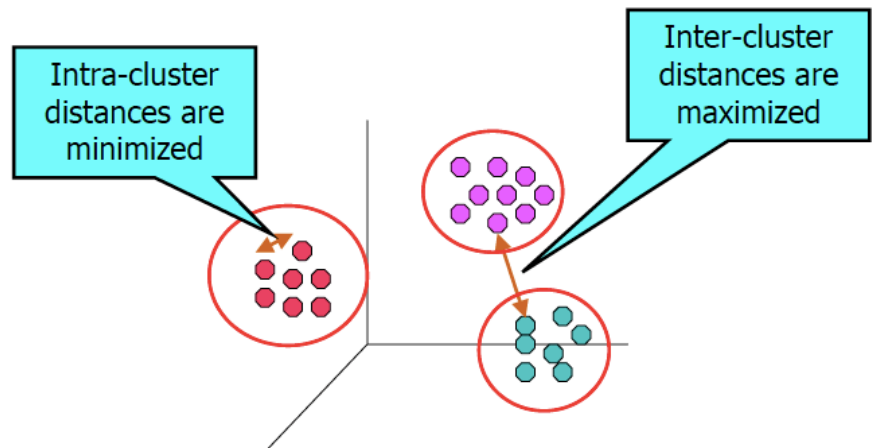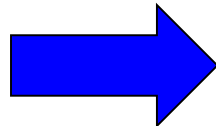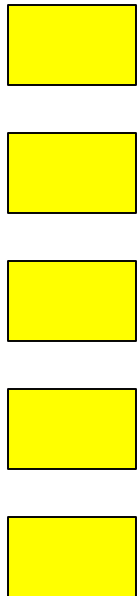# CLUSTERING

# The Problem of Clustering

- Clustering is a technique for finding similarity groups in data, called clusters.
  - Groups data instances that are
    - Similar to (near) each other in one cluster, and
    - Very different (far away) from each other into different clusters.

- Clustering is an <u>unsupervised learning task</u> as no class values denoting an a priori grouping of the data instances are given, which is the case in supervised learning.
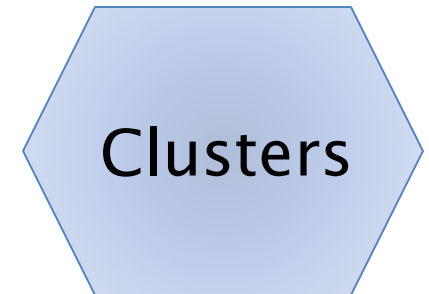
# **Un**supervised learning: clustering

Raw data

Features

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

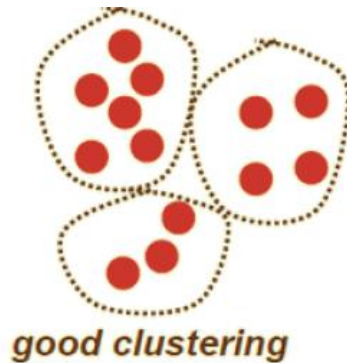extract features

group into classes/clusters

Clusters

**No** "supervision", we're only given data and want to find natural groupings

# What do we need for clustering?

- Proximity measure, either
  - Similarity measure $s(x_1, x_2)$ - large if $x_1, x_2$ are similar
  - Dissimilarity (or distance) measure $d(x_1, x_2)$ – small if $x_1, x_2$ are similar

- Criterion function to evaluate a clustering



good clustering          bad clustering

- Algorithm to compute clustering
  - For example, by optimizing the criterion function

# Distance (dissimilarity) measures

- Euclidean distance

  - $d(x_i, x_j) = \sqrt{\sum_{k=1}^{d}(x_i^{(k)} - x_i^{(k)})^2}$

- Manhattan (city block) distance

  - Approximation to Euclidean distance, cheaper to compute

  - $d(x_i, x_j) = \sum_{k=1}^{d} |x_i^{(k)} - x_j^{(k)}|$

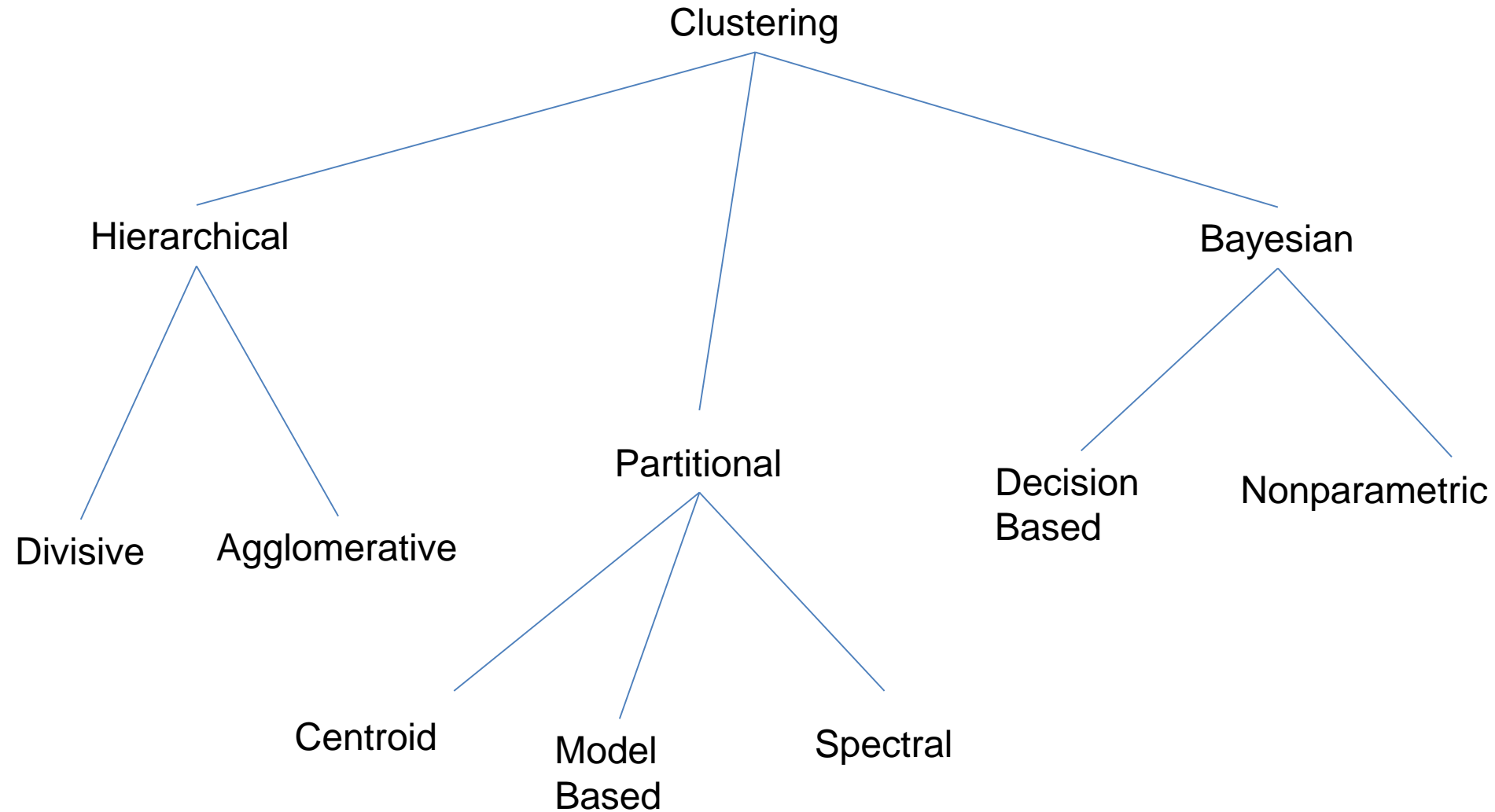- Minkowski distance (there are special cases)

  - $\grave{} d_p(x_i, x_j) = (\sum_{k=1}^{m} |x_{ik} - x_{jk}|^p)^{\frac{1}{p}}$
  - P is a positive integer

# Cluster evaluation (a hard problem)

- **Intra-cluster cohesion**(compactness):
  - Cohesion measures how near the data points in a cluster are to the cluster centroid.
  - Sum of squared error (SSE) is a commonly used measure.

- **Inter-cluster separation**(isolation):
  - Separation means that different cluster centroids should be far away from one another.

- In most applications, expert judgments are still the key

# Clustering Techniques
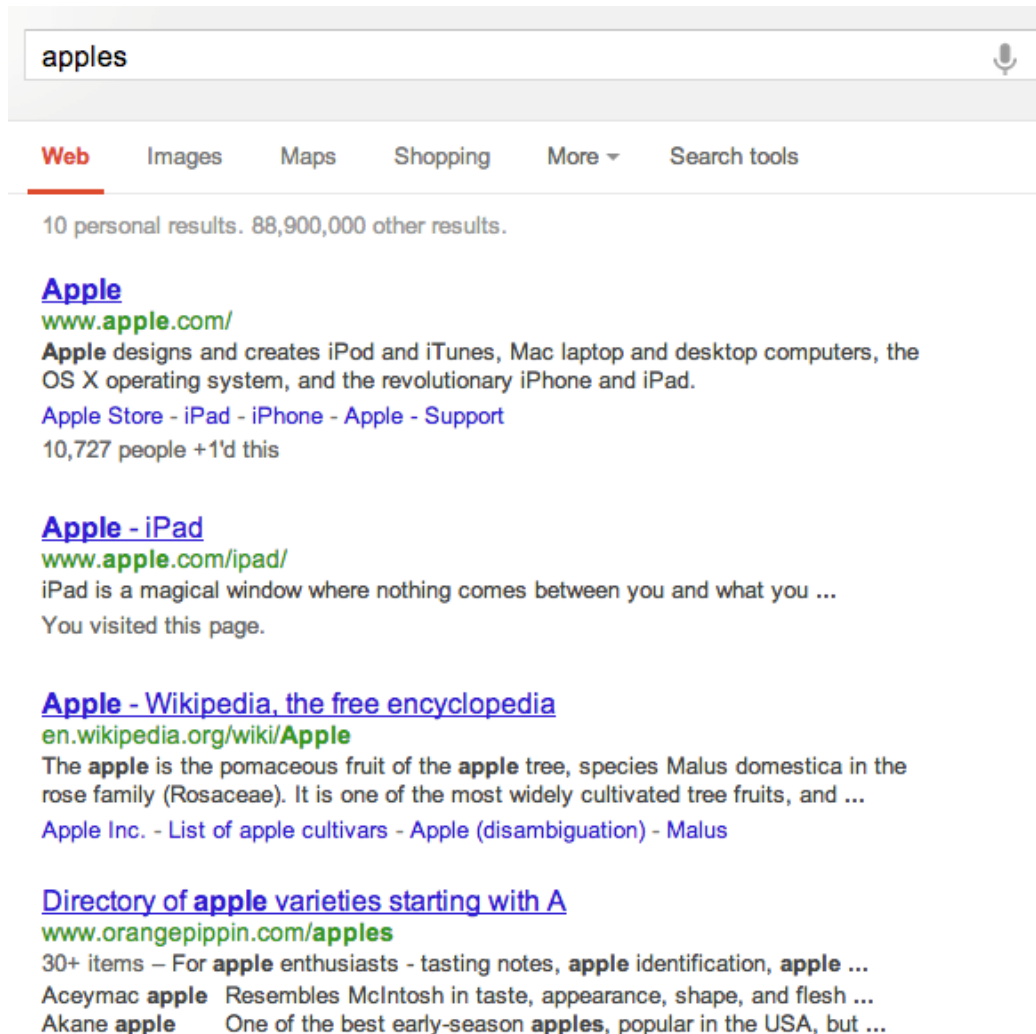
# Applications of clustering

- *Understanding*
  - Grouping objects into conceptually meaningful classes is an important step in analysis.
  - EX: Group related documents for browsing, group genes and proteins that have similar functionality, or group stocks with similar price fluctuations.

- *Summarization*
  - Reduce the size of large data sets
  - EX: PCA requires space complexity $O(m^2)$ so may not be practice for large datasets. We can cluster dataset and apply PCA on the cluster prototypes.

# Applications (clustering for IR)



Documents or webpages in the same cluster are likely to be similar.

# Application – Clustering for image search



Given a collection of unlabeled objects(images), cluster them into groups.
So that given an image, we call pull similar images from the right cluster.

# Gene expression data



Data from Garber et al. PNAS (98), 2001.

# Applications - Face Clustering

# Applications - Face clustering

# Clustering is a hard problem!

# Why is it hard?

- Clustering in two dimensions looks easy
- Clustering small amounts of data looks easy
- And in most cases, looks are not deceiving

- Many applications involve not 2, but 10 or 10,000 dimensions
- High-dimensional spaces look different: Almost all pairs of points are at about the same distance

# Overview: Methods of Clustering

- **Hierarchical:**
  - **Agglomerative** (bottom up):
    - Initially, each point is a cluster
    - Repeatedly combine the two "nearest" clusters into one
  - **Divisive** (top down):
    - Start with one cluster and recursively split it
- **Partitional**
  - Usually start with a random (partial) partitioning
  - Refine it iteratively
    - *K* means clustering
    - Model based clustering
- **Bayesian**
  - Try to generate a posteriori distribution over the collection of all partitions of the data.

# Hard vs. soft clustering

***Hard clustering:*** Each example belongs to exactly one cluster

***Soft clustering:*** An example can belong to more than one cluster (probabilistic)

- Makes more sense for applications like creating browsable hierarchies
- You may want to put a pair of sneakers in two clusters: (i) sports apparel and (ii) shoes

# Hierarchical Clustering

- Key operation: Repeatedly combine two nearest clusters

- Three important questions:
  - 1) How do you represent a cluster of more than one point?
  - 2) How do you determine the "nearness" of clusters?
  - 3) When to stop combining clusters?

# Hierarchical Clustering

- Key operation: Repeatedly combine two nearest clusters

(1) How to represent a cluster of many points?

  - Key problem: As you merge clusters, how do you represent the "location" of each cluster, to tell which pair of clusters is closest
  - Euclidean case: each cluster has a  centroid = average of its (data)points

(2) How to determine "nearness" of clusters?

  - Measure cluster distances by distances of centroids

# Example: Hierarchical clustering



(1,2)

o

x (1.5,1.5)

x (1,1)    o (2,1)

o (0,0)

(5,3)

o

x (4.7,1.3)

o (4,1)

x (4.5,0.5)

o (5,0)

**Data:**
o … data point
x … centroid

**Dendrogram**

# And in the Non-Euclidean Case?

- What about the Non-Euclidean case?
  - The only "locations" we can talk about are the points themselves
    - i.e., there is no "average" of two points

- Approach 1:
  - (1) How to represent a cluster of many points?
    clustroid  = (data)point "closest" to other points

  - (2) How do you determine the "nearness" of clusters? Treat clustroid as if it were centroid, when computing inter-cluster distances

# "Closest" Point?

- (1) How to represent a cluster of many points?
  clustroid  = point "closest" to other points

- Possible meanings of "closest":
  - Smallest maximum distance to other points
  - Smallest average distance to other points
  - Smallest sum of squares of distances to other points
    - For distance metric d clustroid c of cluster C is:

$$\min_{c} \sum_{x \in C} d(x,c)^2$$

Datapoint

**Centroid**

X

**Clustroid**

Cluster on
3 datapoints

**Centroid** is the avg. of all (data)points in the cluster. This means centroid is an "artificial" point.
**Clustroid** is an **existing** (data)point that is "closest" to all other points in the cluster.

# Implementation

- Naïve implementation of hierarchical clustering:
  - At each step, compute pairwise distances between all pairs of clusters, then merge
  - $O(N^3)$

- Careful implementation using priority queue can reduce time to $O(N^2 \log N)$
  - Still too expensive for really big datasets that do not fit in memory

# *K*-MEANS CLUSTERING

Demo of K-means clustering :

http://www.onmyphd.com/?p=kmeans.clustering&ckattempt=2

# K-means

- Given *k*, the *k-means* algorithm works as follows:
1. Choose *k*(random) data points (seeds) to be the initial centroids, cluster centers
2. Assign each data point to the closest centroid
3. Re-compute the centroids using the current cluster memberships
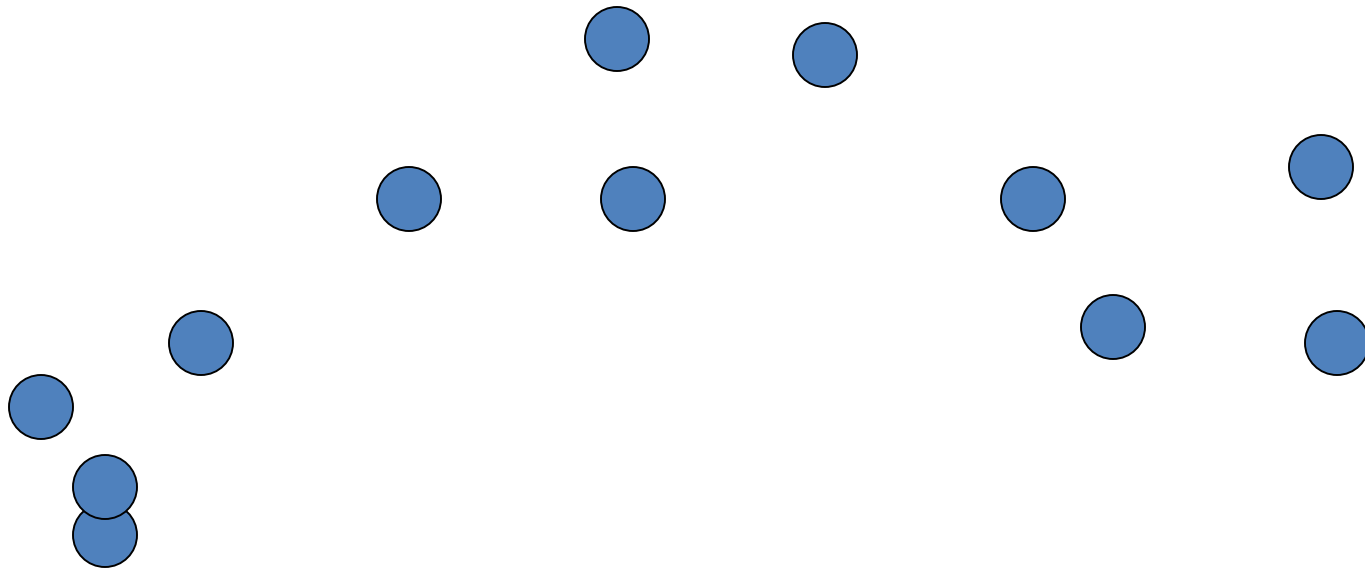4. If a convergence criterion is not met, repeat steps 2 and 3

# K-means convergence (stopping) criterion

- No (or minimum) re-assignments of data points to different clusters, *or*
- No (or minimum) change of centroids, or
- Minimum decrease in the **sum of squared error**(SSE),

  - $SSE = \sum_{j=1}^{k} \sum_{x \in C_j} d(x, m_j)^2$

  - $C_j$ is the $j$th cluster,

  - **m**$_j$ is the centroid of cluster $C_j$ (the mean vector of all the data points in $C_j$),

  - $d$(**x**, **m**$_j$) is the (Euclidian) distance between data point **x** and centroid **m**$_j$
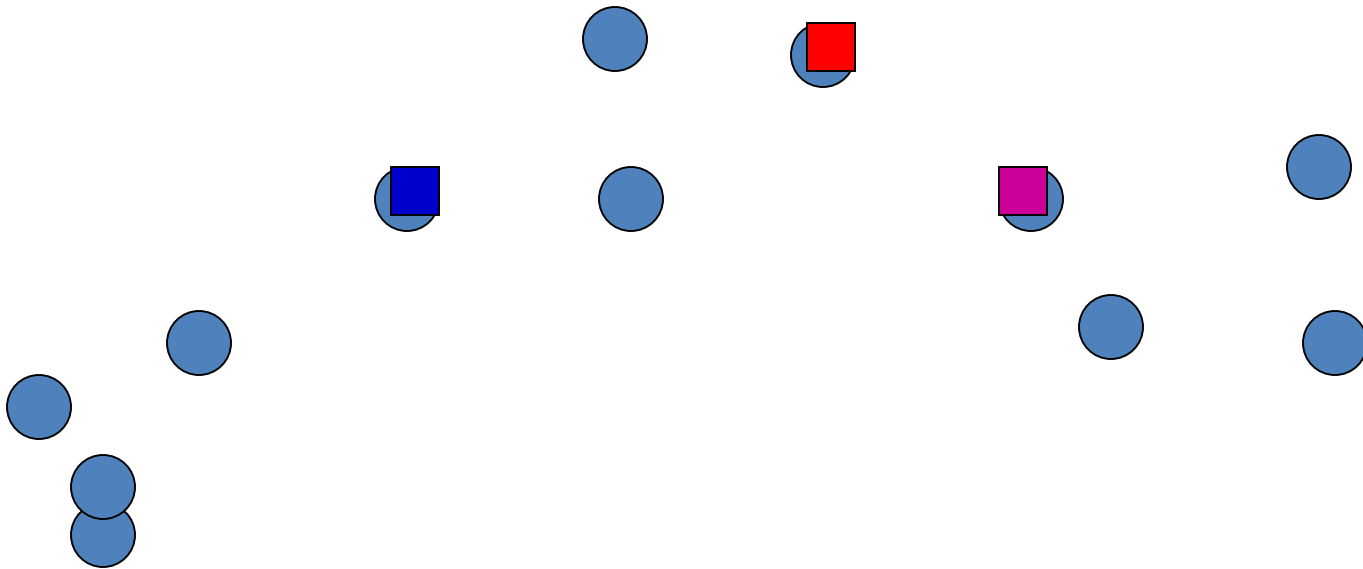
# K-means: an example

# Randomly Initialize the cluster centers

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster
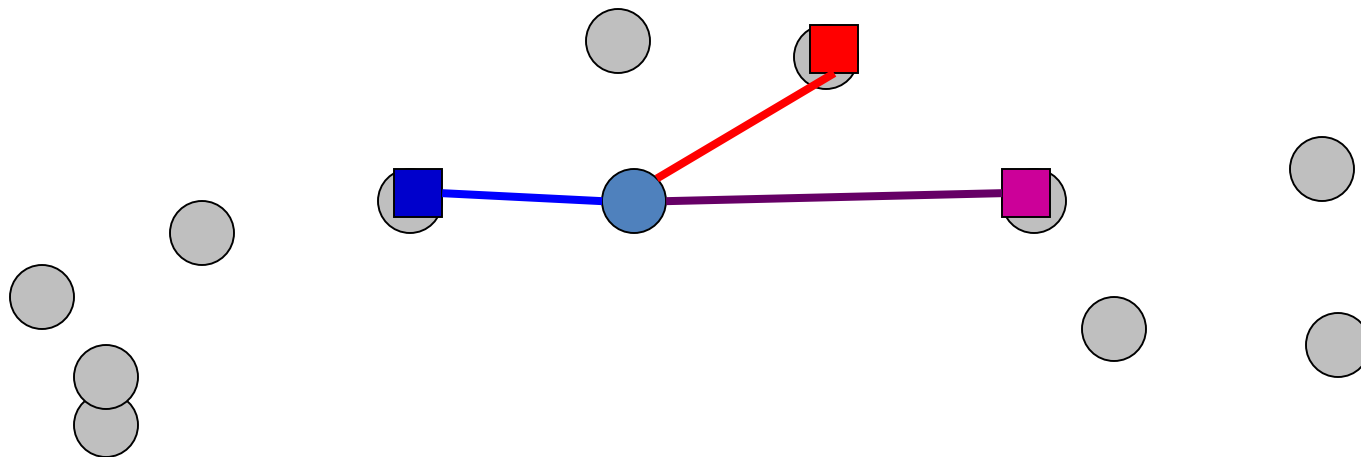
# Iteration # 1
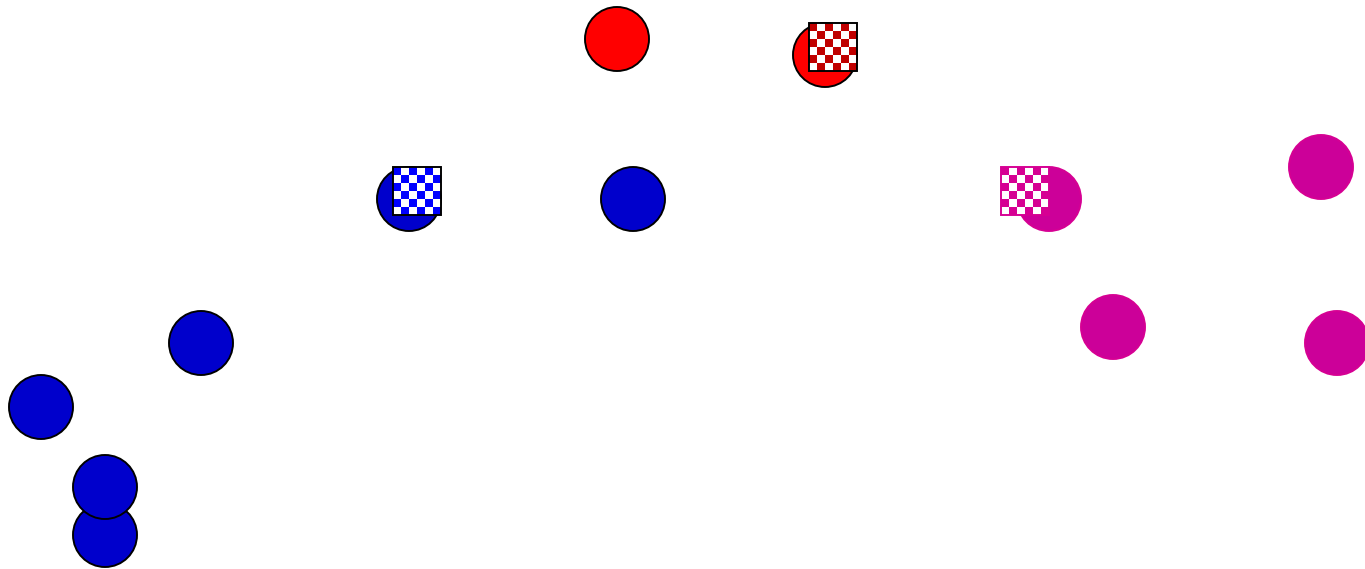
- **Assign each item to closest cluster center**

   iterate over each point:
   - get distance to each cluster center
   - assign to closest center (hard cluster)

- Recalculate centers as the mean of the points in a cluster

# After assigning points to nearest center
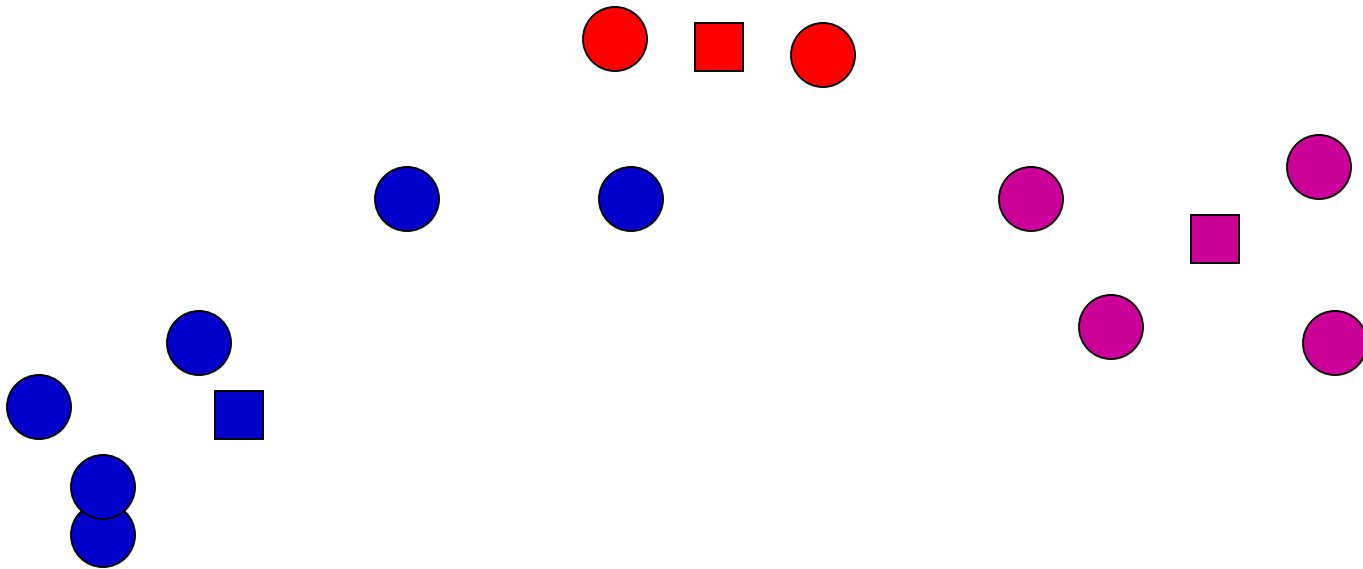
# Iteration # 1

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster
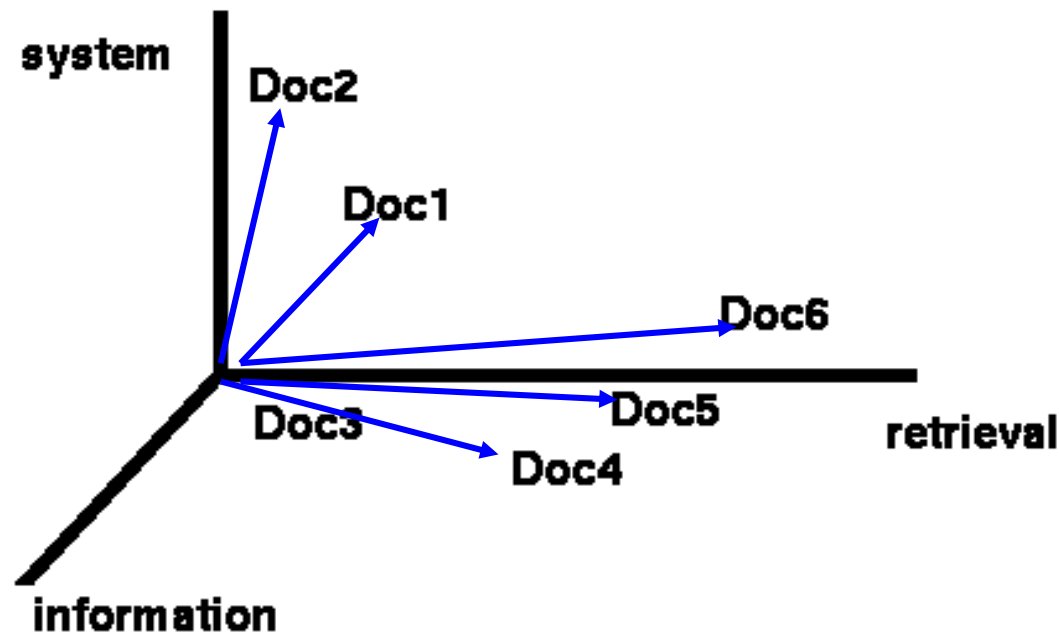
# Distance measures

Euclidean:

$$d(x, y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

good for spatial data

# Clustering documents

One feature for each word.  The value is the number of times that word occurs.

Documents are points or vectors in this space
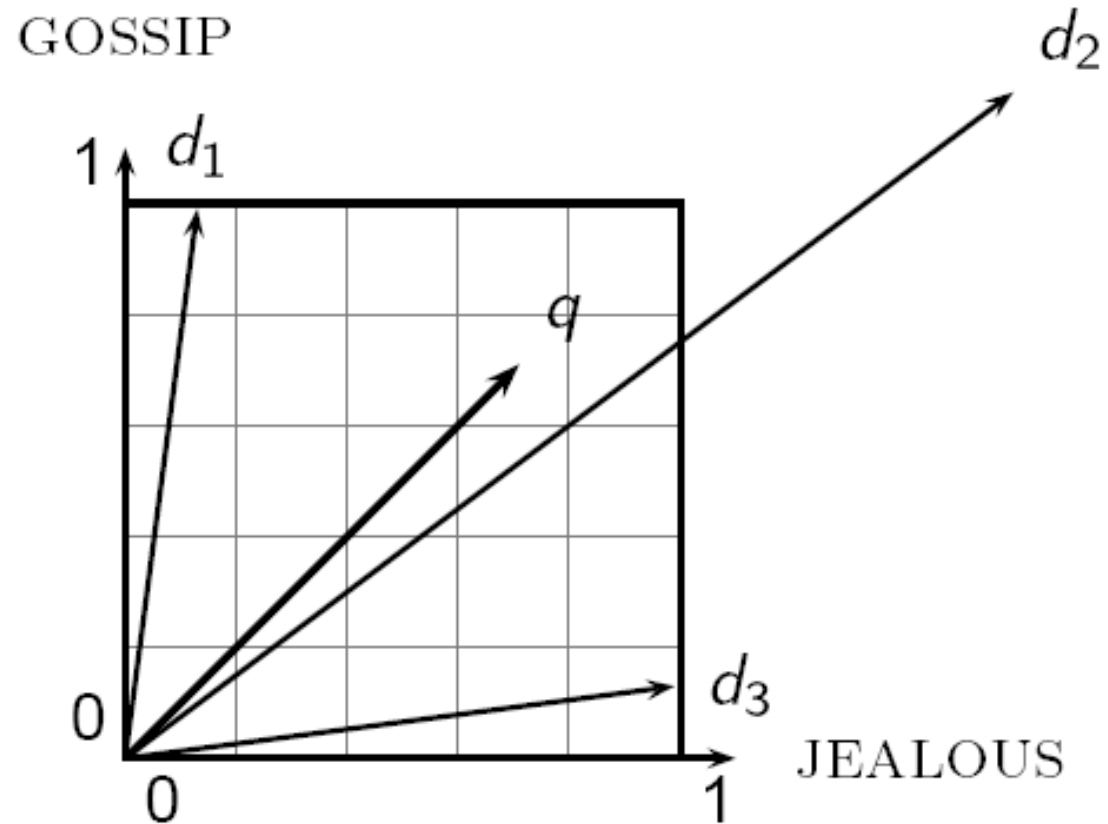
# When Euclidean distance doesn't work

Which document is closest to q using Euclidian distance?

Which do you think should be closer?

# Issues with Euclidian distance

- the Euclidean distance between $q$ and $d_2$ is large

- but, the distribution of terms in the query $q$ and the distribution of terms in the document $d_2$ are very similar

- This is not what we want!

# cosine similarity

$$sim(x, y) = \frac{x \cdot y}{|x||y|} = \frac{x}{|x|} \cdot \frac{y}{|y|} = \frac{\sum_{i=1}^{n} x_i y_i}{\sqrt{\sum_{i=1}^{n} x_i^2} \sqrt{\sum_{i=1}^{n} y_i^2}}$$

correlated with the
angle between two vectors

# cosine distance

cosine similarity is a similarity between 0 and 1, with things that are similar 1 and not 0

We want a distance measure, cosine distance:
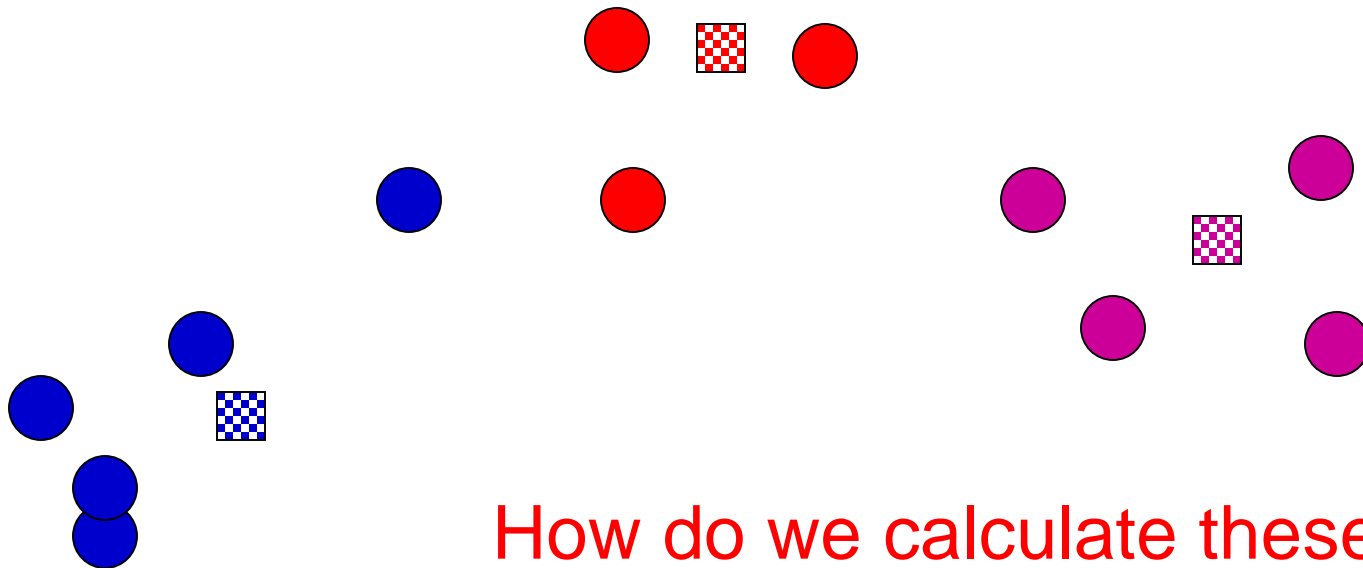
$$d(x, y) = 1 - sim(x, y)$$

- good for text data and many other "real world" data sets
- is computationally friendly since we only need to consider features that have non-zero values in **both** examples

# K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

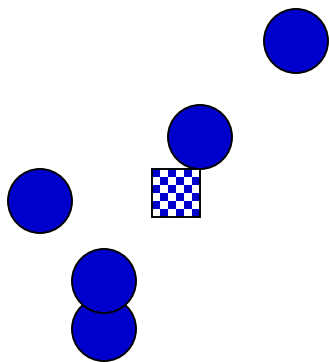How do we calculate these?

# K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

Mean of the points in the cluster:

$$m(\mathrm{C}) = \frac{1}{|C|} \mathring{\mathrm{a}}_{x\hat{\mathrm{I}}\ C} x$$

where:

$$x + y = \mathring{\mathrm{a}}_{i=1}^{n} x_i + y_i \qquad \frac{x}{|C|} = \mathring{\mathrm{a}}_{i=1}^{n} \frac{x_i}{|C|}$$

# Comments on the K-Means Method

- Strength
  - Relatively efficient: O(tkn),
    - where n is # objects,
    - k is # clusters, and
    - t  is # iterations.
    - Normally, k, t << n.
  - Often terminates at a local optimum.
- Weakness
  - Applicable only when mean is defined, then what about categorical data?
  - Need to specify k, the number of clusters, in advance
  - Unable to handle noisy data and outliers
  - Not suitable to discover clusters with non-convex shapes

# K-means loss function

K-means tries to minimize what is called the "k-means" loss function:

$$loss = \sum_{i=1}^{n} d(x_i, m_k)^2 \quad \text{where } m_k \text{ is cluster center for } x_i$$

that is, the <u>sum of the squared distances </u>from each point to the associated cluster center

# Minimizing k-means loss

Iterate:

1. Assign/cluster each example to closest center

2. Recalculate centers as the mean of the points in a cluster

$$loss = \sum_{i=1}^{n} d(x_i, m_k)^2 \quad \text{where } m_k \text{ is cluster center for } x_i$$

Does each step of k-means move towards reducing this loss function (or at least not increasing)?

# Minimizing k-means loss

Iterate:

    1. Assign/cluster each example to closest center

    2. Recalculate centers as the mean of the points in a cluster

$$loss = \sum_{i=1}^{n} d(x_i, m_k)^2 \quad \text{where } m_k \text{ is cluster center for } x_i$$

Intuition:

1. Any other assignment would end up in a larger loss

1. The mean of a set of values minimizes the squared error

# Minimizing k-means loss

Iterate:

    1. Assign/cluster each example to closest center

    2. Recalculate centers as the mean of the points in a cluster

---

$$loss = \sum_{i=1}^{n} d(x_i, m_k)^2 \quad \text{where } m_k \text{ is cluster center for } x_i$$

**Does this mean that k-means will always find the minimum loss/clustering?**

# Minimizing k-means loss

Iterate:

    1. Assign/cluster each example to closest center

    2. Recalculate centers as the mean of the points in a cluster

---

$$loss = \sum_{i=1}^{n} d(x_i, m_k)^2 \quad \text{where } m_k \text{ is cluster center for } x_i$$

NO!  It will find *a minimum*.

Unfortunately, the k-means loss function is generally not convex and for most problems has many, many minimums

We're only guaranteed to find one of them

# K-means variations/parameters

Start with some initial cluster centers

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

What are some other variations/parameters we haven't specified?
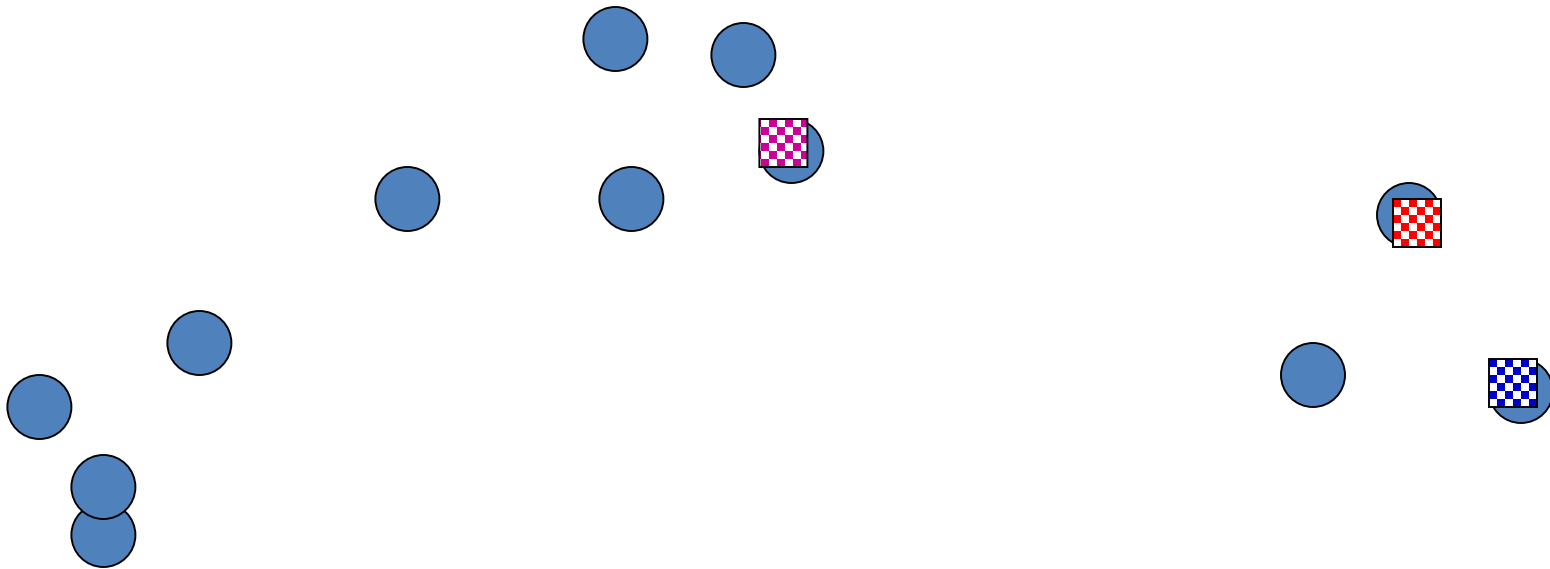
# K-means variations/parameters

Initial (seed) cluster centers

Convergence

- A fixed number of iterations
- partitions unchanged
- Cluster centers don't change

K!

# K-means: Initialize centers randomly



**What would happen here?**

**Seed selection ideas?**

# Seed choice

Results can vary drastically based on random seed selection

Some seeds can result in poor convergence rate, or convergence to sub-optimal clusterings

Common heuristics
- Random centers in the space
- Randomly pick examples
- Points least similar to any existing center (furthest centers heuristic)
- **Try out multiple starting points**
- Initialize with the results of another clustering method

# Furthest centers heuristic
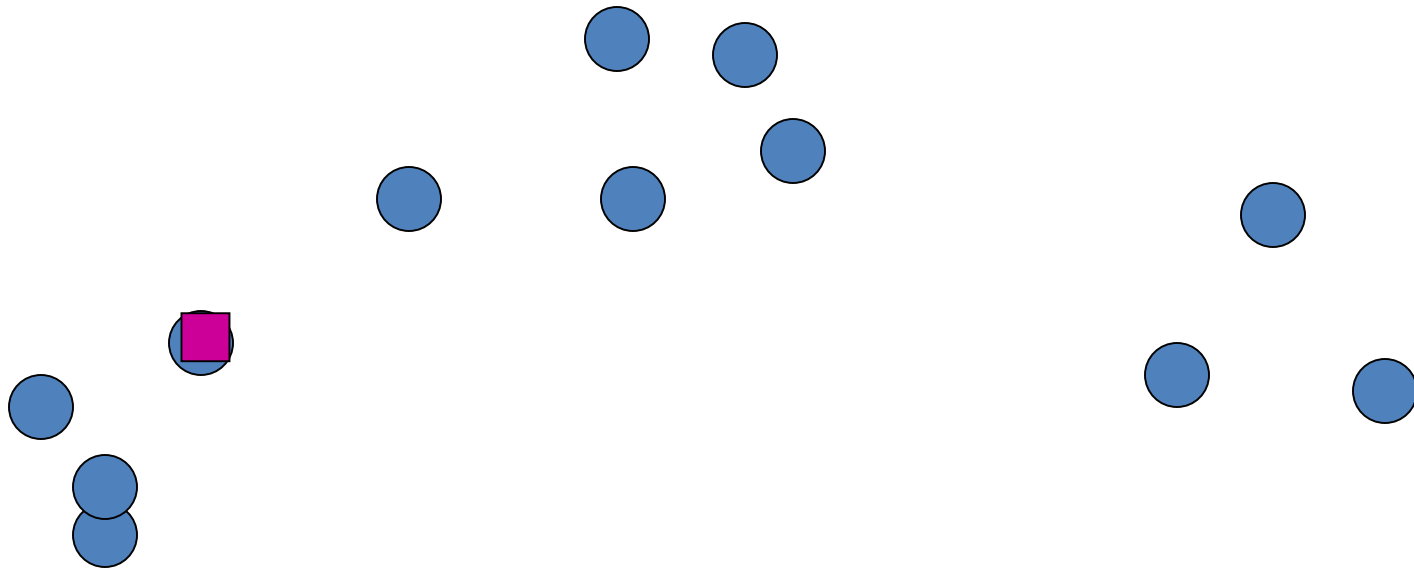
$\mu_1$ = pick random point

for i = 2 to K:

   $\mu_i$ = point that is furthest from **any** previous centers

$$m_i = \underset{x}{\mathrm{argmax}} \; \underset{m_j : 1 < j < i}{\mathrm{min}} \; d(x, m_j)$$

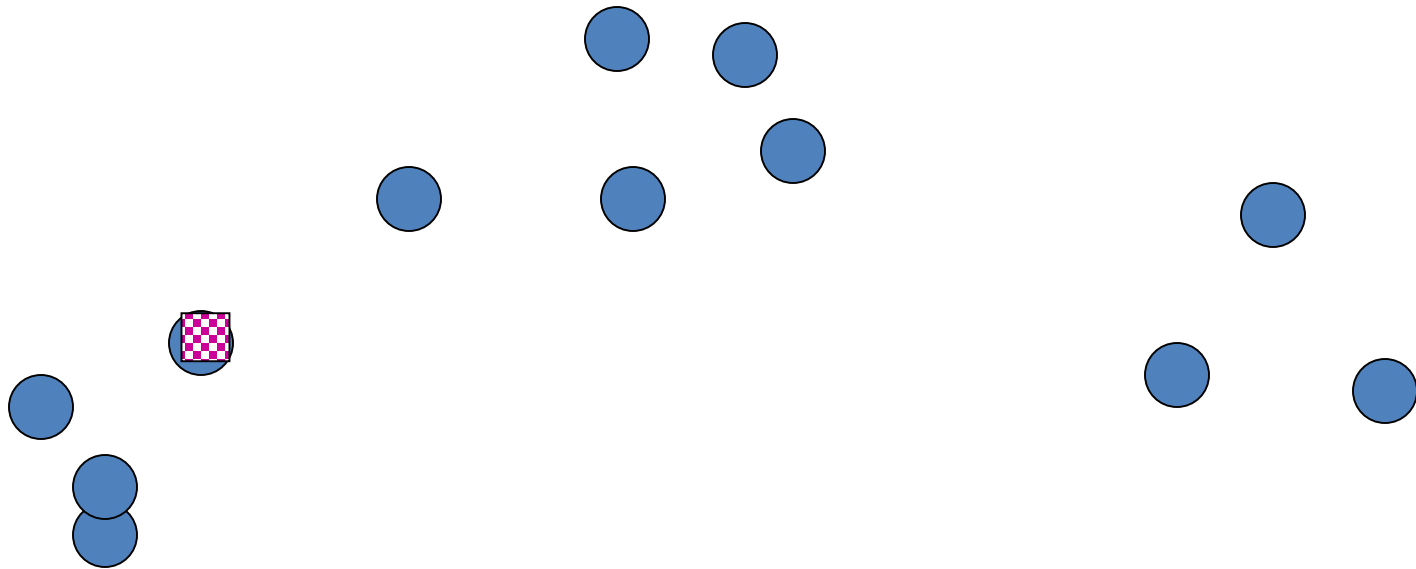point with the largest distance to any previous center

smallest distance from x to any previous center

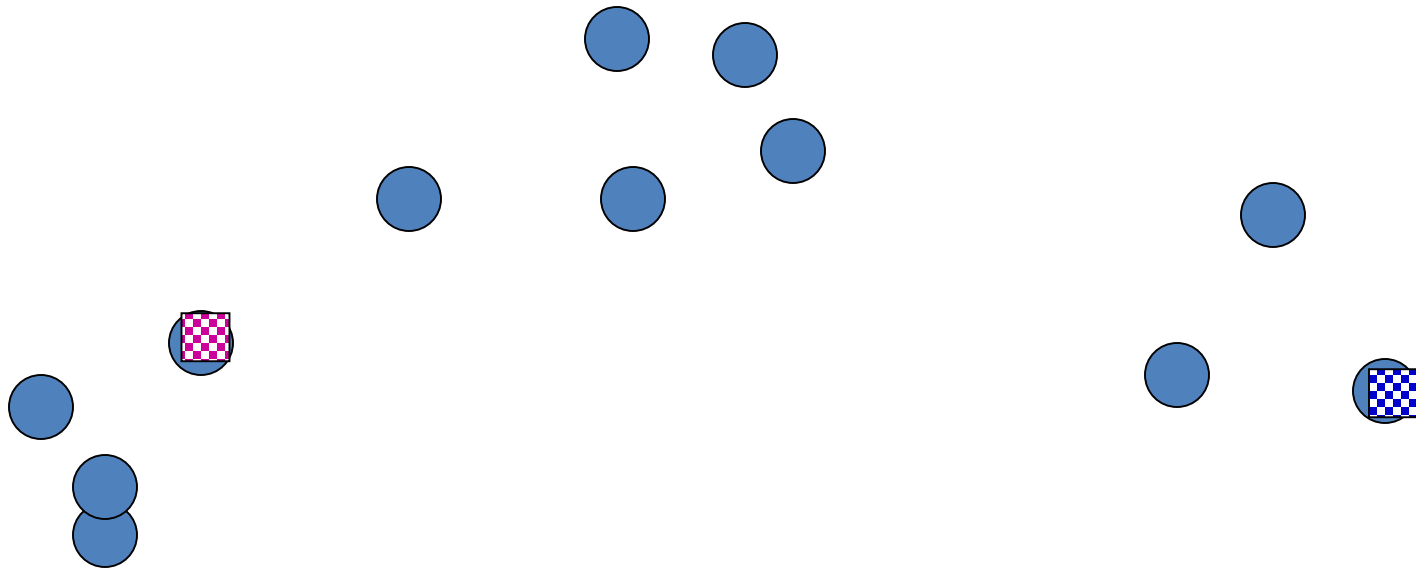# K-means: Initialize furthest from centers



Pick a random point for the first center

# K-means: Initialize furthest from centers
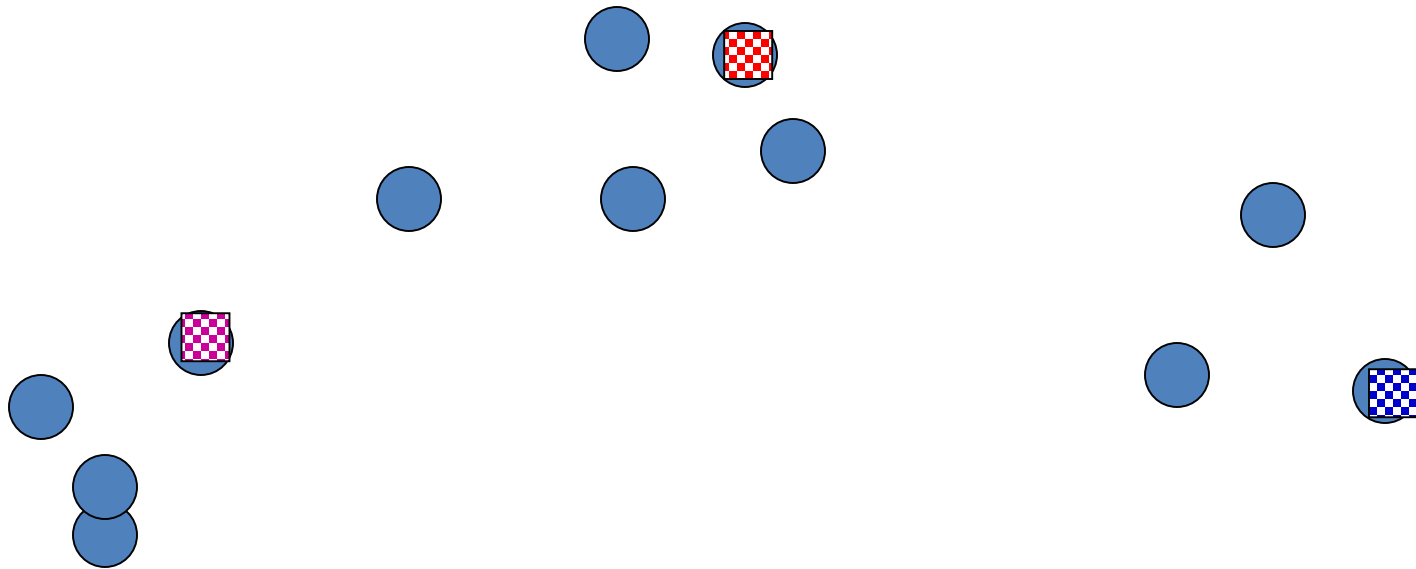
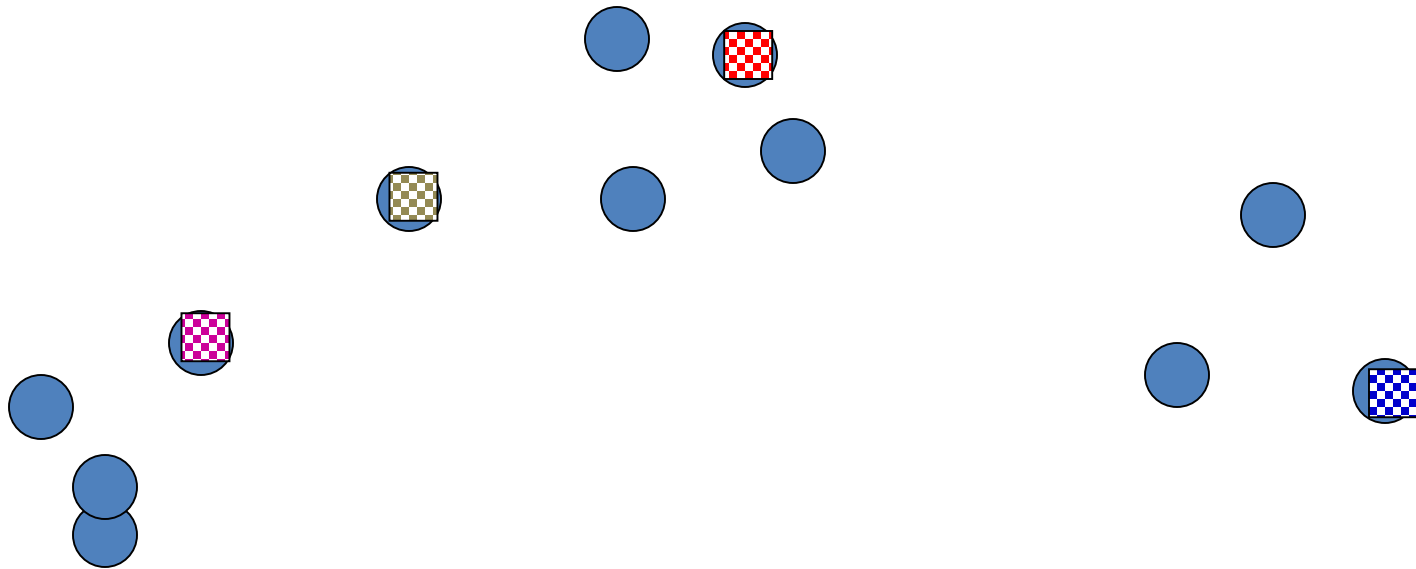What point will be chosen next?

# K-means: Initialize furthest from centers



Furthest point from center

What point will be chosen next?

# K-means: Initialize furthest from centers



Furthest point from center
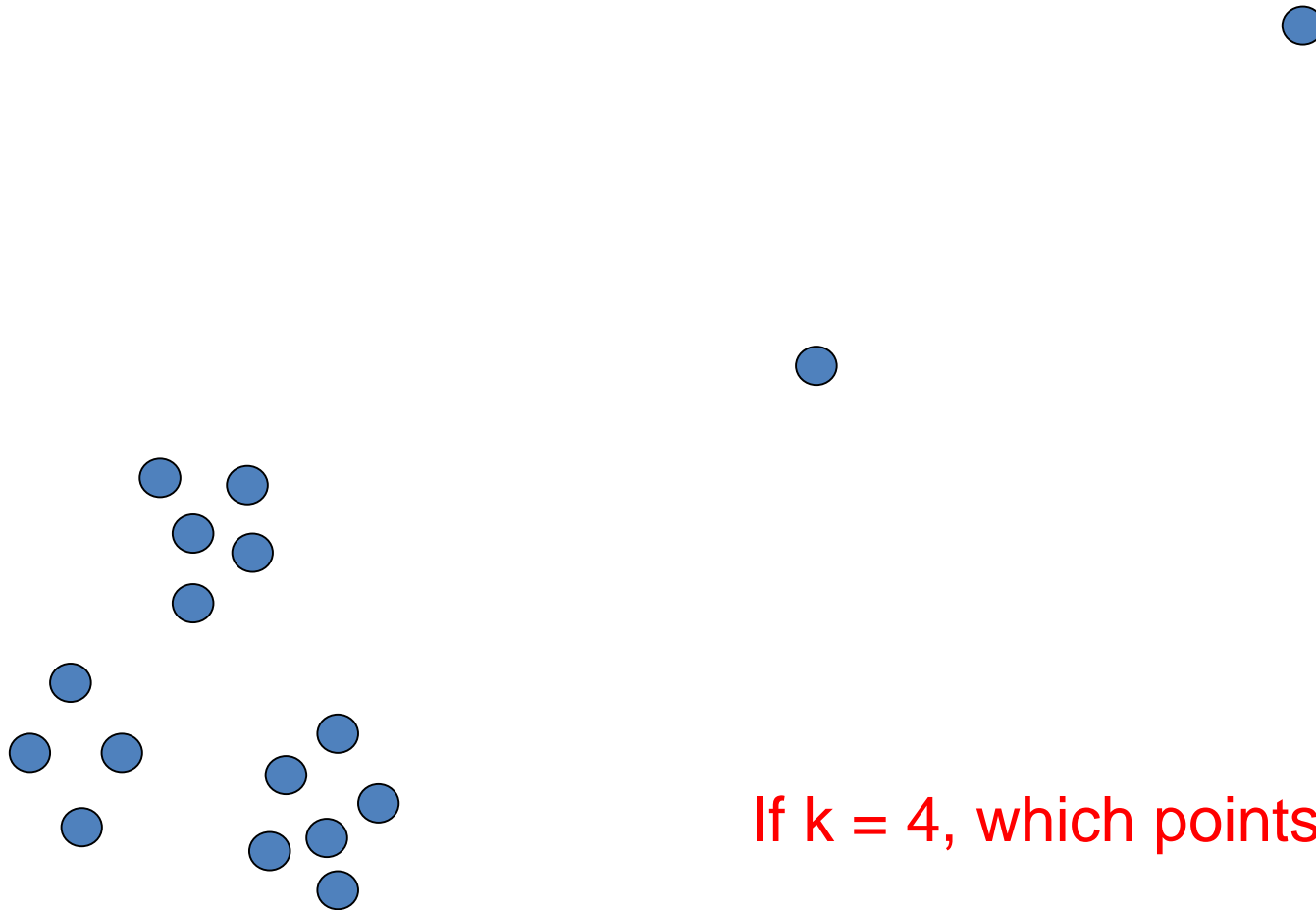
What point will be chosen next?

# K-means: Initialize furthest from centers



Furthest point from center

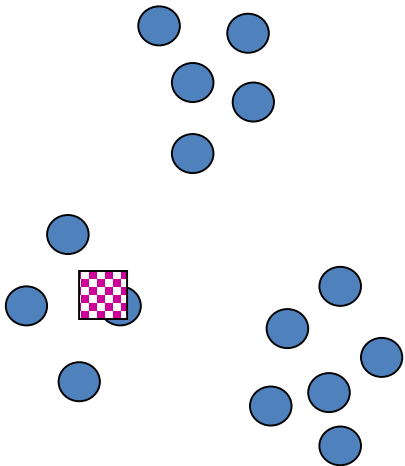Any issues/concerns with this approach?

# Furthest points concerns



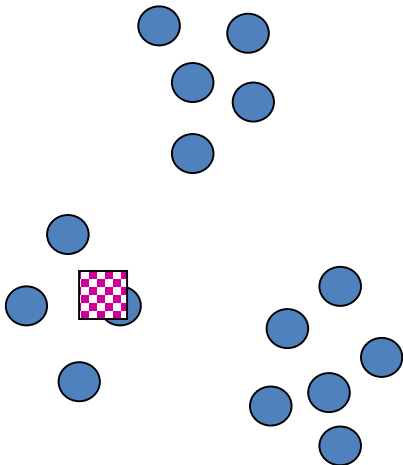If k = 4, which points will get chosen?

# Furthest points concerns

If we do a number of trials, will we get different centers?

# Furthest points concerns

Doesn't deal well with outliers

# K-means++

$\mu_1$ = pick random point

for k = 2 to **K**:
  for i = 1 to **N**:
    $s_i$ = min $d(x_i, \mu_{1\ldots k-1})$ // smallest distance to any center

  $\mu_k$ = randomly pick point ***proportionate*** to ***s***

How does this help?

# K-means++

$\mu_1$ = pick random point

for k = 2 to **K**:

    for i = 1 to **N**:
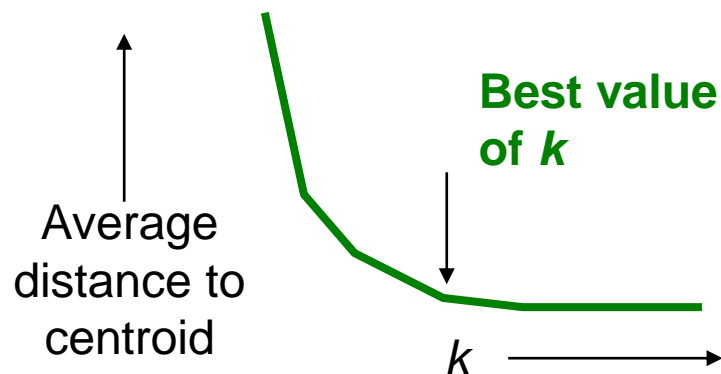
        $s_i$ = min d($x_i$, $\mu_{1\ldots k-1}$) // smallest distance to any center

    $\mu_k$ = randomly pick point ***proportionate*** to ***s***

- Makes it possible to select other points
  - if #points >> #outliers, we will pick good points
- Makes it non-deterministic, which will help with random runs
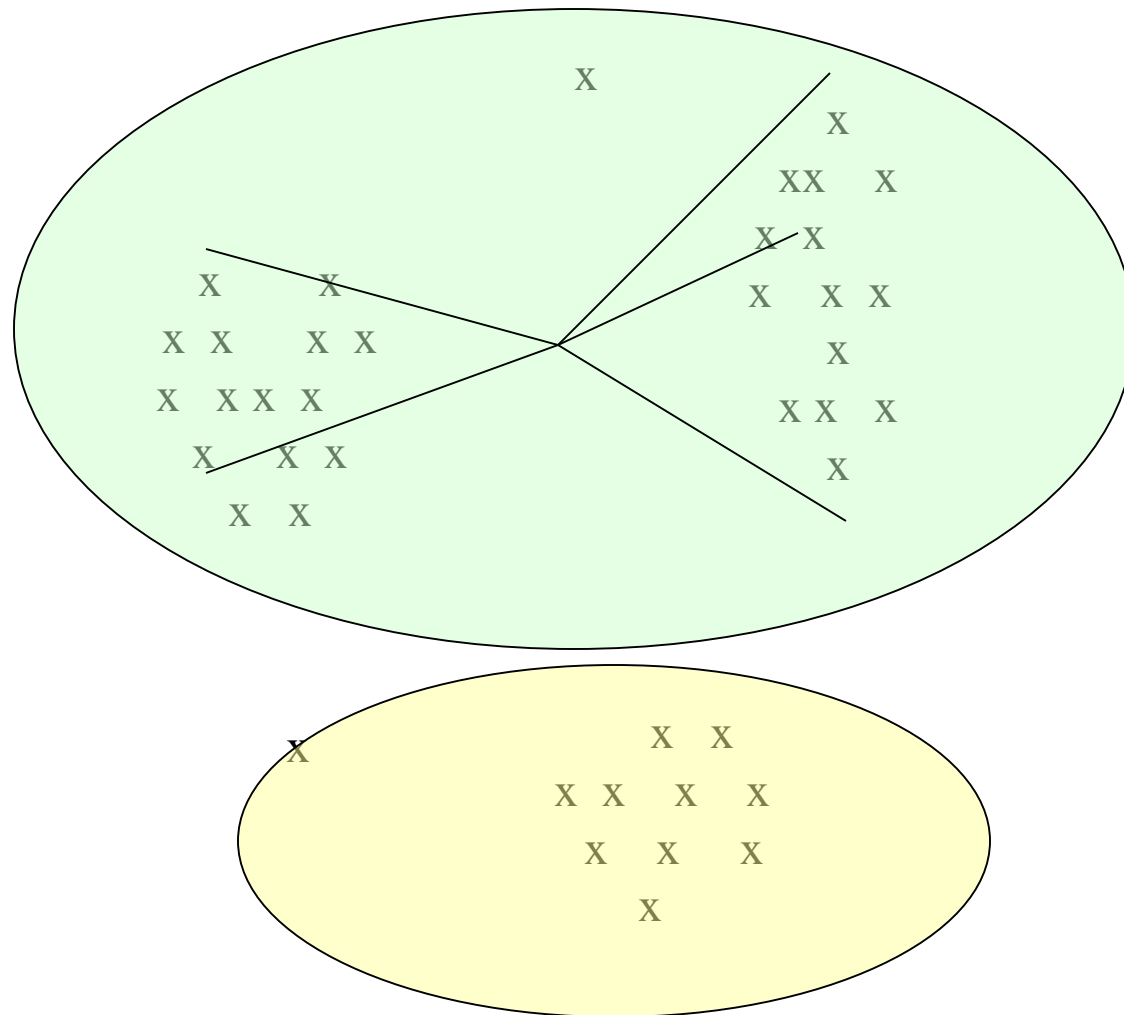- Nice theoretical guarantees!

# Getting the k right

- How to select k?
- Try different k, looking at the change in the average distance to centroid as k increases
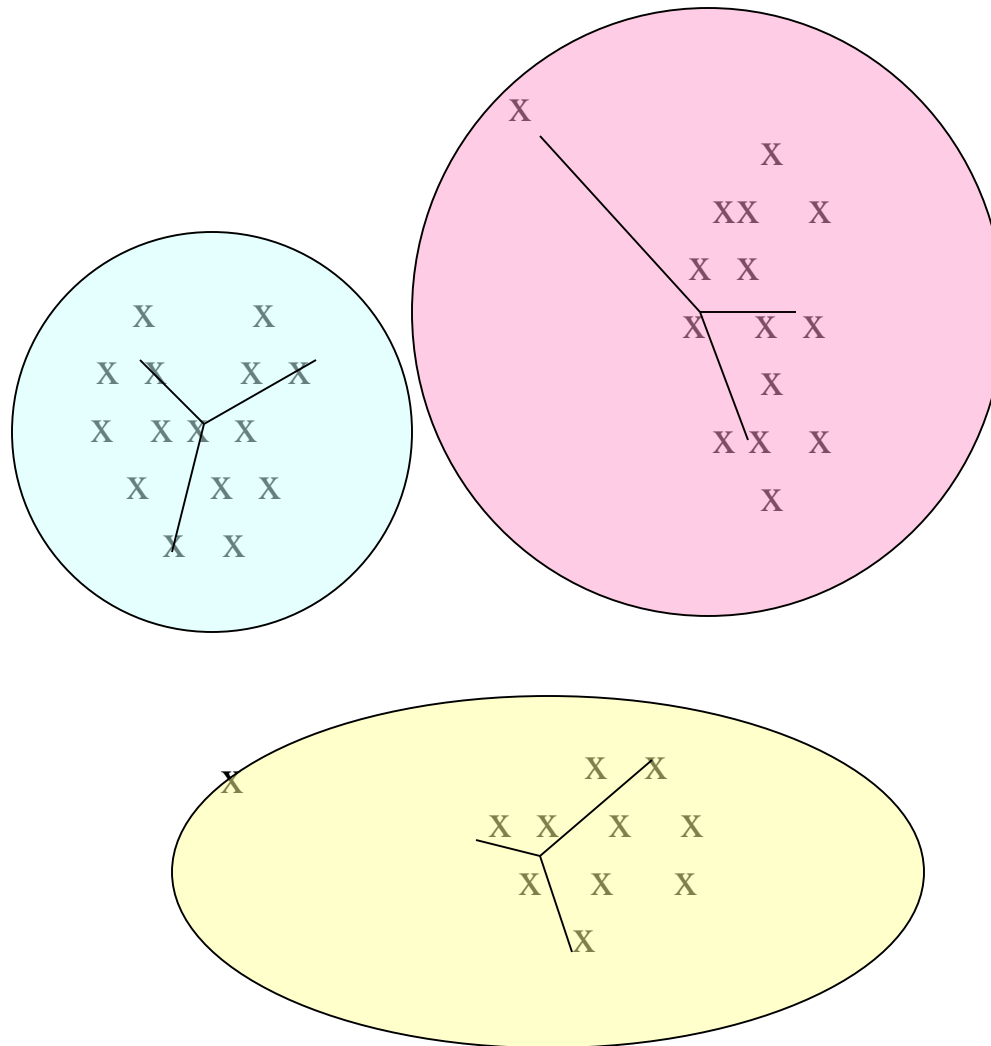- Average falls rapidly until right k, then changes little

# Example: Picking *k*

**Too few;** many long distances to centroid.

# Example: Picking *k*

**Just right;** distances rather short.

# Example: Picking *k*

**Too many;**
little improvement
in average
distance.