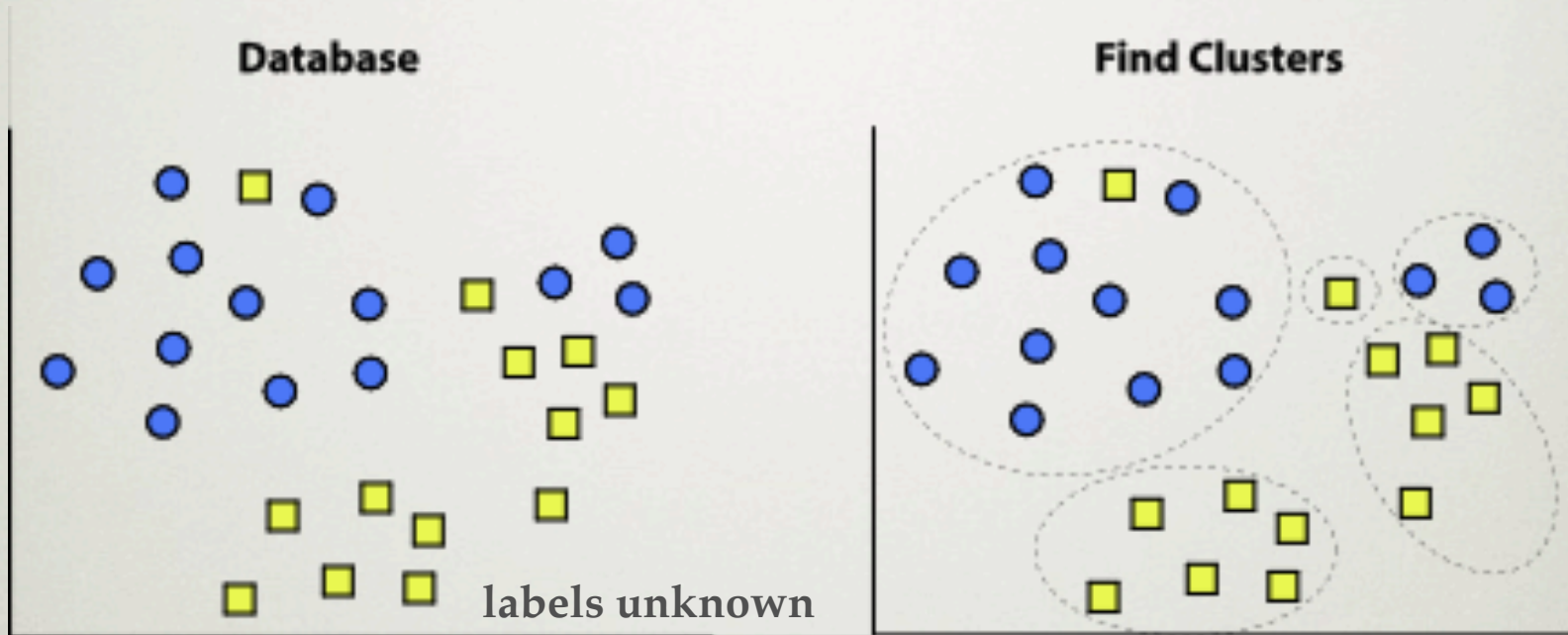
Predict Test Instances



Classification: *Supervised* learning (labels given)

# CLASSIFICATION VS. CLUSTERING

---



Clustering: *Unsupervised* learning

No labels, find “natural” grouping of instances

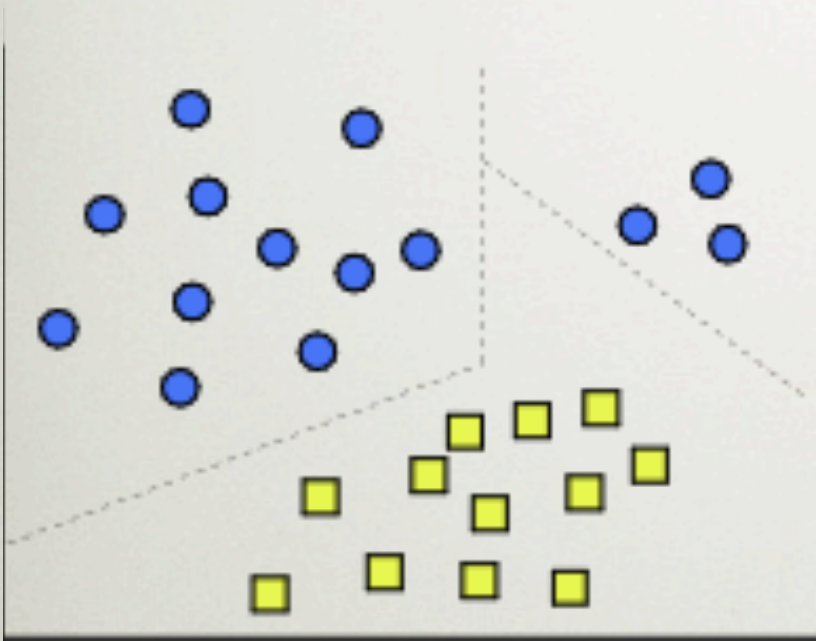
# MANY APPLICATIONS!

---

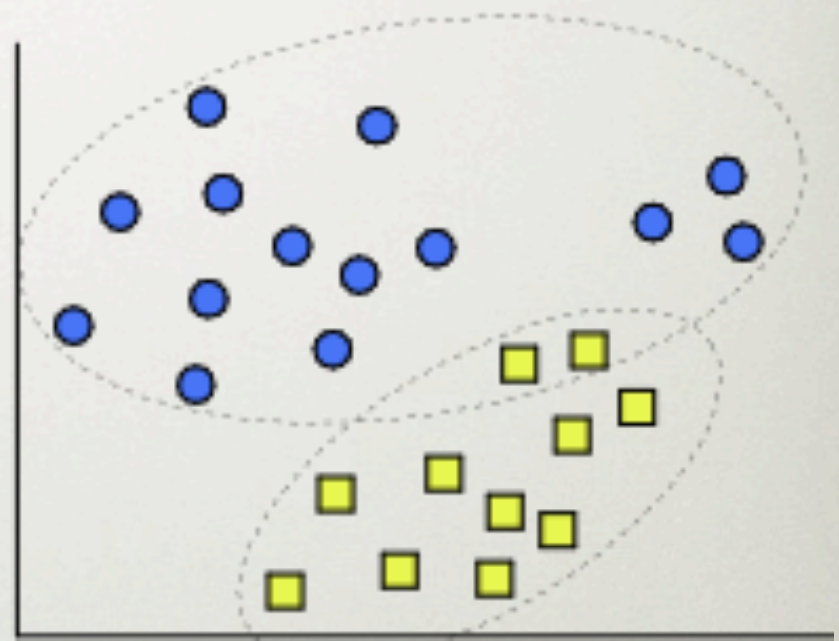
- Basically, everywhere labels are unknown / uncertain / too expensive
  - **Marketing:** find groups of similar customers
  - **Astronomy:** find groups of similar stars, galaxies
  - **Earth-quake studies:** cluster earth quake epicenters along continent faults
  - **Genomics:** find groups of genes with similar expressions

# CLUSTERING METHODS: TERMINOLOGY

---



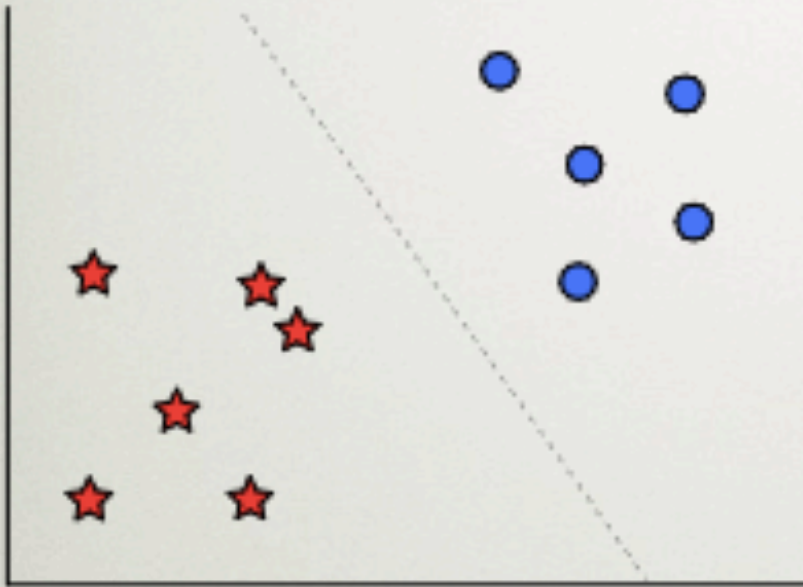
Non-overlapping



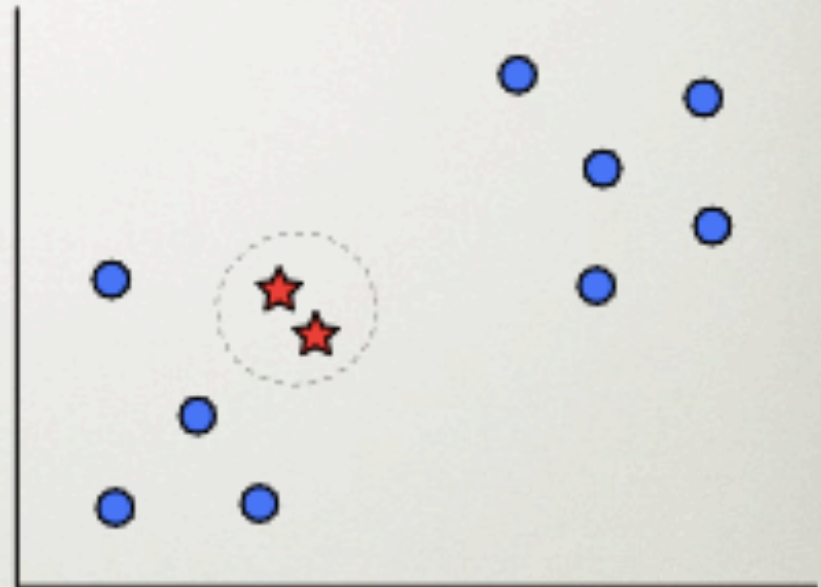
Overlapping

# CLUSTERING METHODS: TERMINOLOGY

---



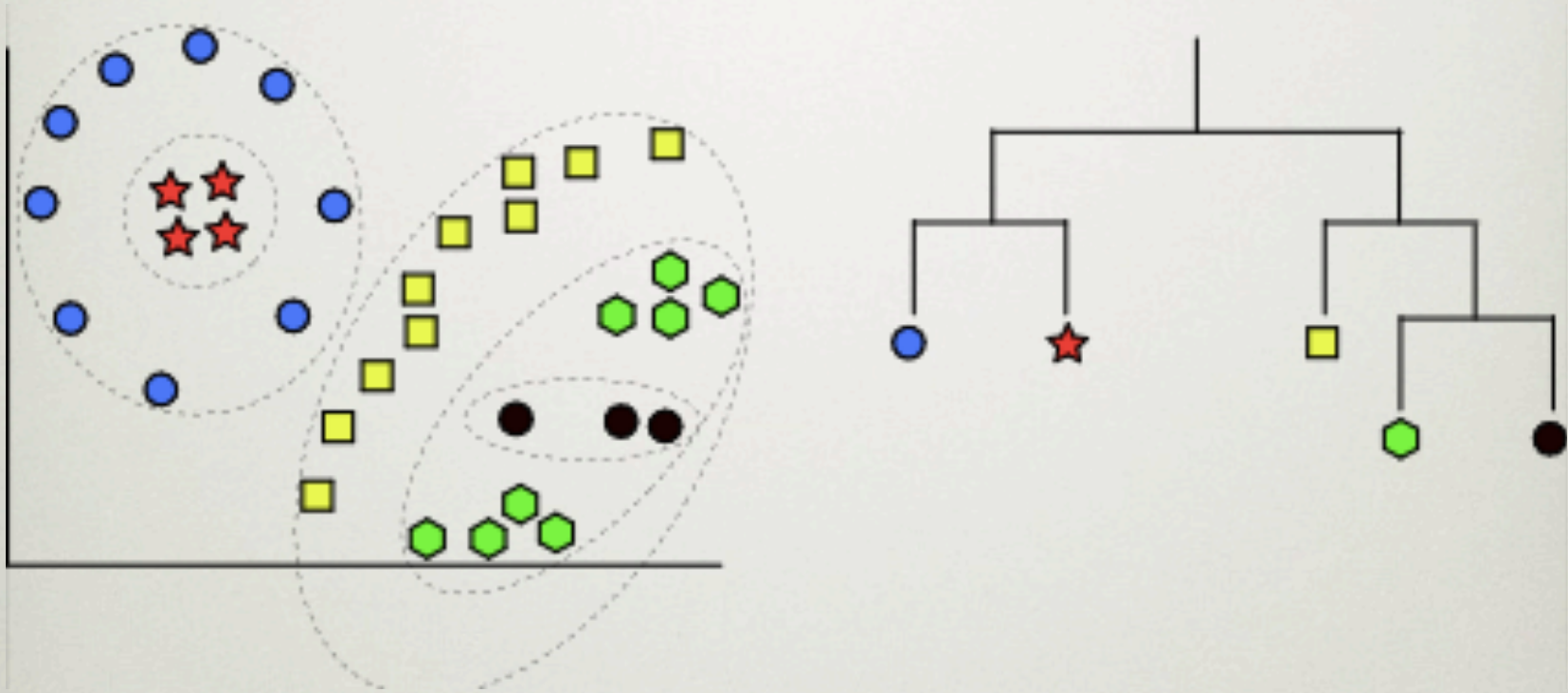
Top-down



Bottom-up  
(agglomerative)

# CLUSTERING METHODS: TERMINOLOGY

---

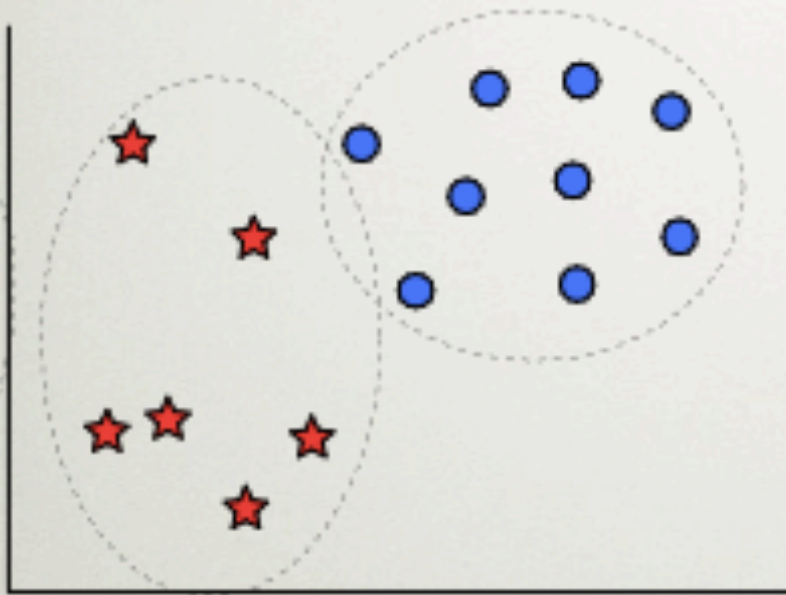


Hierarchical

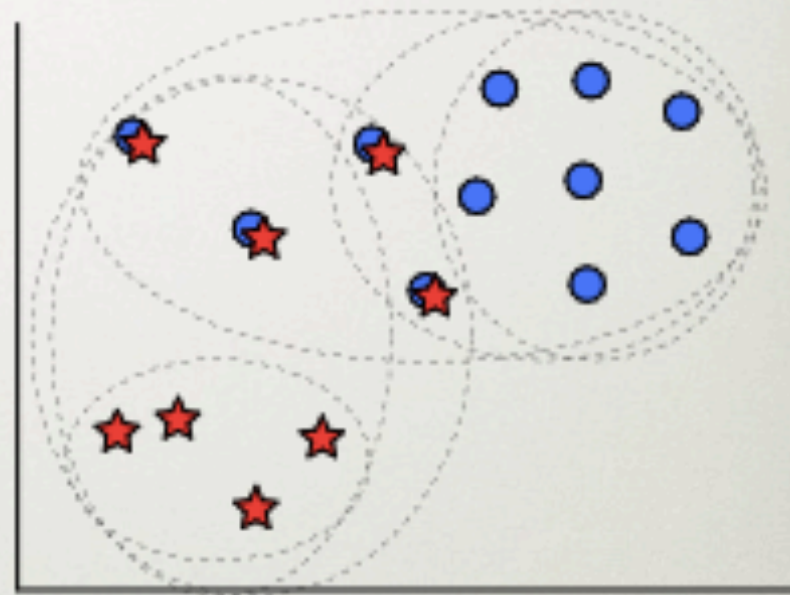


# CLUSTERING METHODS: TERMINOLOGY

---



Deterministic

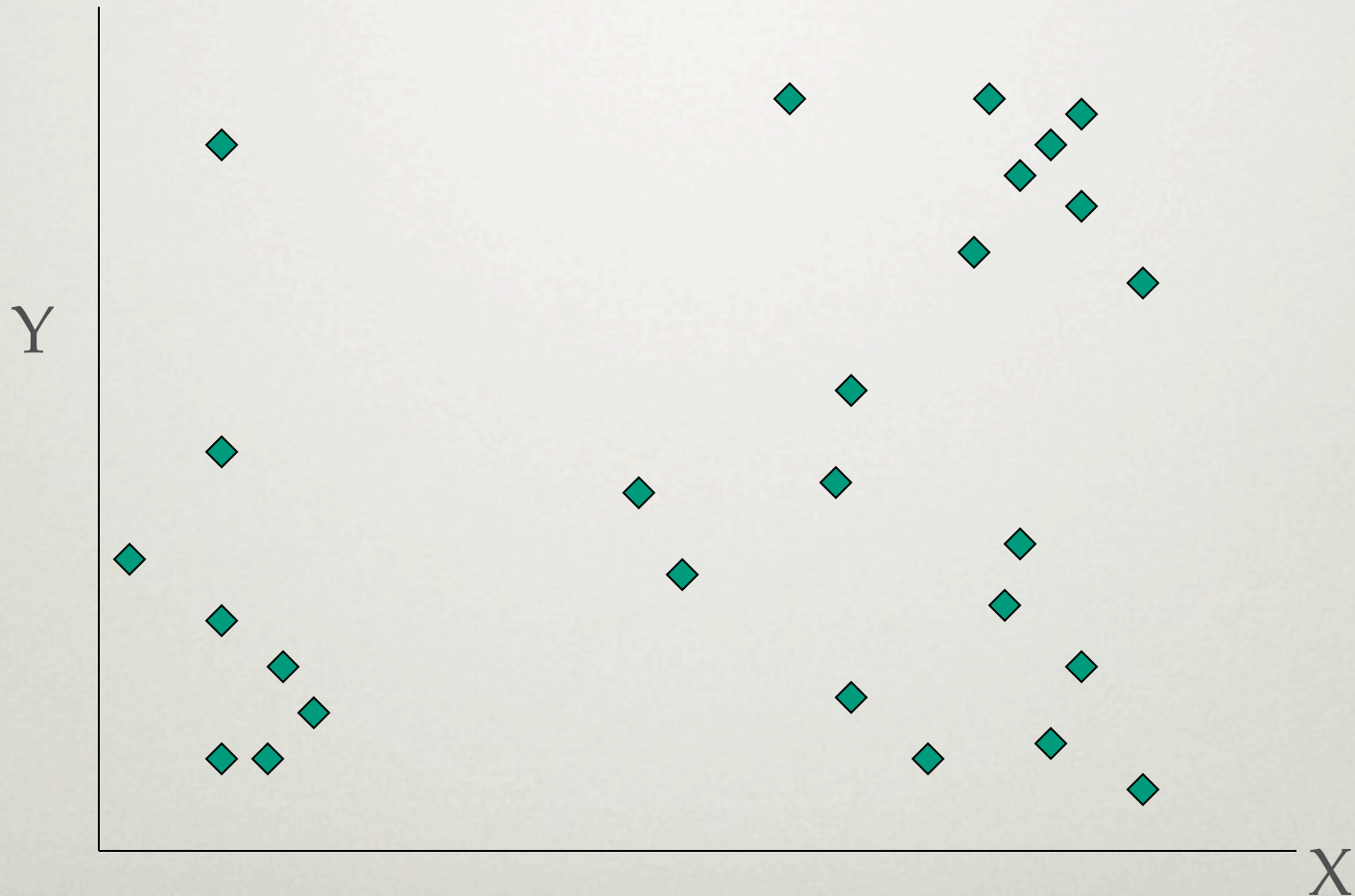


Probabilistic

# K-MEANS CLUSTERING

# K-MEANS CLUSTERING (K=3)

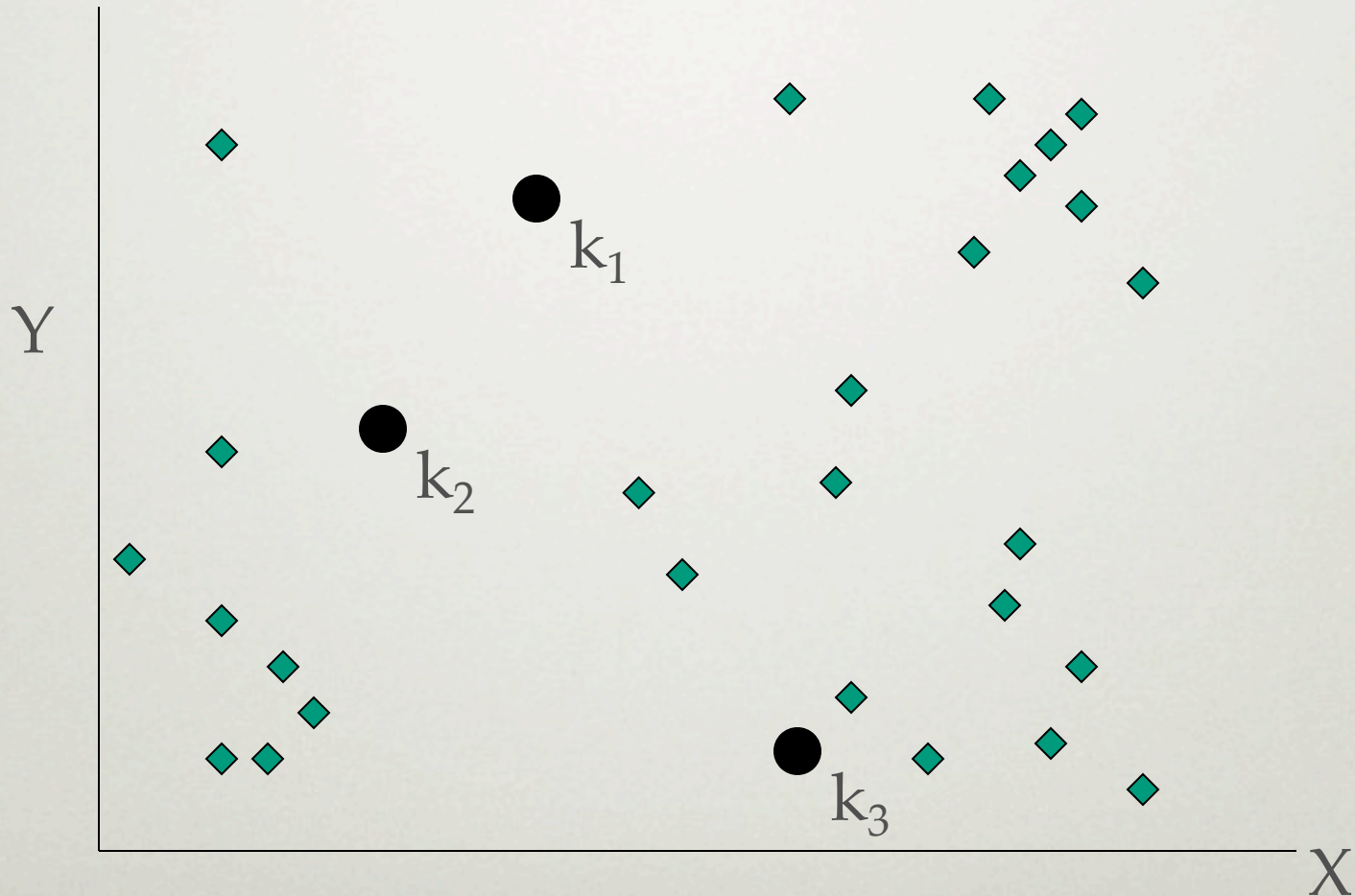
---



Pick k random points: initial cluster centers

# K-MEANS CLUSTERING (K=3)

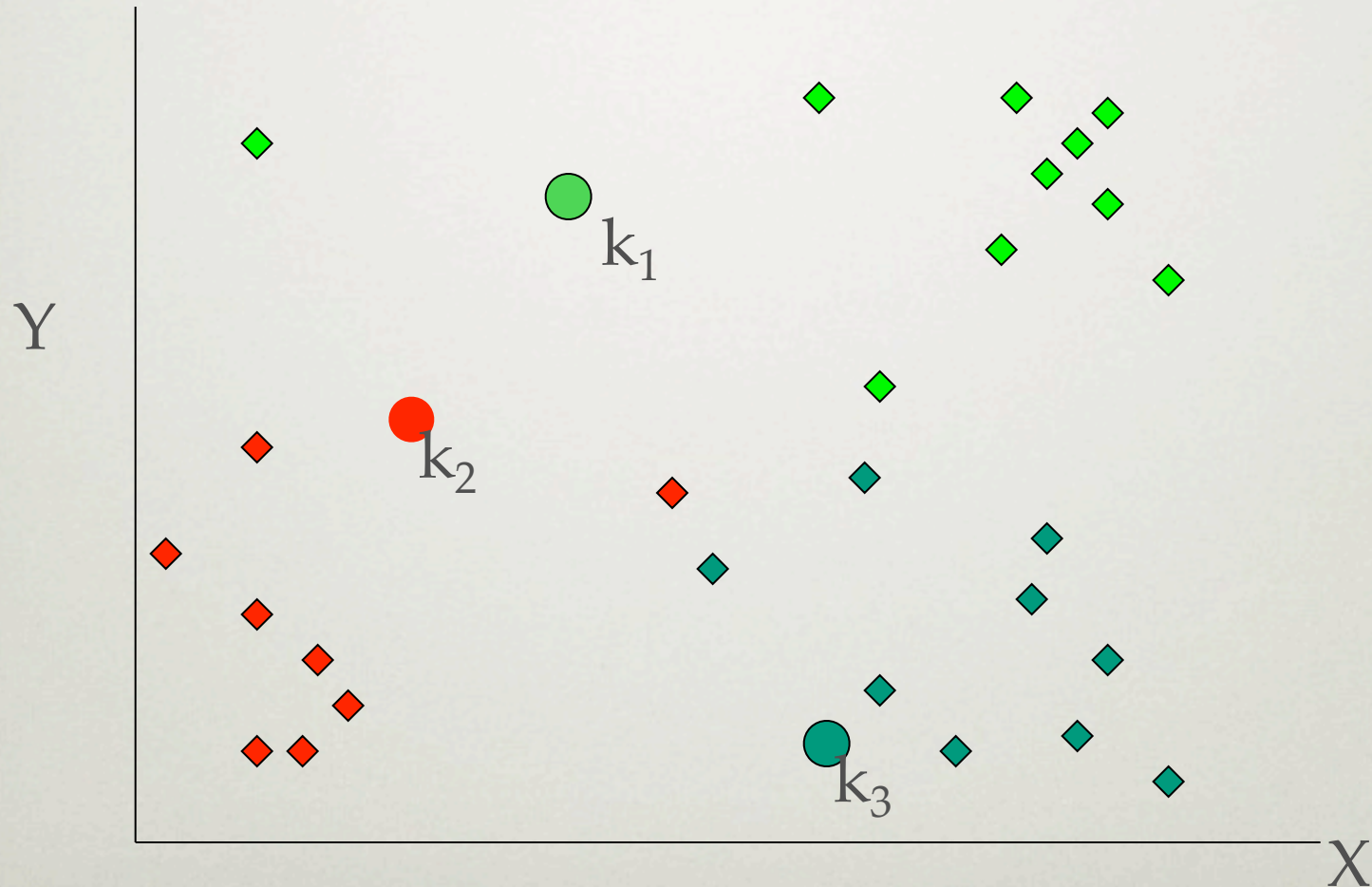
---



Pick  $k$  random points: initial cluster centers

# K-MEANS CLUSTERING (K=3)

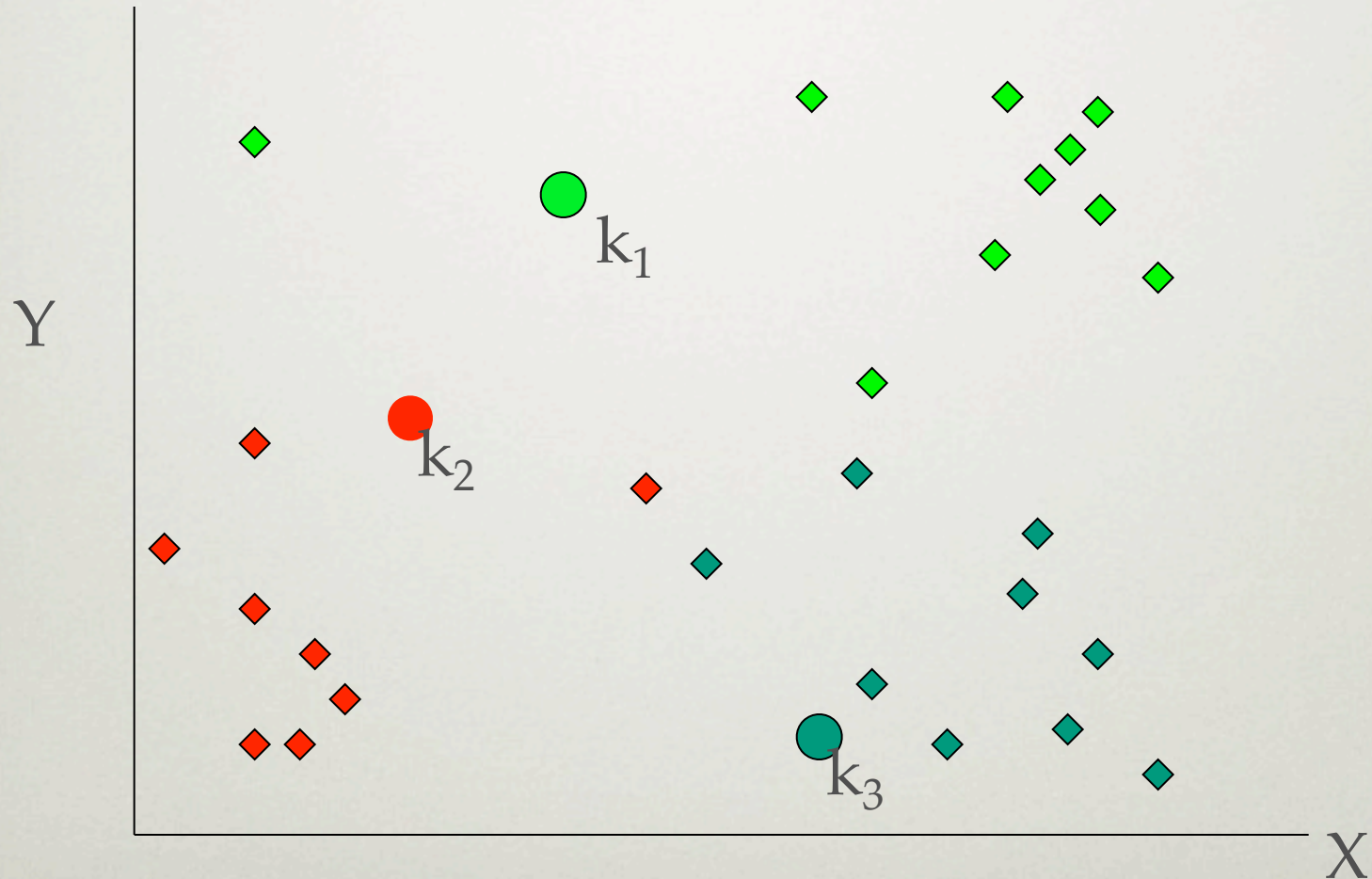
---



Assign each point to nearest cluster center

# K-MEANS CLUSTERING (K=3)

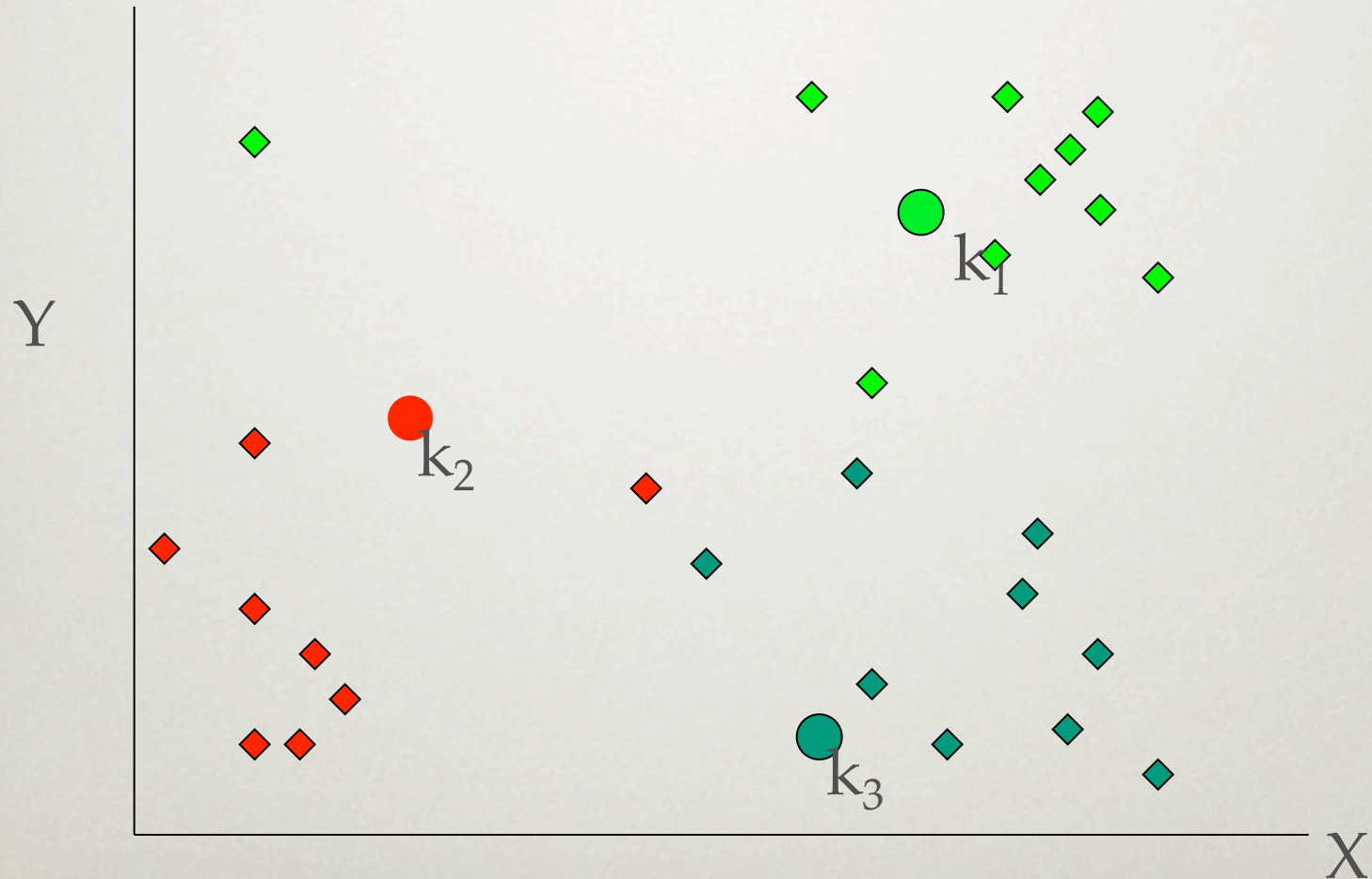
---



Move cluster centers to *mean* of each cluster

# K-MEANS CLUSTERING (K=3)

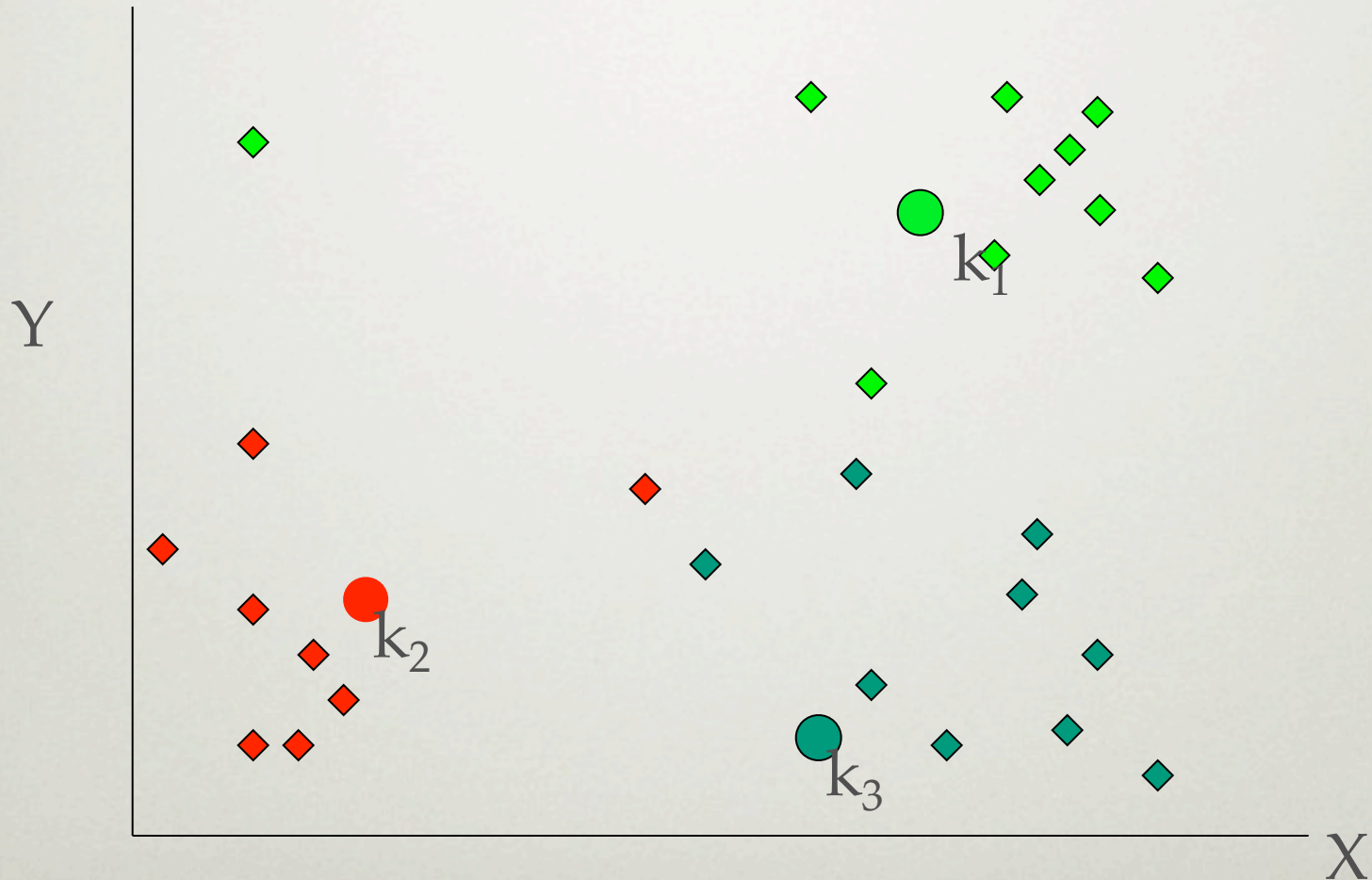
---



Move cluster centers to *mean* of each cluster

# K-MEANS CLUSTERING (K=3)

---

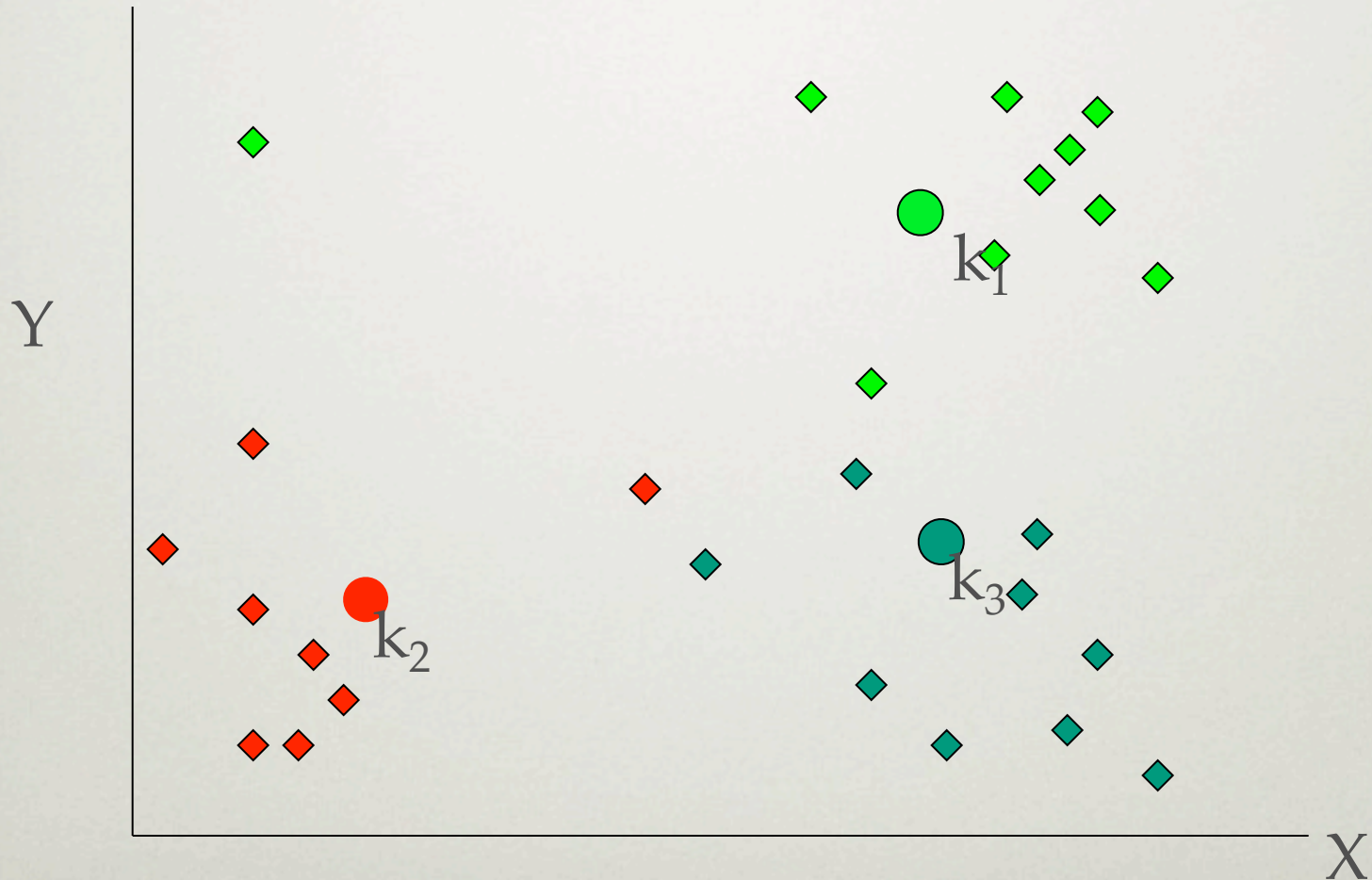


Move cluster centers to *mean* of each cluster



# K-MEANS CLUSTERING (K=3)

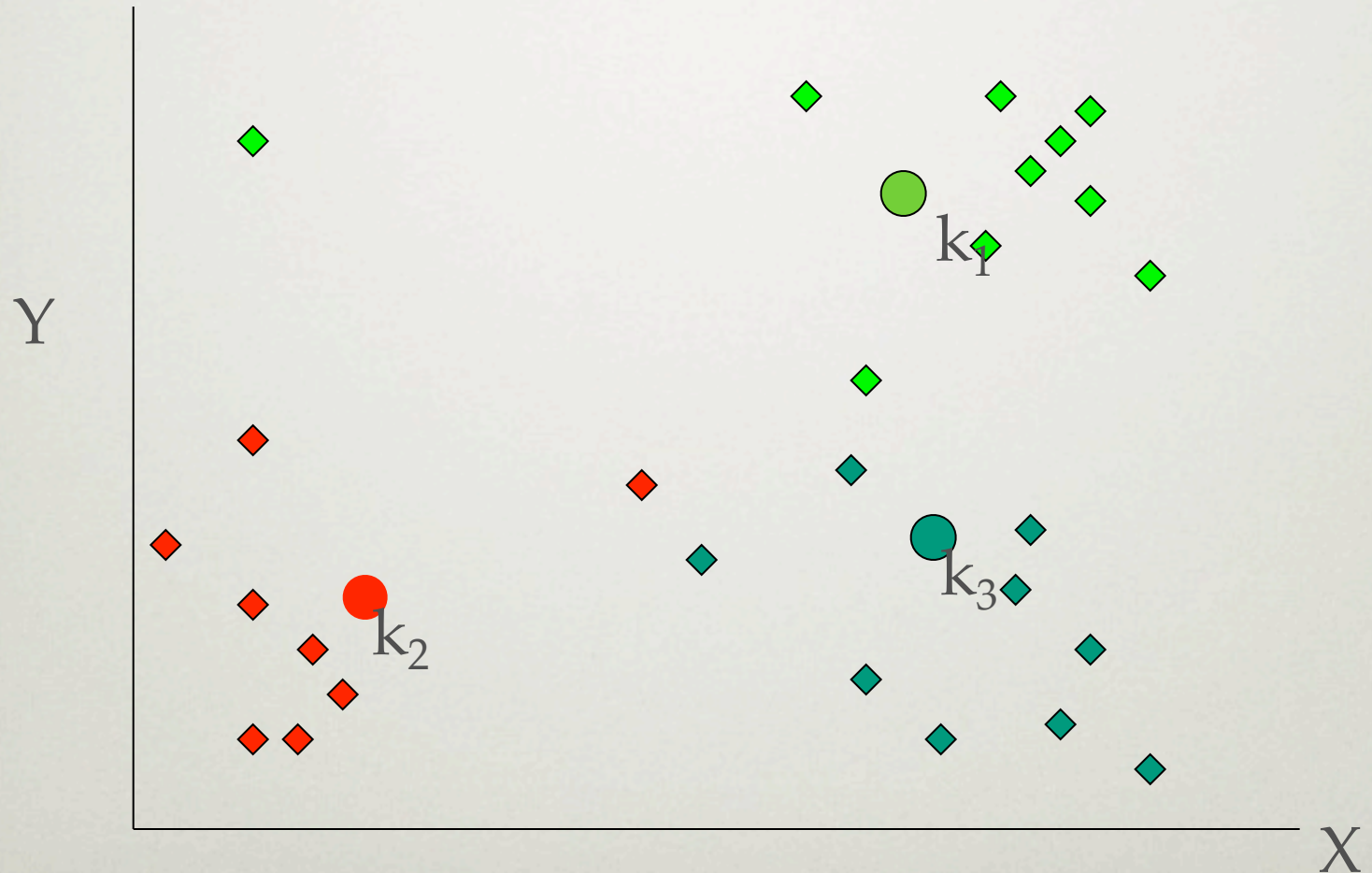
---



Move cluster centers to *mean* of each cluster

# K-MEANS CLUSTERING (K=3)

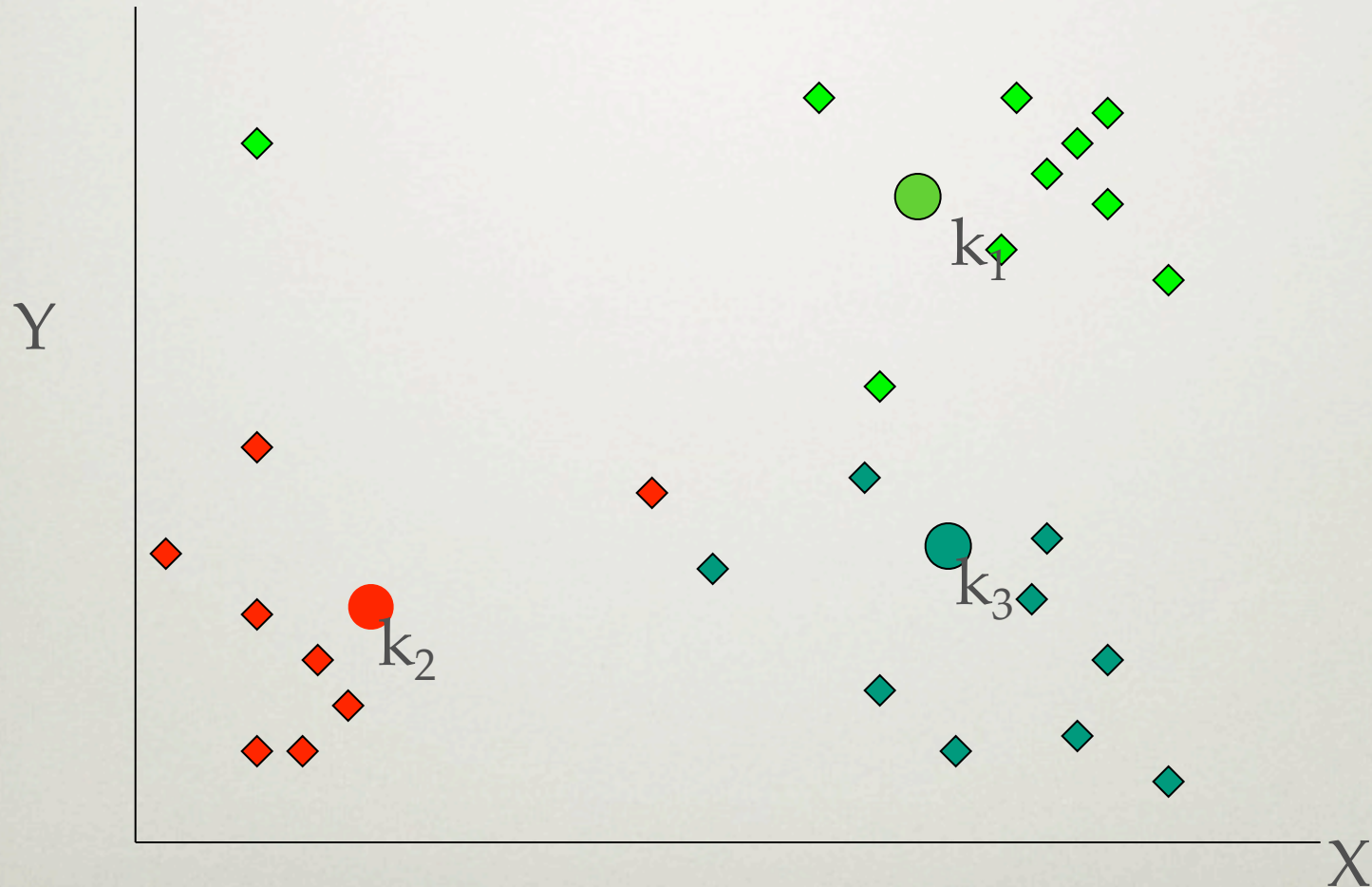
---



Reassign points to nearest cluster center

# K-MEANS CLUSTERING (K=3)

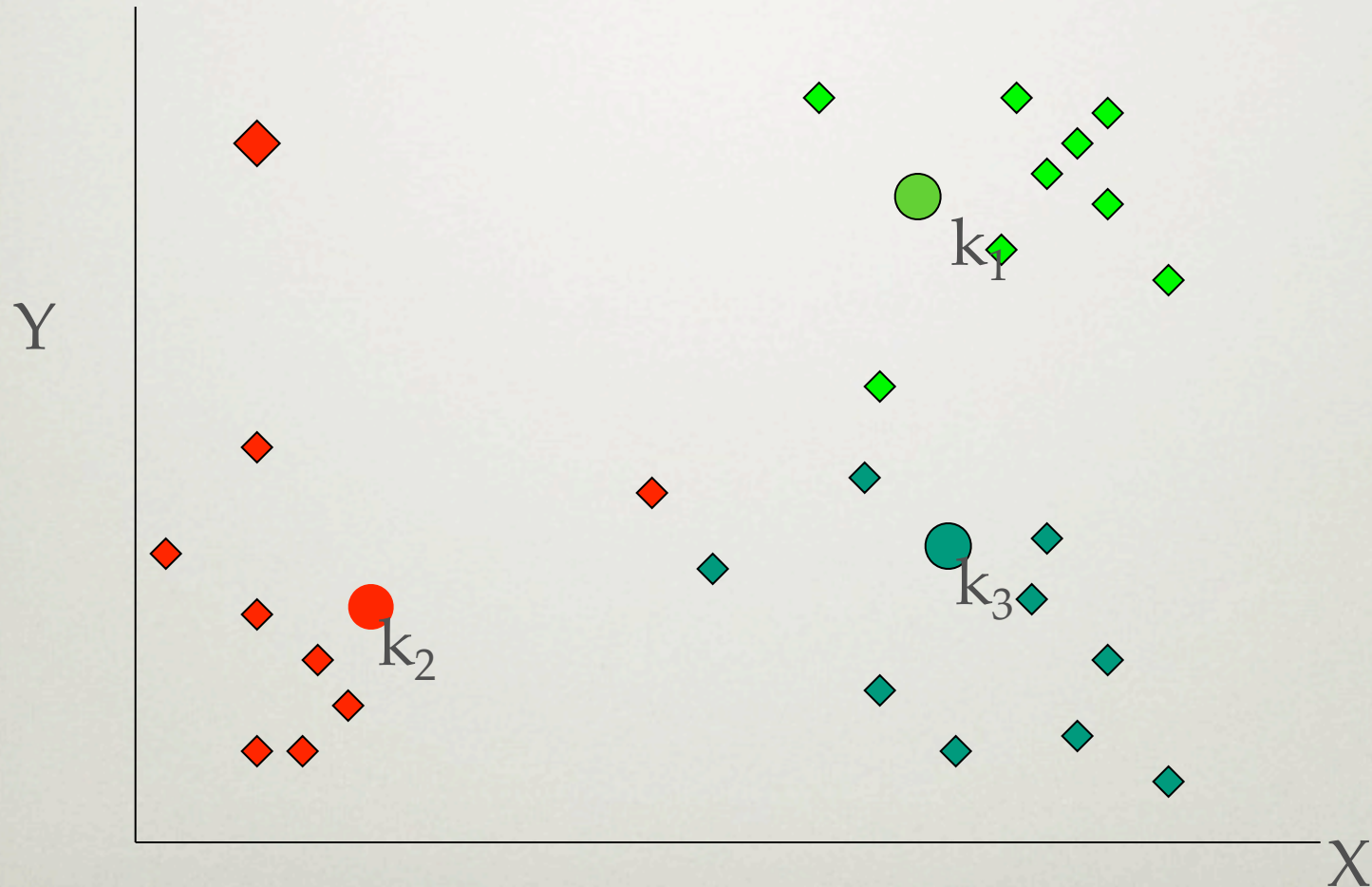
---



Reassign points to nearest cluster center

# K-MEANS CLUSTERING (K=3)

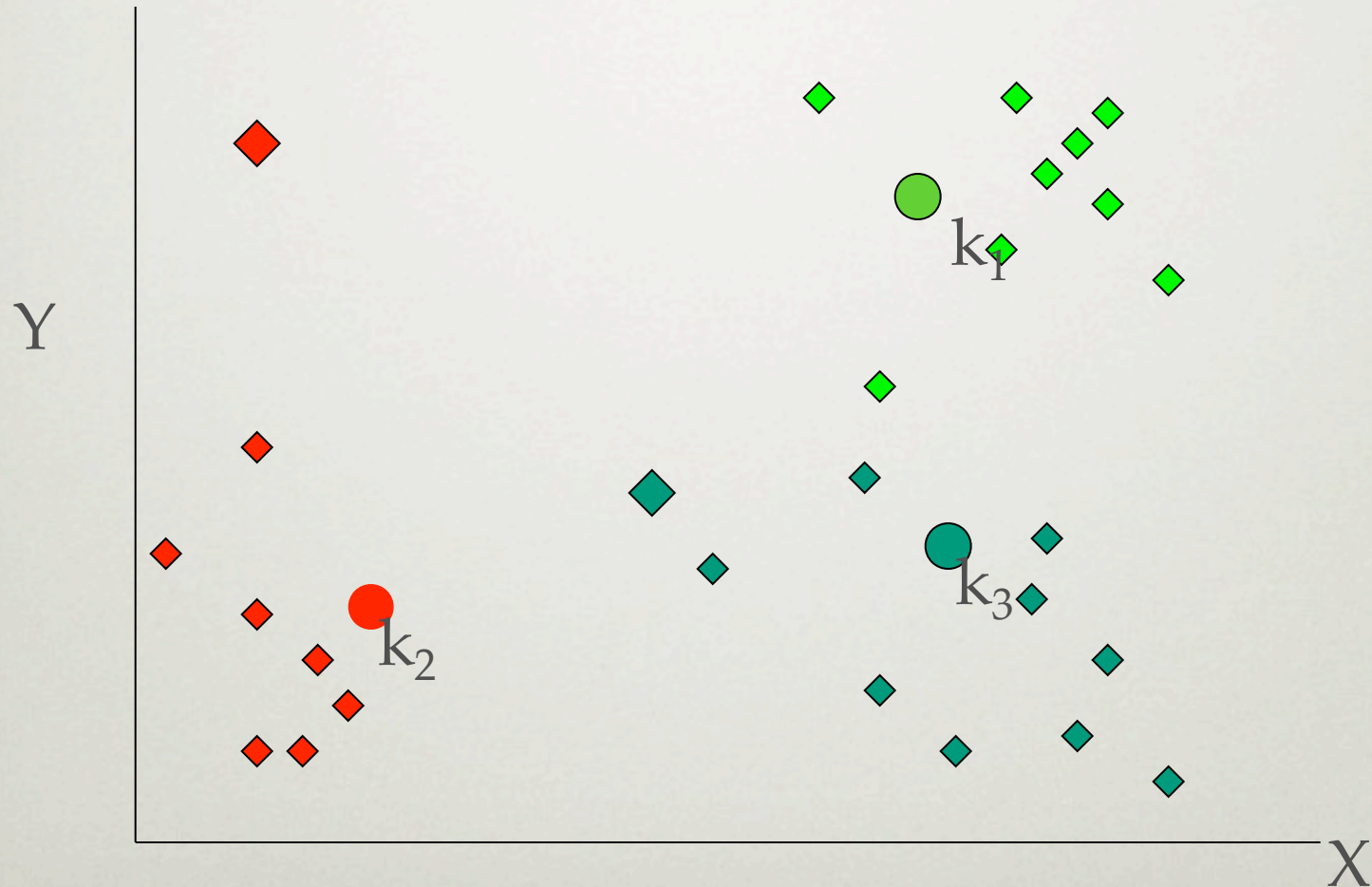
---



Reassign points to nearest cluster center

# K-MEANS CLUSTERING (K=3)

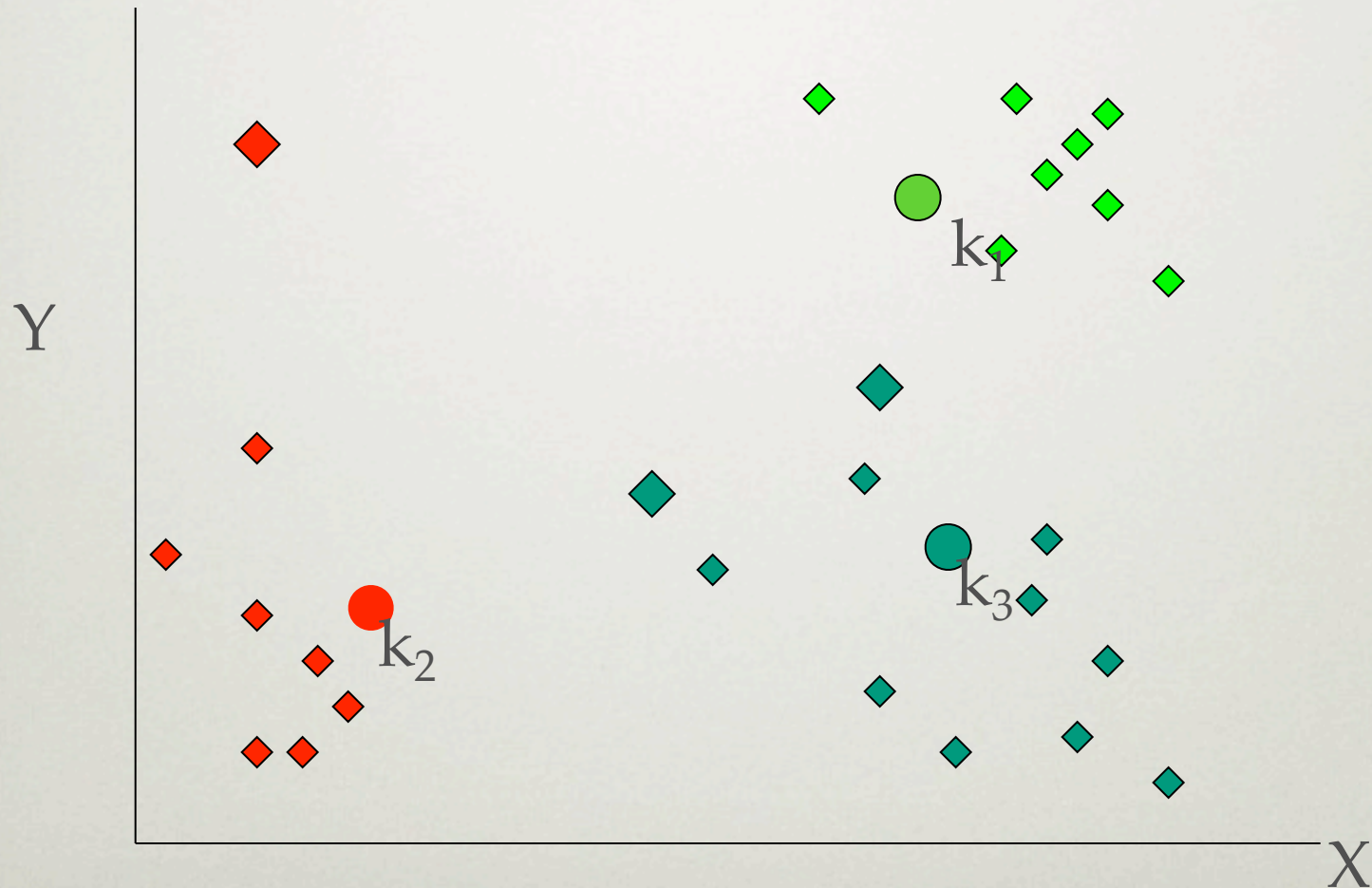
---



Reassign points to nearest cluster center

# K-MEANS CLUSTERING (K=3)

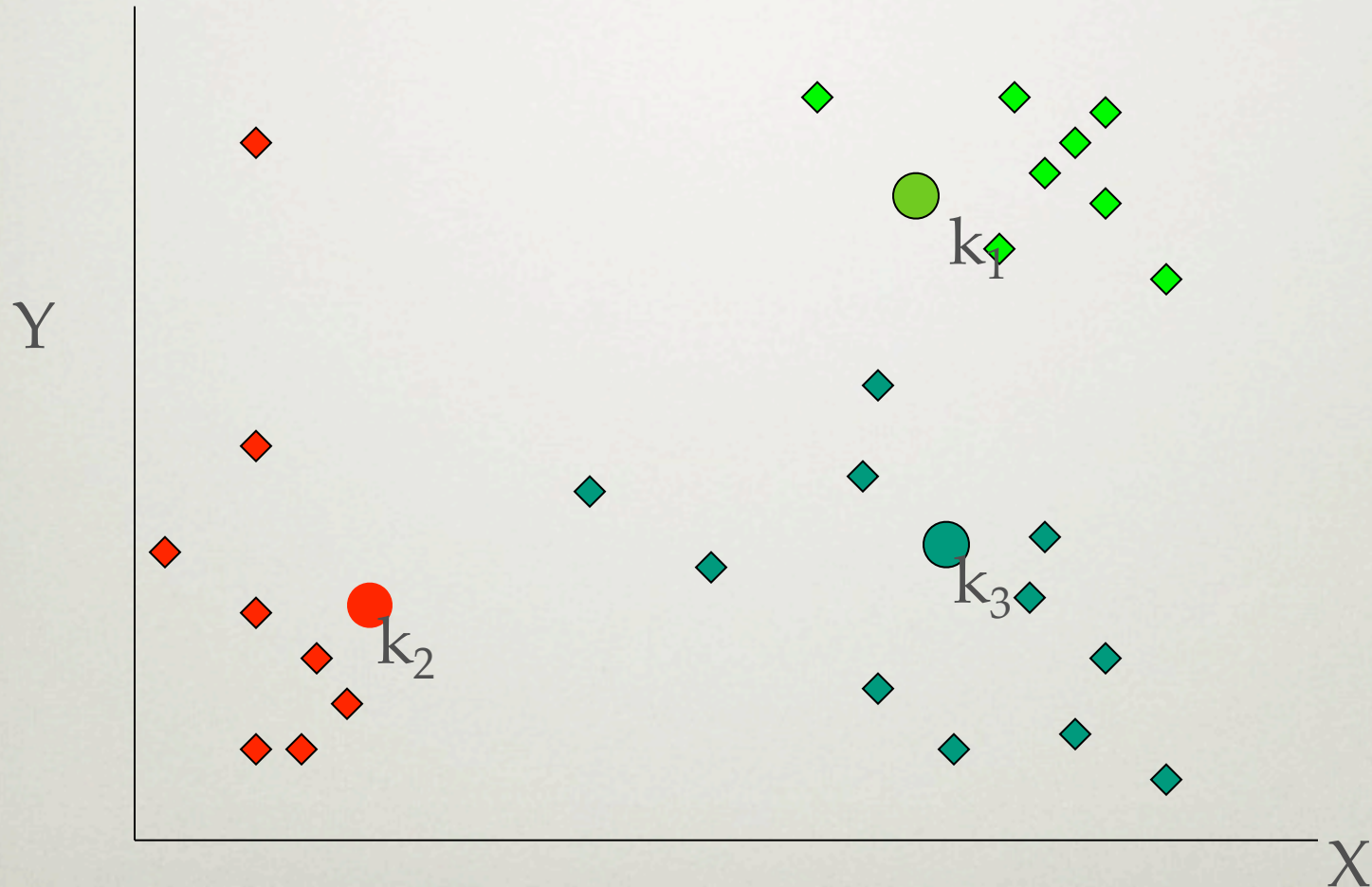
---



Reassign points to nearest cluster center

# K-MEANS CLUSTERING (K=3)

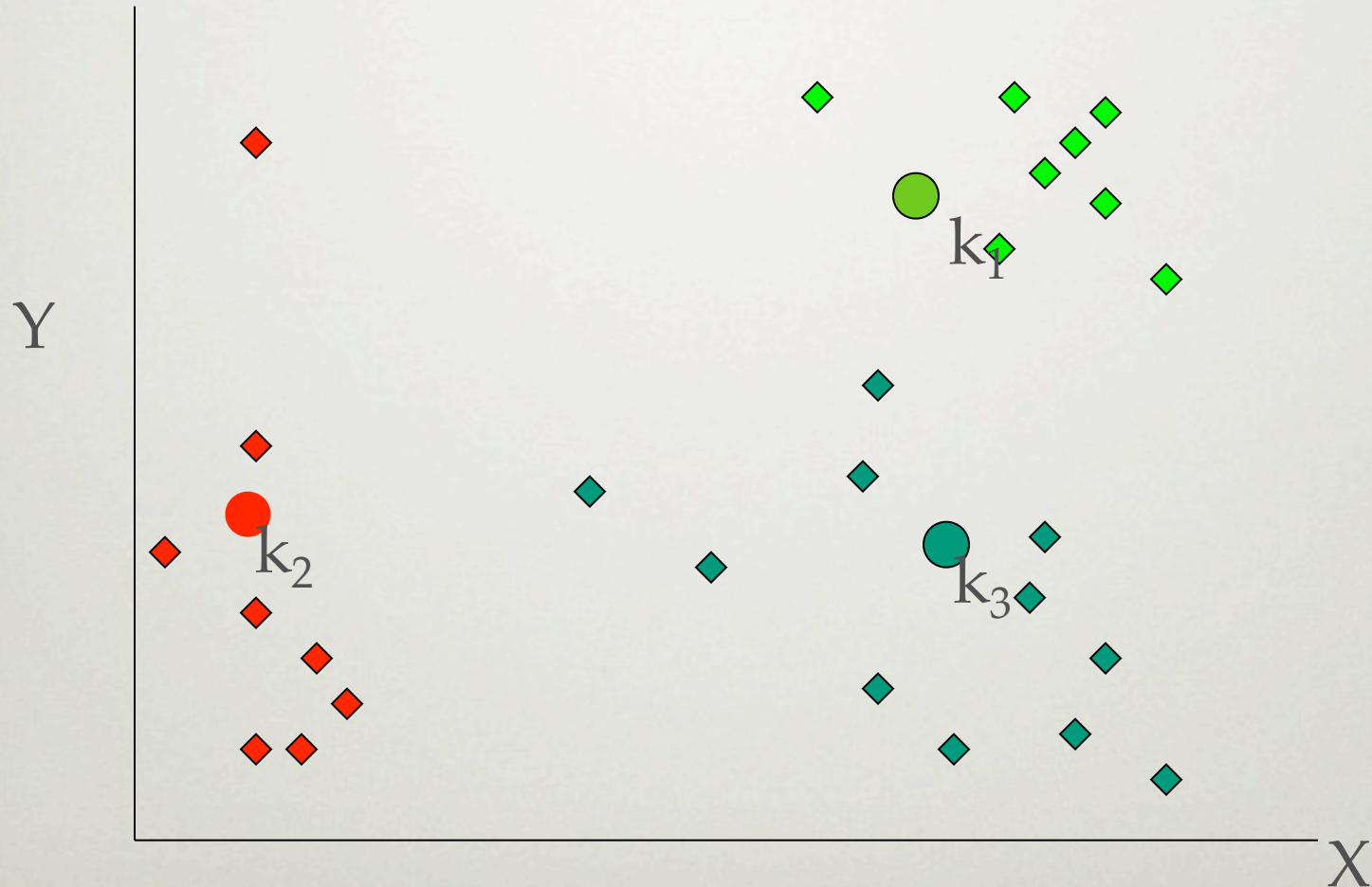
---



Repeat step 3-4 until cluster centers converge (don't/hardly move)

# K-MEANS CLUSTERING (K=3)

---

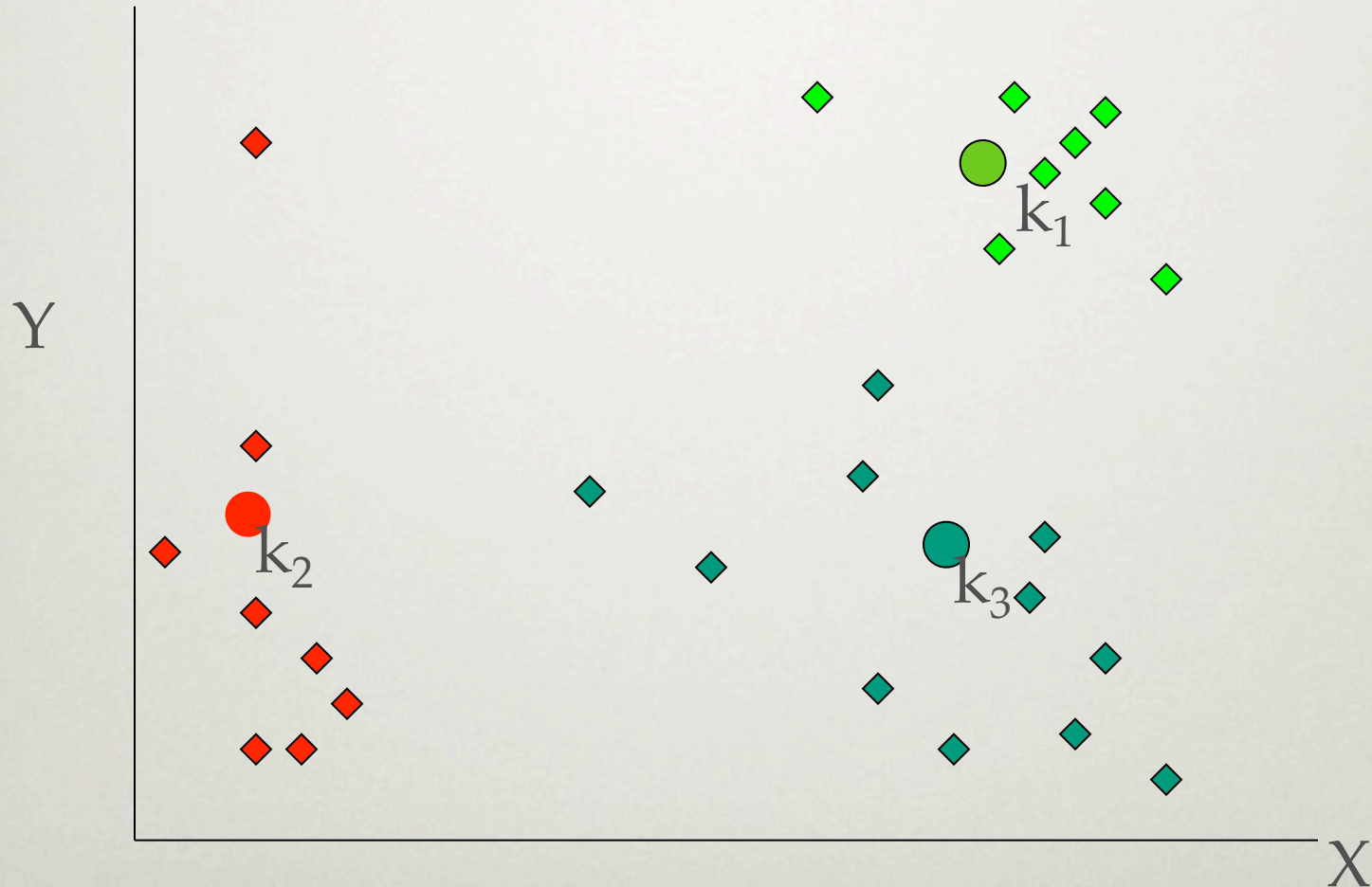


Repeat step 3-4 until cluster centers converge (don't/hardly move)



# K-MEANS CLUSTERING (K=3)

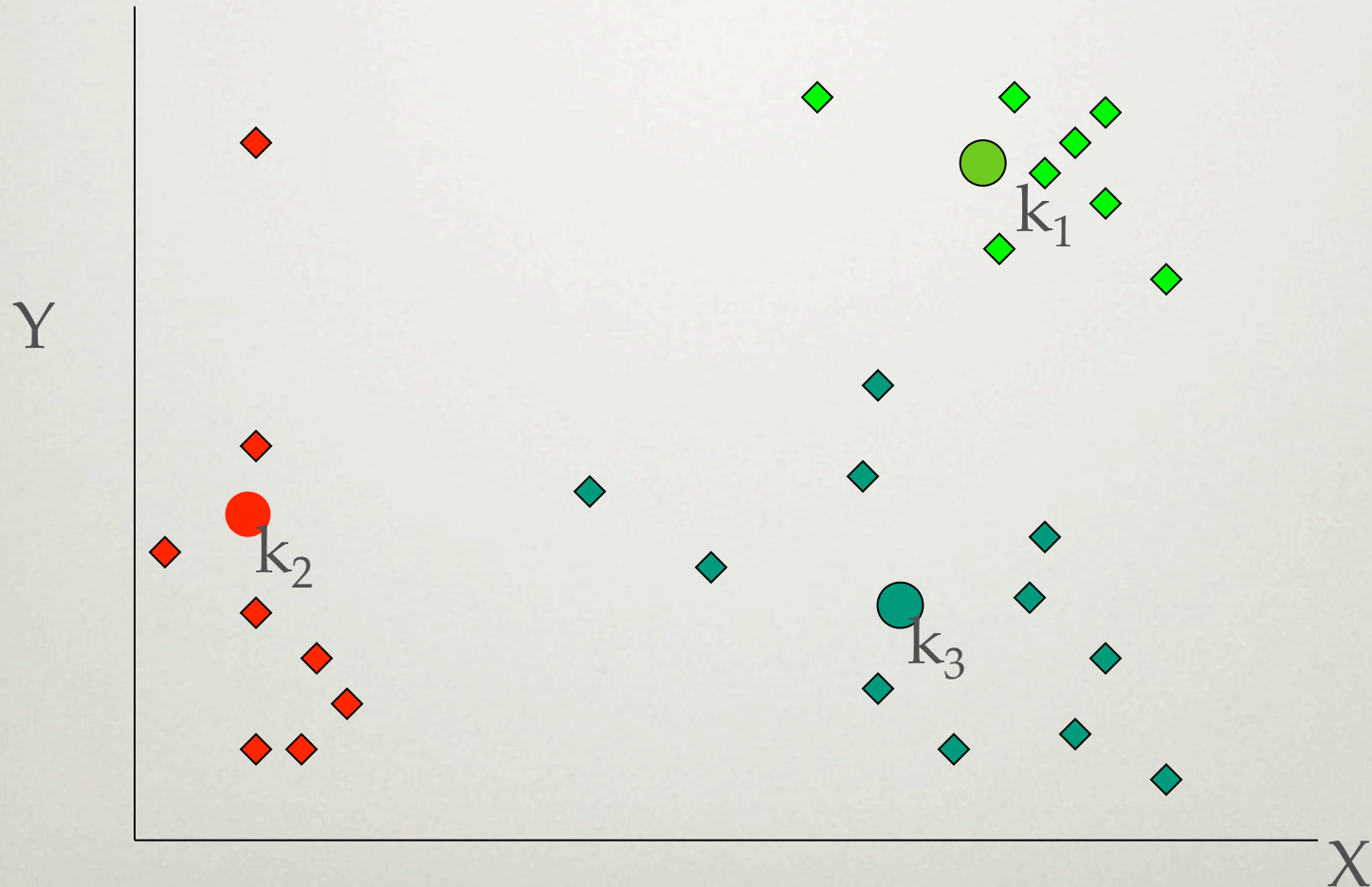
---



Repeat step 3-4 until cluster centers converge (don't/hardly move)

# K-MEANS CLUSTERING (K=3)

---



Repeat step 3-4 until cluster centers converge (don't/hardly move)

# K-MEANS

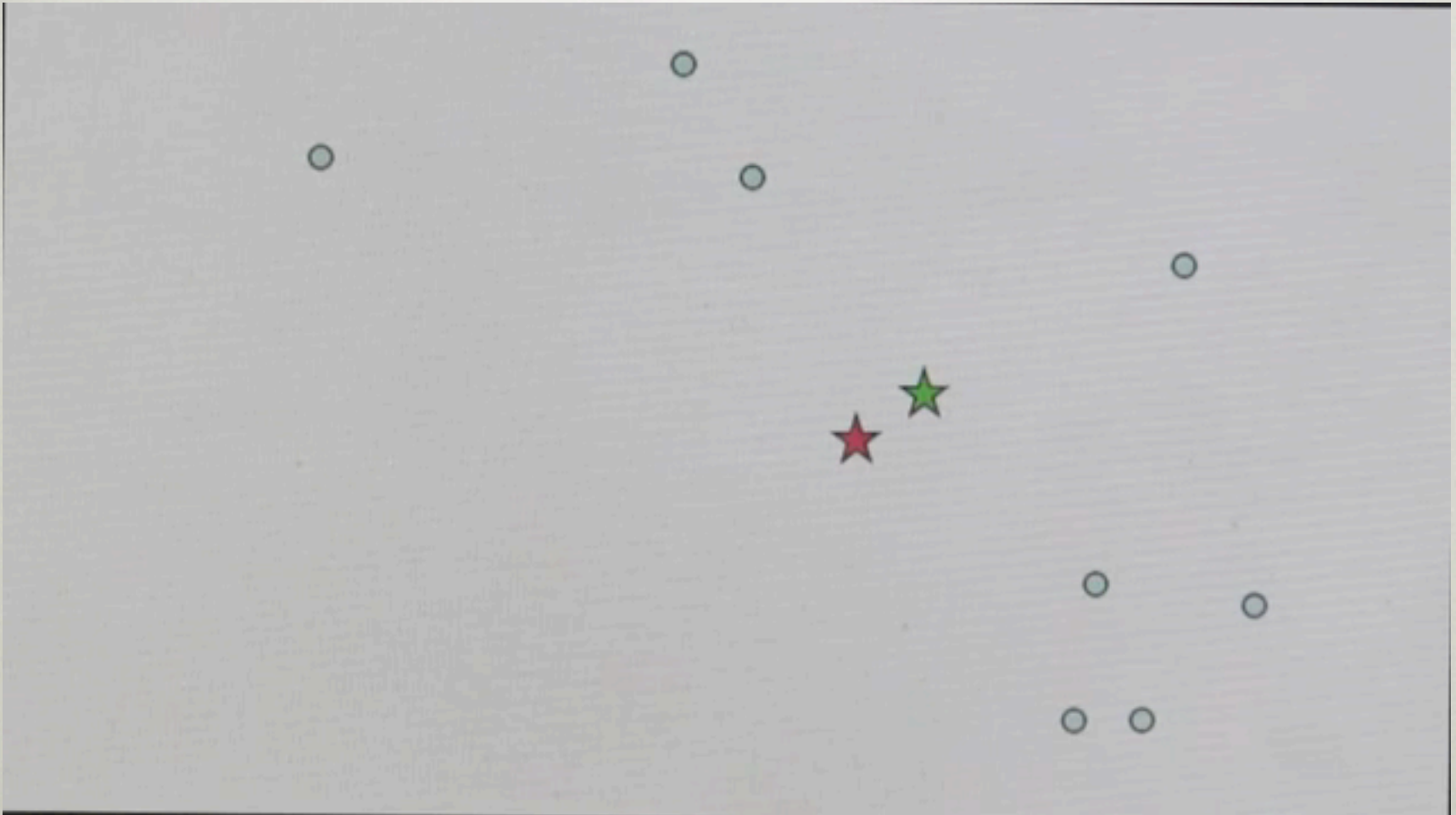
---

Works with numeric data only

- 1) Pick K random points: initial cluster centers
- 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 CLUSTERING: ANOTHER EXAMPLE

---

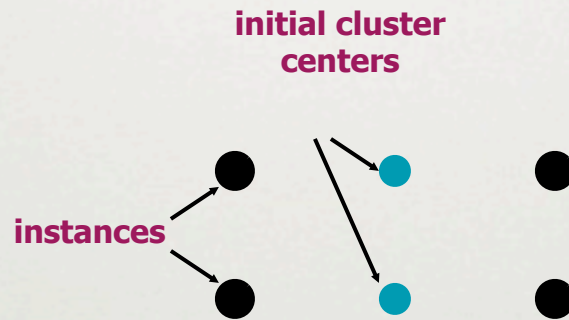


<http://www.youtube.com/watch?v=zaKjh2N8jN4#!>

# DISCUSSION

---

- Result can vary significantly depending on initial choice of centers
- Can get trapped in local minimum
  - Example:



- 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 single cluster
- Sensitive to outliers

# K-MEANS: VARIATIONS

---

- K-medoids – instead of mean, use medians of each cluster
  - Mean of 1, 3, 5, 7, 1009 is
  - Median of 1, 3, 5, 7, 1009 is
- For large databases, use sampling

# K-MEANS: VARIATIONS

---

- K-medoids – instead of mean, use medians of each cluster
  - Mean of 1, 3, 5, 7, 1009 is 205
  - Median of 1, 3, 5, 7, 1009 is
- For large databases, use sampling



# K-MEANS: VARIATIONS

---

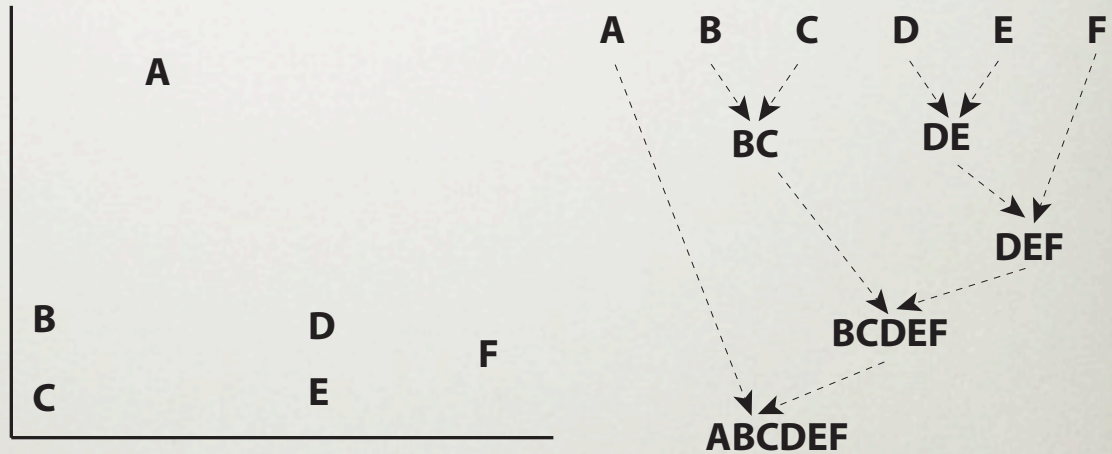
- K-medoids – instead of mean, use medians of each cluster
  - Mean of 1, 3, 5, 7, 1009 is 205
  - Median of 1, 3, 5, 7, 1009 is 5
- For large databases, use sampling

# HIERARCHICAL CLUSTERING

# BOTTOM-UP VS TOP-DOWN CLUSTERING

---

- Bottom up / Agglomerative
  - Start with single-instance clusters
  - At each step, join two “closest” clusters

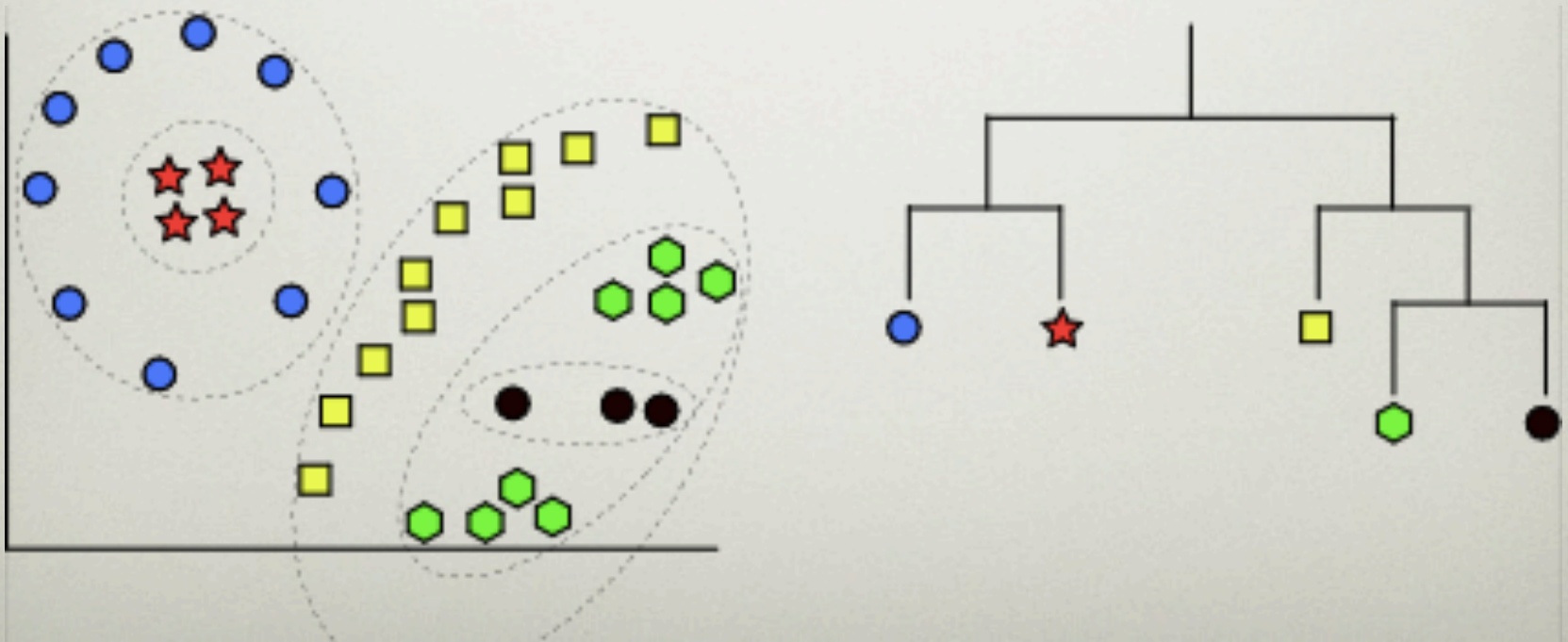


- Top down
  - Start with one universal cluster
  - Split in two clusters
  - Proceed recursively on each subset

# HIERARCHICAL CLUSTERING

---

- Hierarchical clustering represented in *dendrogram*
  - tree structure containing hierarchical clusters
  - clusters in leafs, union of child clusters in nodes

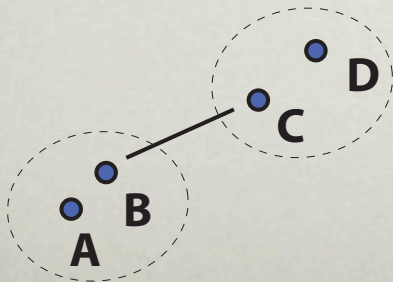


# DISTANCE BETWEEN CLUSTERS

- *Centroid*: distance between centroids
- Sometimes hard to compute (e.g. mean of molecules?)
- *Single Link*: smallest distance between points
- *Complete Link*: largest distance between points
- *Average Link*: average distance between points

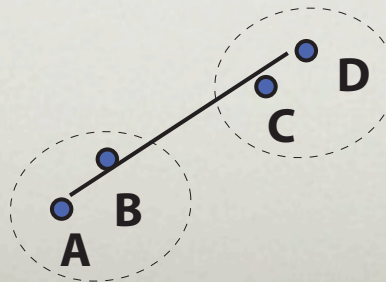
## single link

distance = 1



## complete link

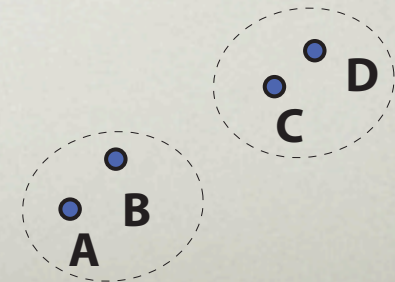
distance = 2



## average link

distance = 1.5

$$\frac{(d(A,C)+d(A,D)+d(B,C)+d(B,D))}{4}$$



# DISTANCE BETWEEN CLUSTERS

---

- *Centroid*: distance between centroids
  - Sometimes hard to compute (e.g. mean of molecules?)
- *Single Link*: smallest distance between points
- *Complete Link*: largest distance between points
- *Average Link*: average distance between points
- *Group-average*: group two clusters into one, then take average distance between all points (*incl.  $d(A,B)$  &  $d(C,D)$* )

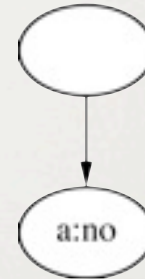
# INCREMENTAL CLUSTERING

# CLUSTERING WEATHER DATA

---

ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>False</b>
F	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

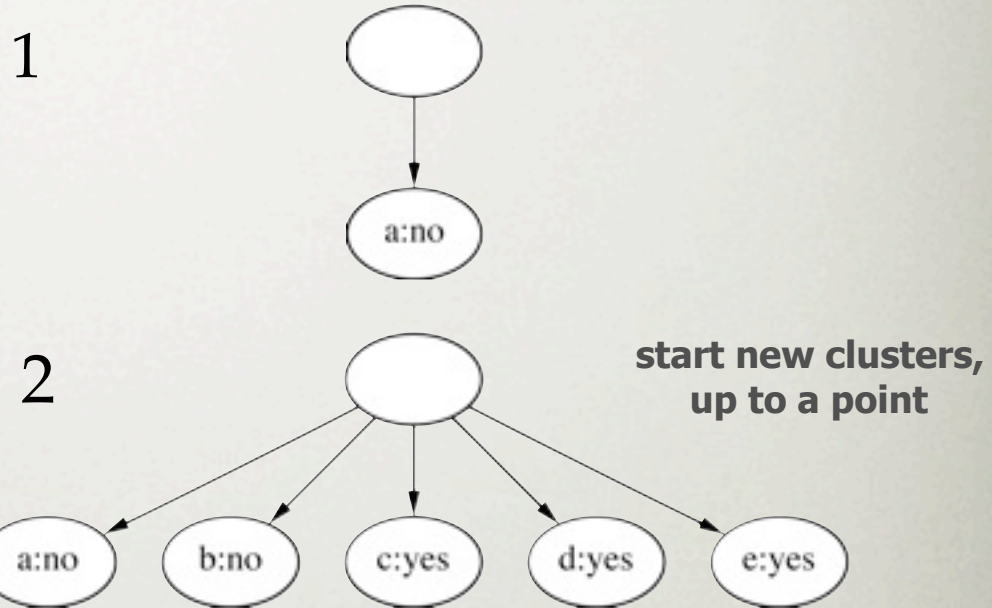
1





# CLUSTERING WEATHER DATA

ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	Rainy	Cool	Normal	False
F	Rainy	Cool	Normal	True
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True



# CATEGORY UTILITY

---

- Category utility: overall quality of clustering
- Quadratic loss function
  - nominal: clusters  $C_i$ , attributes  $a_i$ , values  $v_{ij}$ :

$$CU(C_1, C_2, \dots, C_k) = \frac{\sum_l \Pr[C_l] \sum_i \sum_j (\Pr[a_i = v_{ij} | C_l]^2 - \Pr[a_i = v_{ij}]^2)}{k}$$

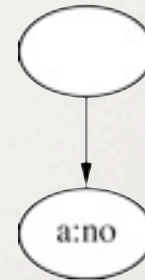
- numeric: similar, assume Gaussian distribution
- Intuitively:
  - good clusters allow to predict value of new data points:  
 $\Pr[a_i=v_{ij} | C_i] > \Pr[a_i=v_{ij}]$
  - $1/k$  factor: penalty for using many clusters (avoids overfitting)

# CLUSTERING WEATHER DATA

---

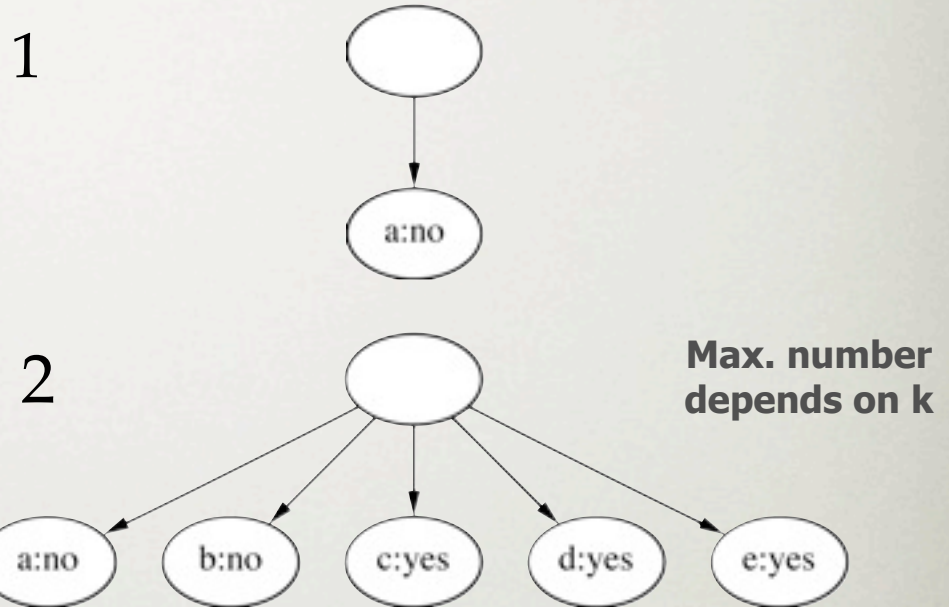
ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>False</b>
F	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

1



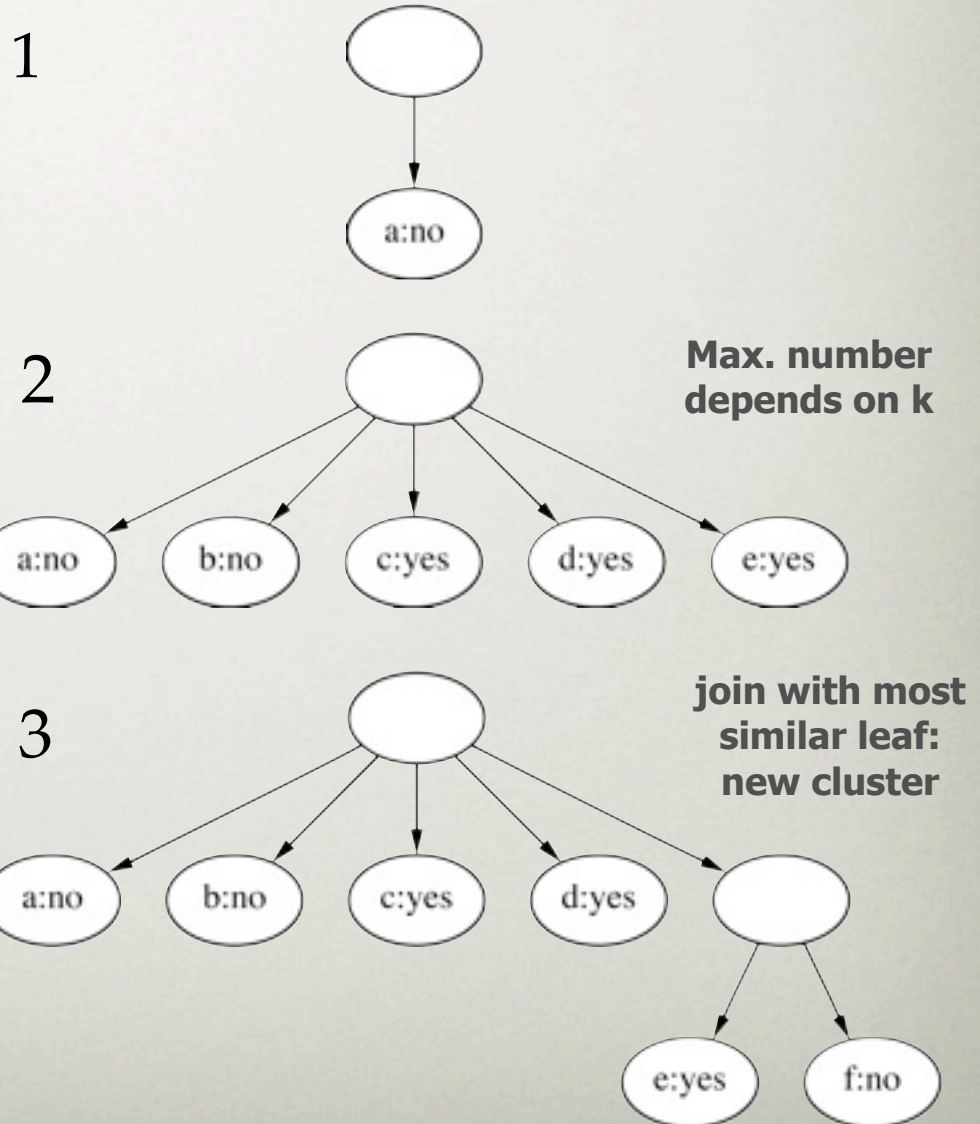
# CLUSTERING WEATHER DATA

ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	Rainy	Cool	Normal	False
F	Rainy	Cool	Normal	True
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True



# CLUSTERING WEATHER DATA

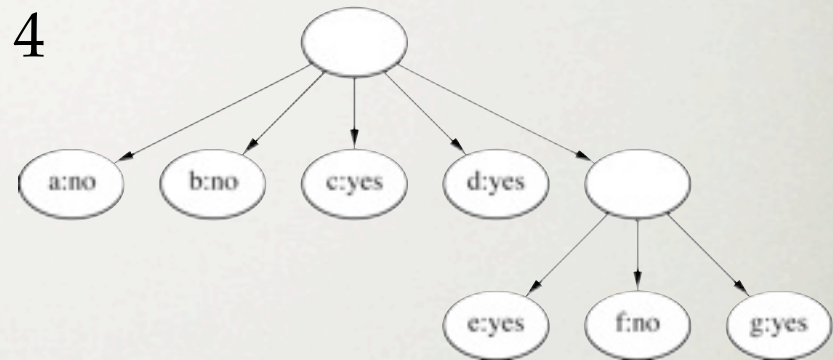
ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	Rainy	Cool	Normal	False
F	Rainy	Cool	Normal	True
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True



# CLUSTERING WEATHER DATA

ID	Outlook	Temp.	Humidity	Windy
A	<i>Sunny</i>	<i>Hot</i>	<i>High</i>	<i>False</i>
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	<i>Rainy</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
E	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>False</b>
F	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
G	<b>Overcast</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
H	<i>Sunny</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

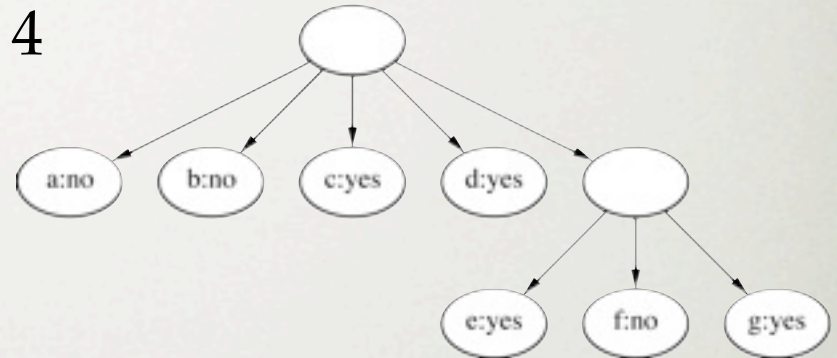
4



# CLUSTERING WEATHER DATA

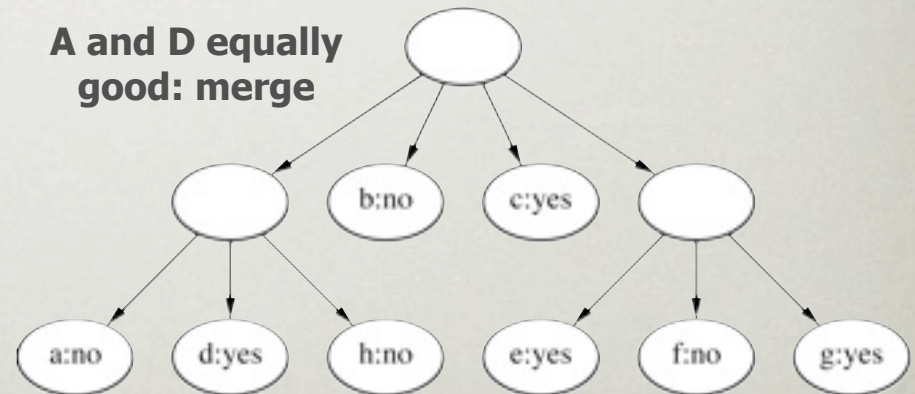
ID	Outlook	Temp.	Humidity	Windy
A	<i>Sunny</i>	<i>Hot</i>	<i>High</i>	<i>False</i>
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	<i>Rainy</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
E	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>False</b>
F	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
G	<b>Overcast</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
H	<i>Sunny</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

4



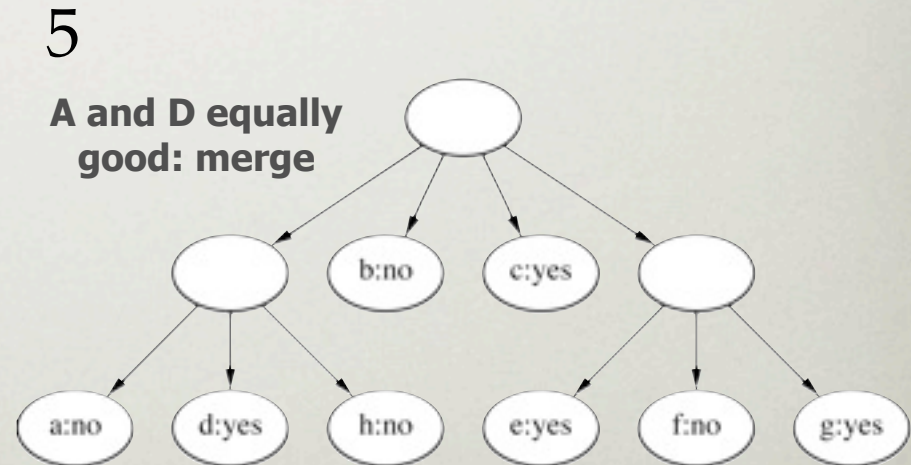
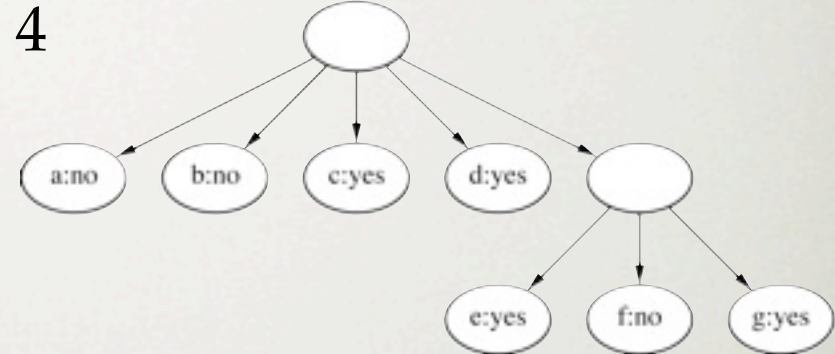
5

**A and D equally good: merge**



# CLUSTERING WEATHER DATA

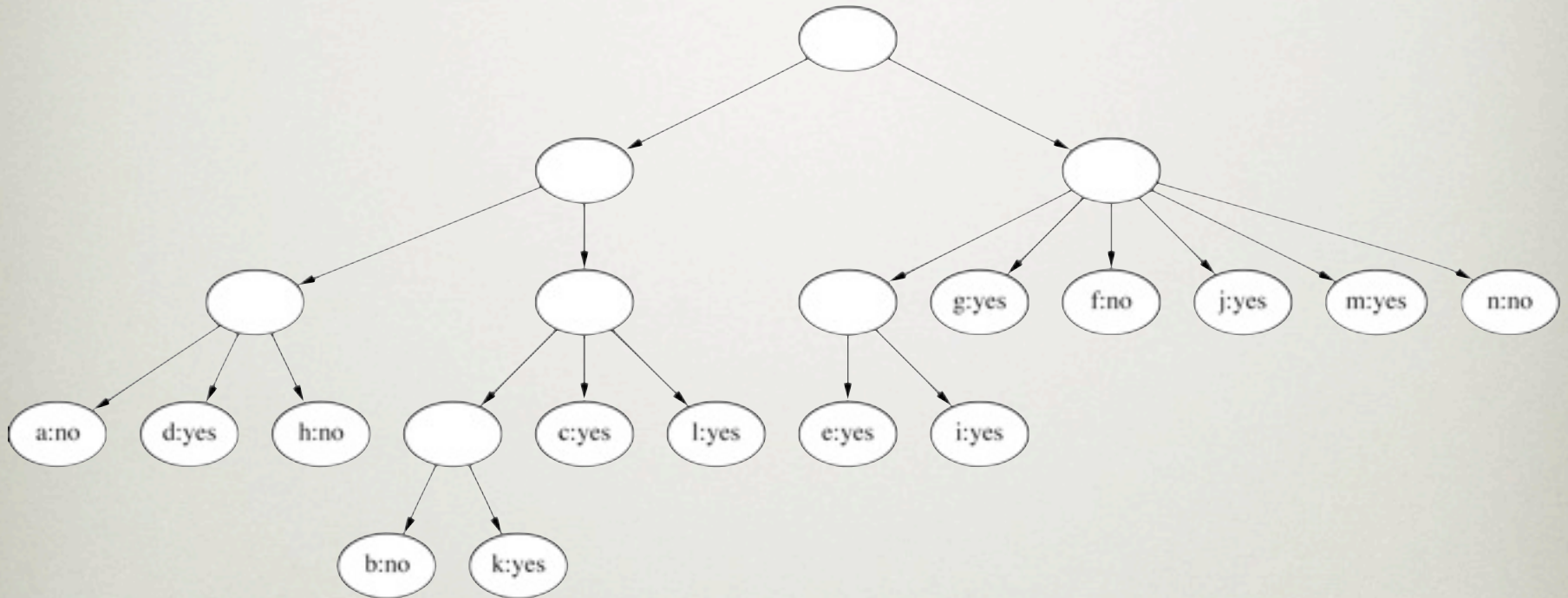
ID	Outlook	Temp.	Humidity	Windy
A	<i>Sunny</i>	<i>Hot</i>	<i>High</i>	<i>False</i>
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	<i>Rainy</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
E	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>False</b>
F	<b>Rainy</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
G	<b>Overcast</b>	<b>Cool</b>	<b>Normal</b>	<b>True</b>
H	<i>Sunny</i>	<i>Mild</i>	<i>High</i>	<i>False</i>
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True



**Consider splitting the best host if merging doesn't help**



# FINAL HIERARCHY



ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False

**Note that a and b are actually very similar, but end up in different clusters**

# INCREMENTAL CLUSTERING

---

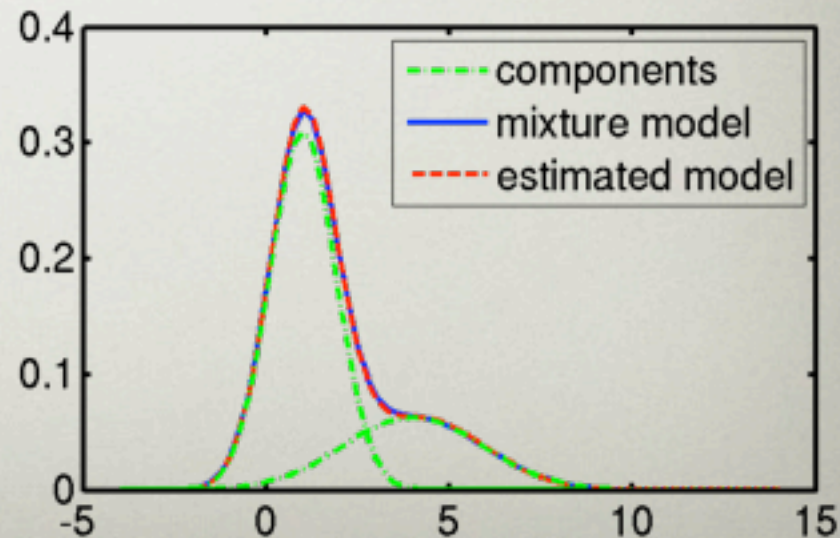
- For large, regularly updated databases
  - start with tree and empty root node
  - add instances one by one
  - update tree appropriately at each stage
    - form new leaf
    - join instance with most similar leaf: new node (cluster)
    - *merge* existing leafs (move down one level)
    - *split* node into leafs (move up one level)
- Best decision: *category utility*

# PROBABILITY-BASED CLUSTERING

# PROBABILITY-BASED CLUSTERING

---

- Given  $k$  clusters, each instance belongs to *all* clusters, with a certain probability
  - *mixture model: set of  $k$  distributions (one per cluster)*
  - *also: each cluster has prior likelihood*
- If correct clustering known, we know parameters  $\mu, \sigma$  and  $P(C_i)$  for each cluster: calculate  $P(C_i | x)$  using Bayes' rule
- Estimate the unknown parameters
  - How?



# EM

## EXPECTATION MAXIMIZATION

---

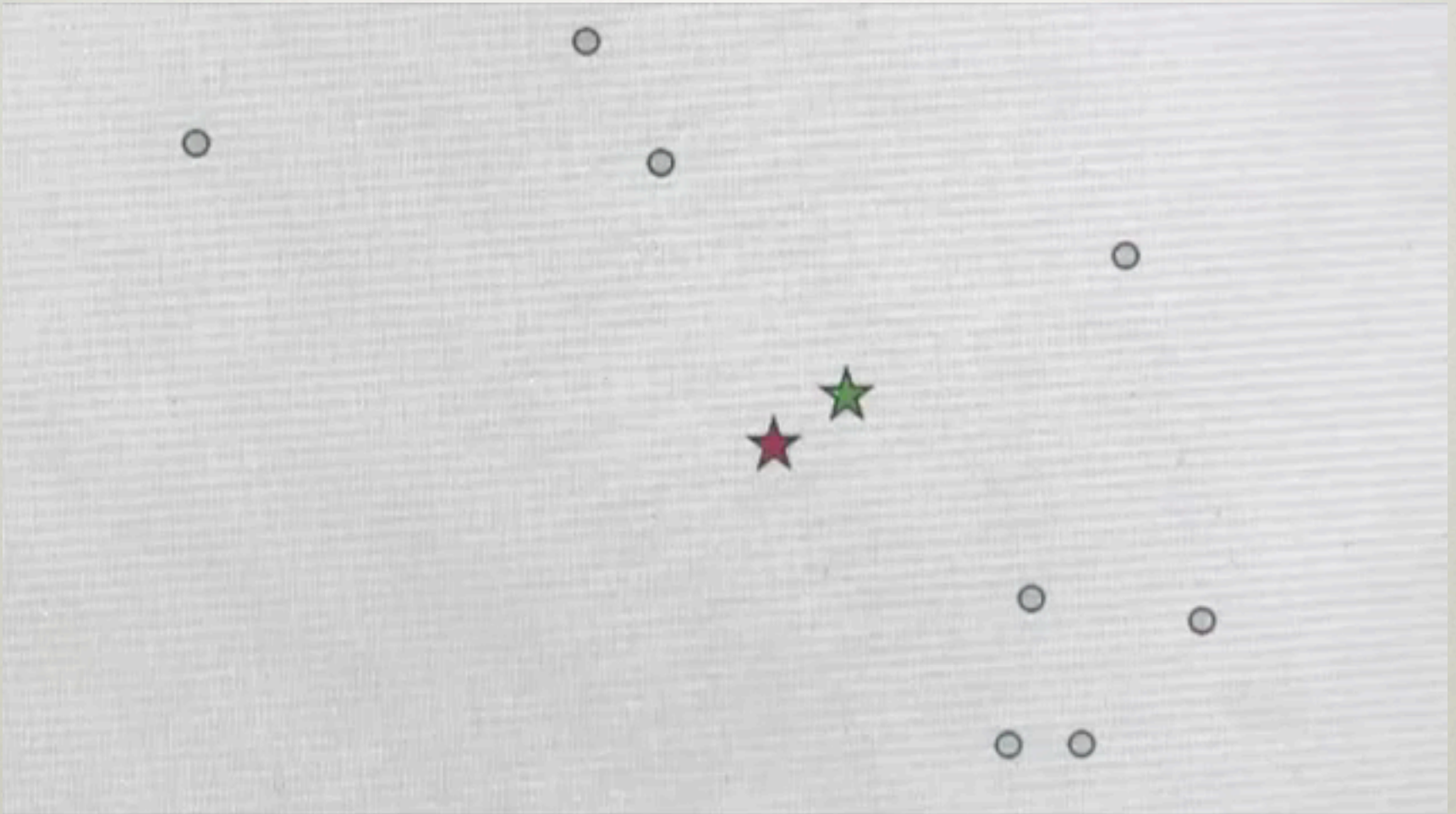
- Finds the parameters  $\mu, \sigma$  for the distributions and the cluster membership

$$\mu_A = \frac{w_1 x_1 + \dots + w_n x_n}{w_1 + \dots + w_n}$$

- (Random) initialization
  - Initial parameters  $\mu, \sigma, P(C_i)$  for each cluster
- Iterative algorithm:
  - Expectation step: with current parameters, calculate  $P(C | x)$
  - Maximization step: update parameters using  $P(C | x)$ : new  $\mu, \sigma, P(C_i)$
- Iterate until converged to local optimum

# EM VS K-MEANS

---



<http://www.youtube.com/watch?v=1CWDWmF0i2s>

# QUIZ

EM Quiz

$\mu, \Sigma$   
is  $\Sigma$

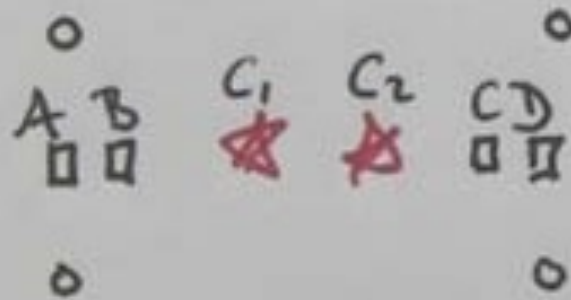
- circular
- elongated



# QUIZ

Quiz

EM versus K-Means



EM

$C_1 \rightarrow A$

$C_1 \rightarrow B$

K MEANS

$C_1 \rightarrow A$

$C_1 \rightarrow B$



# CLUSTERING EVALUATION

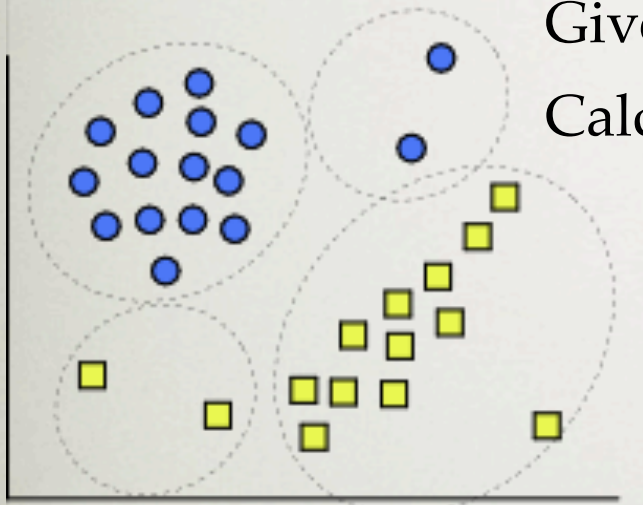
---

- Manual inspection
- Benchmarking on existing labels
- Cluster quality measures
  - distance measures
  - high similarity within a cluster, low across clusters

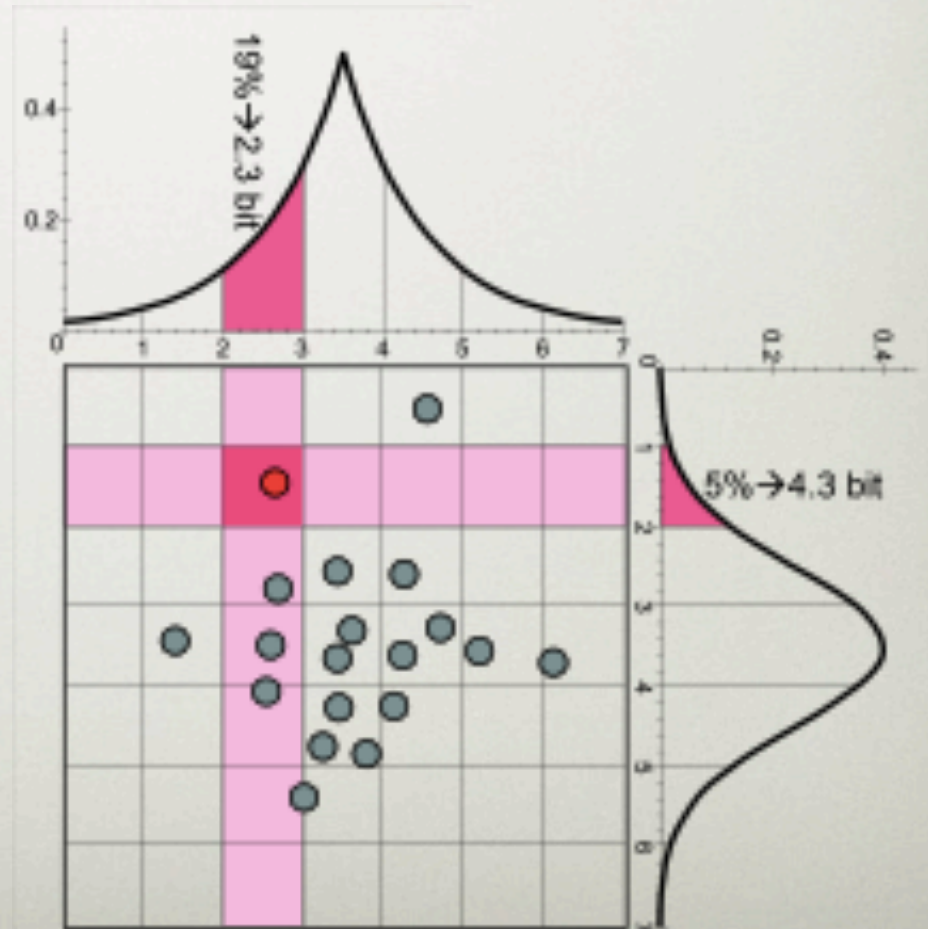
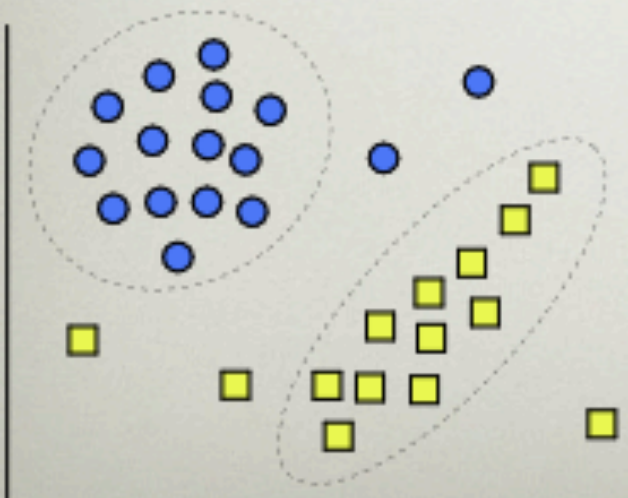
# GOODNESS OF FIT

Given a function that defines the cluster

Calculate for each point how well it fits the cluster



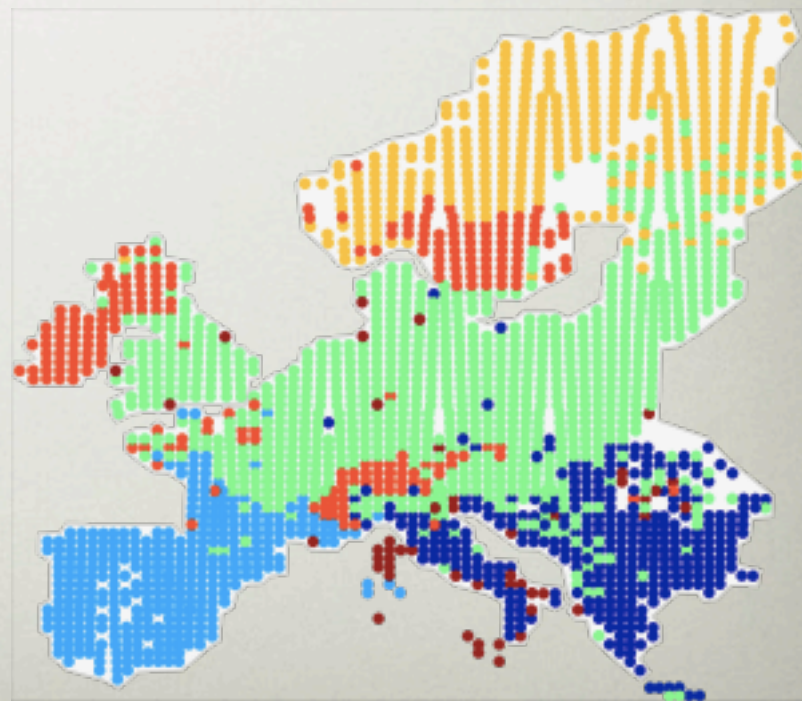
..., or better?



# HOW TO CHOOSE K?

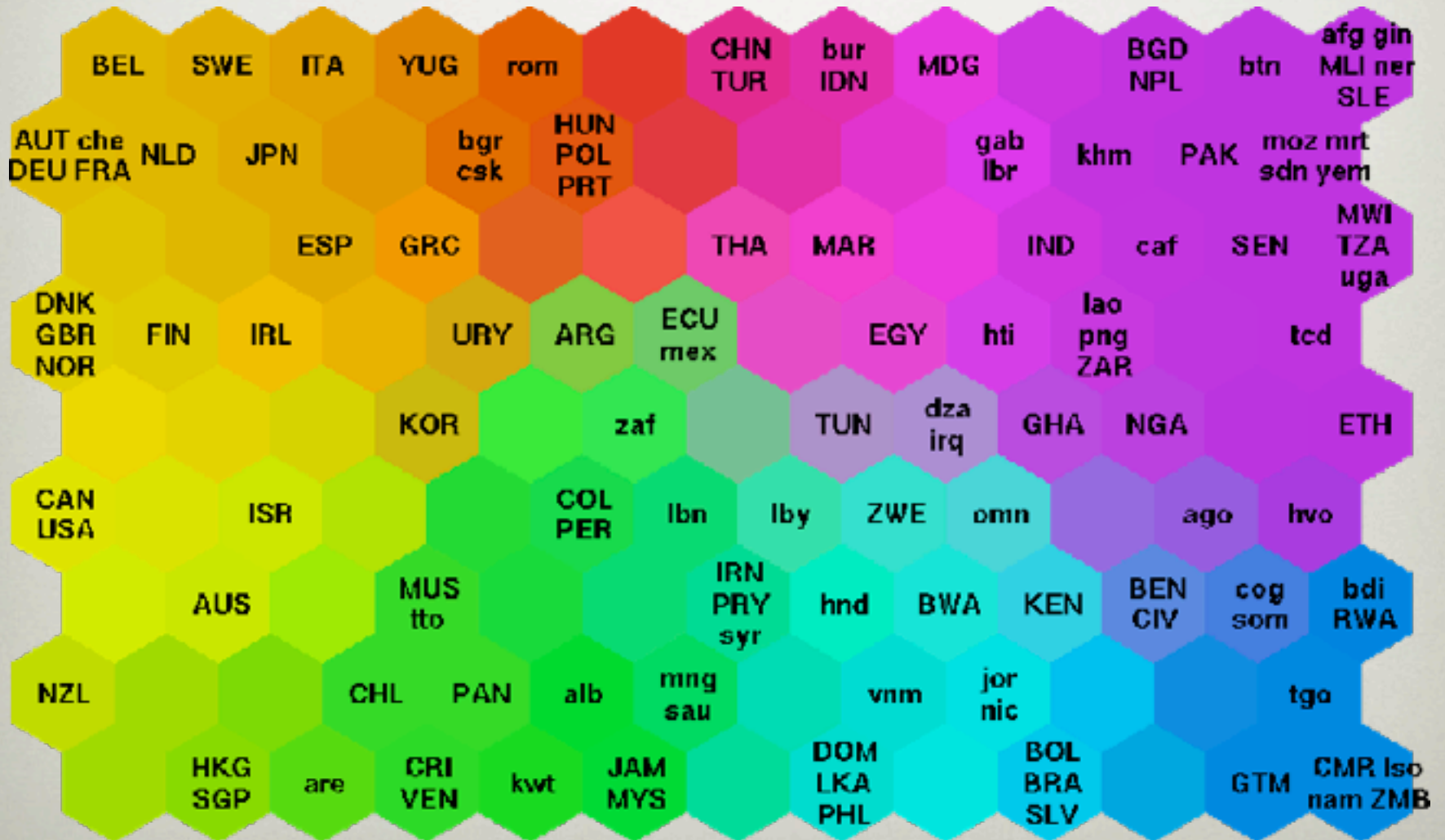
---

- One important parameter  $k$ , but how to choose?
  - Domain dependent, we simply want  $k$  clusters
- Alternative: repeat for several values of  $k$  and choose the best
- Example:
  - cluster mammal properties
  - $k$  different clusters
  - Use an MDL based encoding
  - Alternative to category theory
  - Each additional cluster introduces a penalty
  - Optimal for  $k=6$



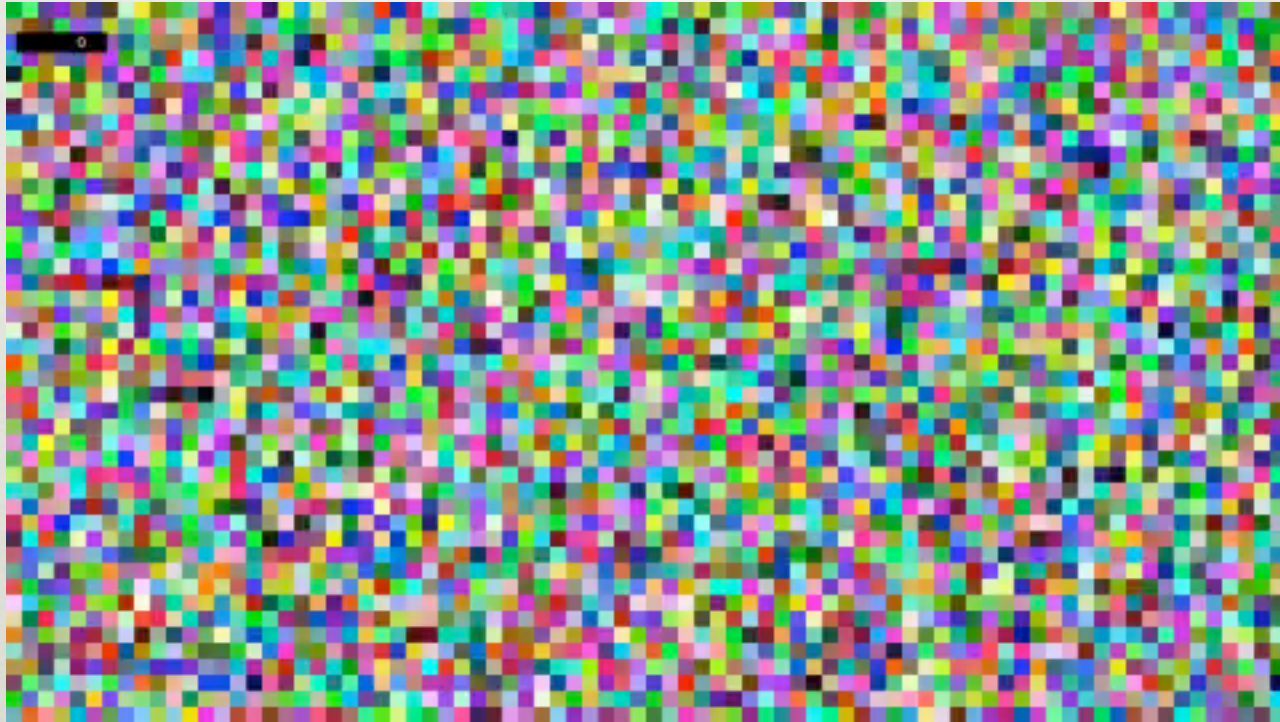
# SELF-ORGANIZING MAPS

# SELF ORGANIZING MAP



# SELF ORGANIZING MAP

---



<http://www.youtube.com/watch?v=71wmOT4lHWc>

# SELF ORGANIZING MAP

---

- Applications
  - Group similar data together
  - Dimensionality reduction
  - Data visualization technique
- Similar to neural networks
  - Neurons try to mimic the input vectors
  - The winning neuron (and its neighborhood) wins
  - Topology preserving, using Neighborhood function

# SOM LEARNING ALGORITHM

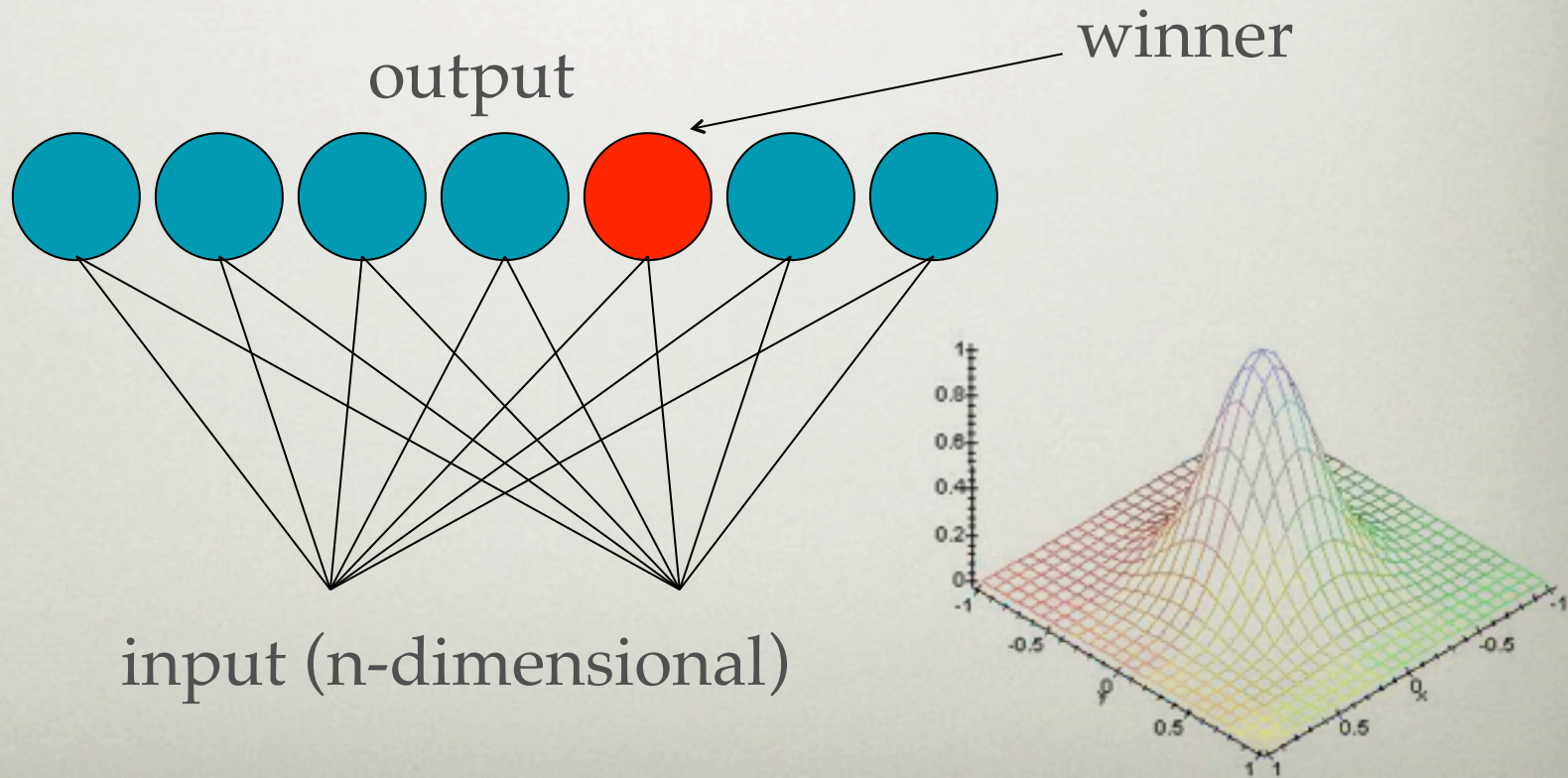
---

- Initialize SOM (random, or such that dissimilar input is mapped far apart)
- For t from 0 to N
  - Randomly select a training instance
  - Get the best matching neuron
    - calculate distance, e.g.  $\sqrt{\sum_{i=0}^n (x_i - w_i)^2}$
  - Scale neighbors
    - Who? decrease over time: Hexagons, squares, gaussian, ...
    - Update of neighbours towards the training instance



# SELF ORGANIZING MAP

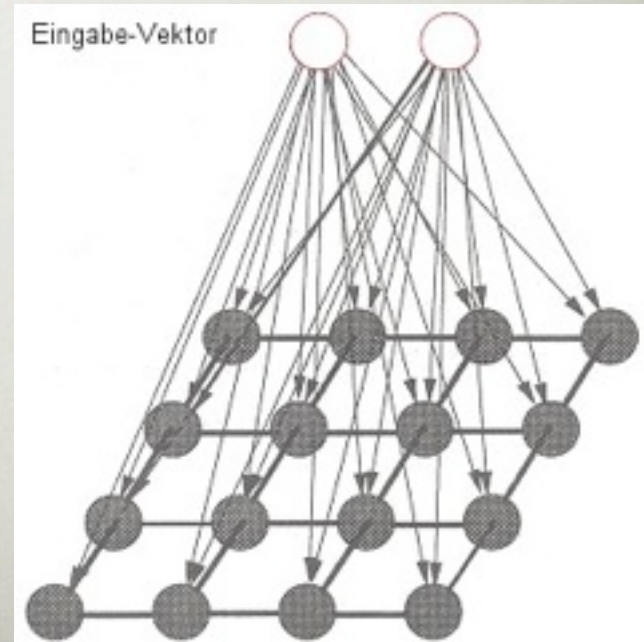
- Neighborhood function to preserve topological properties of the input space
- Neighbors share the prize (but slightly less)



# SELF ORGANIZING MAP

---

- Input: uniformly randomly distributed points
- Output: Map of  $20^2$  neurons
- Training: Starting with a large learning rate and neighborhood size, both are gradually decreased to facilitate convergence
- After learning, neurons with similar weights tend to cluster on the map



# DISCUSSION

---

- Can interpret clusters by using supervised learning
  - learn a classifier based on clusters
- Decrease dependence between attributes?
  - pre-processing step
  - E.g. use *principal component analysis*
- Can be used to fill in missing values
- Key advantage of probabilistic clustering:
  - Can estimate likelihood of data
  - Use it to compare different models objectively