# Clustering and the $k$-means Algorithm

David M. Blei

COS424
Princeton University

March 2, 2012

## Clustering

- Goal: Automatically segment data into groups of similar points

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?

## Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data

## Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data

## Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space

## Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space
- Examples:

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space
- Examples:
    - Customers according to purchase histories

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
  - Automatically organizing data
  - Understanding hidden structure in some data
  - Representing high-dimensional data in a low-dimensional space
- Examples:
  - Customers according to purchase histories
  - Genes according to expression profile

## Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space
- Examples:
    - Customers according to purchase histories
    - Genes according to expression profile
    - Search results according to topic

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space
- Examples:
    - Customers according to purchase histories
    - Genes according to expression profile
    - Search results according to topic
    - Facebook users according to interests

# Clustering

- Goal: Automatically segment data into groups of similar points
- Question: When and why would we want to do this?
- Useful for:
    - Automatically organizing data
    - Understanding hidden structure in some data
    - Representing high-dimensional data in a low-dimensional space
- Examples:
    - Customers according to purchase histories
    - Genes according to expression profile
    - Search results according to topic
    - Facebook users according to interests
    - A museum catalog according to image similarity

# Clustering set-up

- Our data are
$$\mathscr{D} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}.$$

## Clustering set-up

- Our data are
$$\mathscr{D} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}.$$

- Each data point is $p$-dimensional, i.e.,
$$\mathbf{x}_n = \langle x_{n,1}, \ldots, x_{n,p} \rangle.$$

## Clustering set-up

- Our data are
$$\mathscr{D} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}.$$

- Each data point is *p*-dimensional, i.e.,
$$\mathbf{x}_n = \langle x_{n,1}, \ldots, x_{n,p} \rangle.$$

- Define a *distance function* between data, $d(\mathbf{x}_n, \mathbf{x}_m)$.

## Clustering set-up

- Our data are

$$\mathscr{D} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}.$$

- Each data point is *p*-dimensional, i.e.,

$$\mathbf{x}_n = \langle x_{n,1}, \ldots, x_{n,p} \rangle.$$

- Define a *distance function* between data, $d(\mathbf{x}_n, \mathbf{x}_m)$.

- Goal: segment the data into *k* groups

$$\{z_1, \ldots, z_N\} \quad \text{where} \quad z_i \in \{1, \ldots, K\}.$$

# Example data



500 2-dimensional data points: $\mathbf{x}_n = \langle x_{n,1}, x_{n,2} \rangle$

# Example data



- What is a good distance function here?

# Example data



- What is a good distance function here?
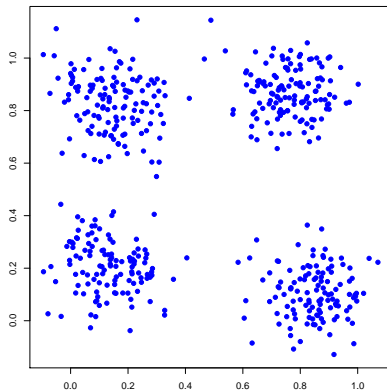- Squared Euclidean distance is reasonable

$$d(\mathbf{x}_n, \mathbf{x}_m) = \sum_{i=1}^{p} (x_{n,i} - x_{m,i})^2 = ||x_n - x_m||^2$$
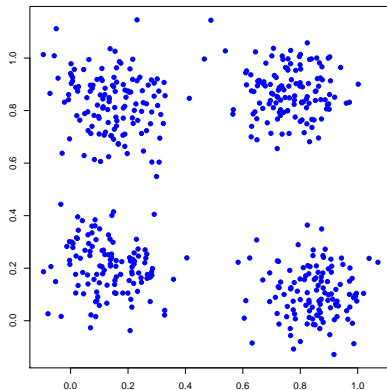
# Example data



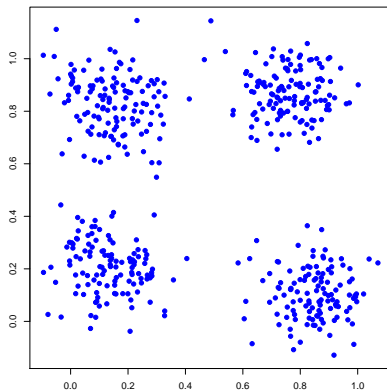- Goal: segment this data into *k* groups.

# Example data



- Goal: segment this data into *k* groups.
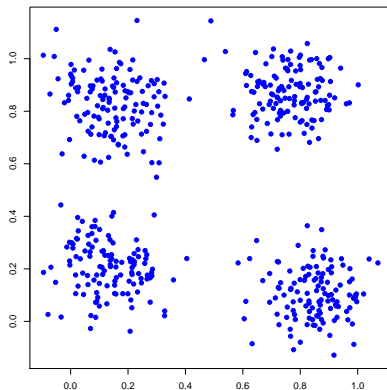- What should *k* be?

# Example data



- Goal: segment this data into *k* groups.
- What should *k* be?
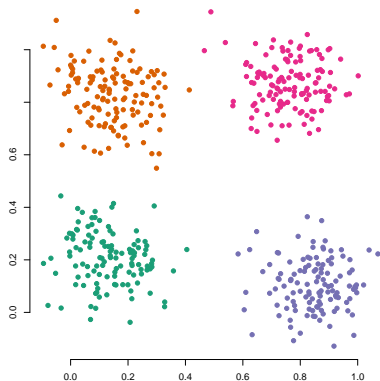- Automatically choosing *k* is complicated; for now, 4.

# *k*-means



- Different clustering algorithms use data and distance in different ways

# *k*-means



- Different clustering algorithms use data and distance in different ways
- We discuss *k*-means, the simplest clustering algorithm

# *k*-means



- Different clustering algorithms use data and distance in different ways
- We discuss *k*-means, the simplest clustering algorithm

- The basic idea is to describe each cluster by its mean value.

- The basic idea is to describe each cluster by its mean value.
- (Note: this works only for distances such that a mean is well-defined.)

- The basic idea is to describe each cluster by its mean value.
- (Note: this works only for distances such that a mean is well-defined.)
- The goal of *k*-means is to assign data to clusters and define these clusters with their means.

# $k$-means algorithm

1. Initialization

# $k$-means algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$

# *k*-means algorithm

① Initialization

- Data are $\mathbf{x}_{1:N}$
- Choose initial cluster means $\mathbf{m}_{1:k}$ (same dimension as data).

# *k*-means algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster means $\mathbf{m}_{1:k}$ (same dimension as data).

2. Repeat

# $k$-means algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster means $\mathbf{m}_{1:k}$ (same dimension as data).

2. Repeat
   1. Assign each data point to its closest mean

   $$z_n = \arg\min_{i \in \{1,\ldots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

# $k$-means algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster means $\mathbf{m}_{1:k}$ (same dimension as data).

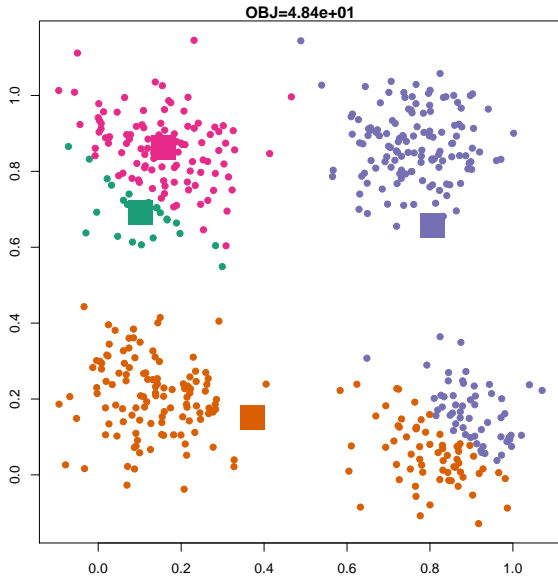2. Repeat
   1. Assign each data point to its closest mean

   $$z_n = \arg \min_{i \in \{1,\ldots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

   2. Compute each cluster mean to be the coordinate-wise average over data points assigned to that cluster,

   $$\mathbf{m}_k = \frac{1}{N_k} \sum_{\{n:\, z_n = k\}} \mathbf{x}_n$$

# *k*-means algorithm

**1** Initialization

- Data are $\mathbf{x}_{1:N}$
- Choose initial cluster means $\mathbf{m}_{1:k}$ (same dimension as data).

**2** Repeat

**1** Assign each data point to its closest mean

$$z_n = \arg \min_{i \in \{1,\dots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

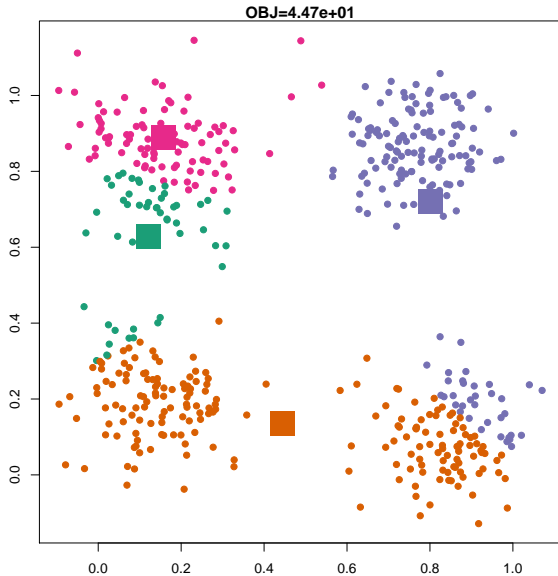**2** Compute each cluster mean to be the coordinate-wise average over data points assigned to that cluster,

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{\{n : z_n = k\}} \mathbf{x}_n$$

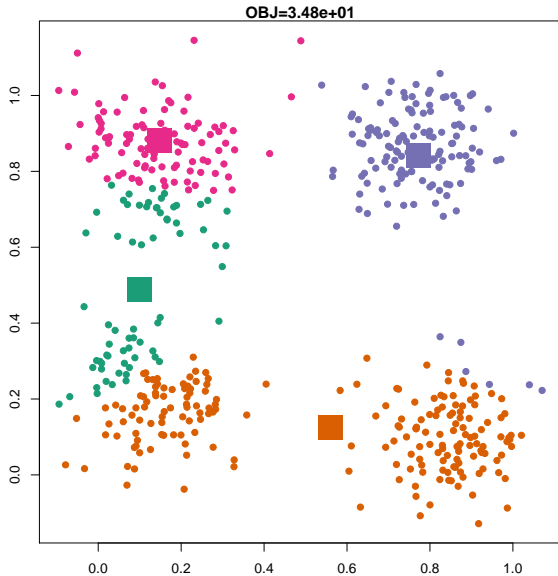**3** Until assignments $\mathbf{z}_{1:N}$ do not change
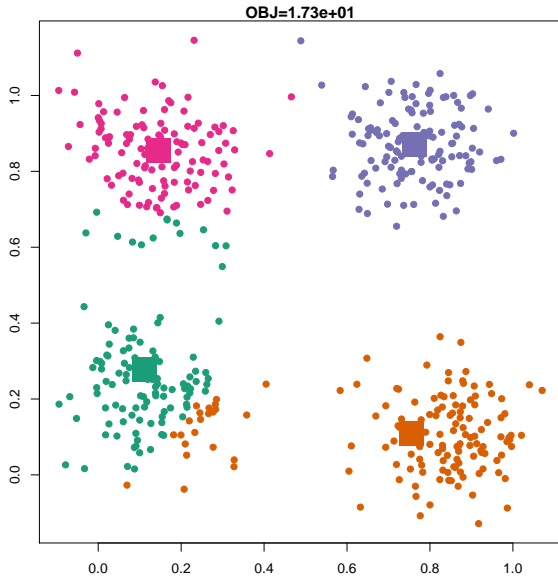
# *k*-means example



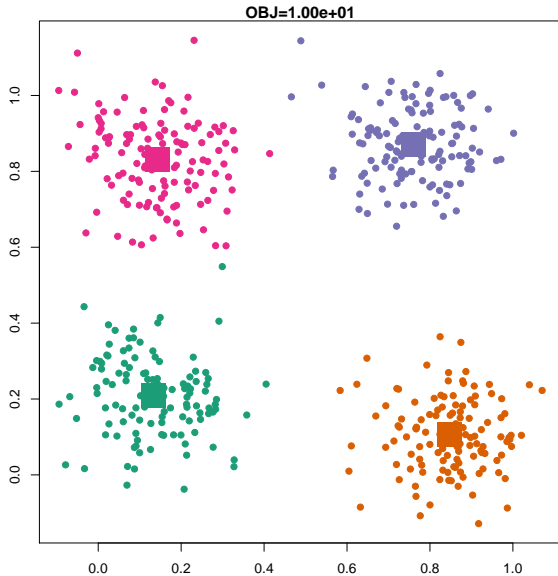OBJ=4.84e+01

# *k*-means example



OBJ=4.47e+01

# *k*-means example



OBJ=3.48e+01

# *k*-means example

# *k*-means example



OBJ=1.00e+01

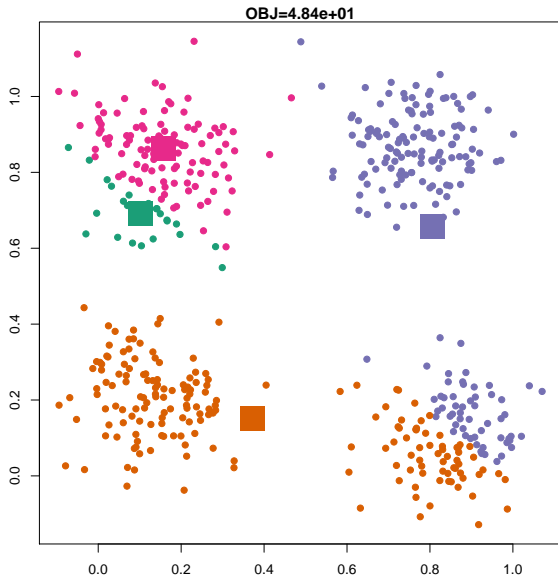# *k*-means example



OBJ=9.97e+00

# *k*-means example

# Objective function

- How can we measure how well our algorithm is doing?

- How can we measure how well our algorithm is doing?
- The $k$-means objective function is the sum of the squared distances of each point to each assigned mean

$$F(z_{1:N}, \mathbf{m}_{1:k}) = \frac{1}{2} \sum_{n=1}^{N} ||\mathbf{x}_n - \mathbf{m}_{z_n}||^2$$

# *k*-means example (look at the objective)



OBJ=4.84e+01

# *k*-means example (look at the objective)



OBJ=4.47e+01

# *k*-means example (look at the objective)



OBJ=3.48e+01

# *k*-means example (look at the objective)



OBJ=1.73e+01

# *k*-means example (look at the objective)



OBJ=1.00e+01

# *k*-means example (look at the objective)



OBJ=9.97e+00

# *k*-means example (look at the objective)



OBJ=9.97e+00

## Coordinate descent

$$F(z_{1:N}, \mathbf{m}_{1:k}) = \frac{1}{2} \sum_{n=1}^{N} \|\mathbf{x}_n - \mathbf{m}_{z_n}\|^2$$

- Holding the means fixed, assigning each point to its closest mean minimizes $F$ with respect to $z_{1:N}$.

# Coordinate descent

$$F(z_{1:N}, \mathbf{m}_{1:k}) = \frac{1}{2} \sum_{n=1}^{N} ||\mathbf{x}_n - \mathbf{m}_{z_n}||^2$$

- Holding the means fixed, assigning each point to its closest mean minimizes $F$ with respect to $z_{1:N}$.
- Holding the assignments fixed, computing the centroids of each cluster minimizes $F$ with respect to $\mathbf{m}_{1:k}$.

# Coordinate descent

$$F(z_{1:N}, \mathbf{m}_{1:k}) = \frac{1}{2} \sum_{n=1}^{N} \|\mathbf{x}_n - \mathbf{m}_{z_n}\|^2$$

- Holding the means fixed, assigning each point to its closest mean minimizes $F$ with respect to $z_{1:N}$.

- Holding the assignments fixed, computing the centroids of each cluster minimizes $F$ with respect to $\mathbf{m}_{1:k}$.

- Thus, $k$-means is a *coordinate descent* algorithm.

# Coordinate descent

$$F\big(z_{1:N}, \mathbf{m}_{1:k}\big) = \frac{1}{2} \sum_{n=1}^{N} \|\mathbf{x}_n - \mathbf{m}_{z_n}\|^2$$

- Holding the means fixed, assigning each point to its closest mean minimizes $F$ with respect to $z_{1:N}$.

- Holding the assignments fixed, computing the centroids of each cluster minimizes $F$ with respect to $\mathbf{m}_{1:k}$.

- Thus, $k$-means is a *coordinate descent* algorithm.

- It finds a *local minimum*. (Multiple restarts are often necessary.)

## Objective for the example data

# Compressing images



- Each pixel is associated with a red, green, and blue value

# Compressing images



- Each pixel is associated with a red, green, and blue value
- A $1024 \times 1024$ image is a collection of 1048576 values $\langle x_1, x_2, x_3 \rangle$, which requires 3M of storage

# Compressing images



- Each pixel is associated with a red, green, and blue value
- A $1024 \times 1024$ image is a collection of 1048576 values $\langle x_1, x_2, x_3 \rangle$, which requires 3M of storage
- How can we use $k$-means to compress this image?

- Replace each pixel $\mathbf{x}_n$ with its assignment $\mathbf{m}_{z_n}$ ("paint by numbers").

# Vector quantization



- Replace each pixel $\mathbf{x}_n$ with its assignment $\mathbf{m}_{z_n}$ ("paint by numbers").
- The $k$ means are called the *codebook*.

# Vector quantization



- Replace each pixel $\mathbf{x}_n$ with its assignment $\mathbf{m}_{z_n}$ ("paint by numbers").
- The $k$ means are called the *codebook*.
- With $k = 100$, we need 7 bits per pixel plus $100 \times 3$ bits $\approx$ 897K.

2 means

4 means

8 means

16 means

32 means

64 means

128 means

256 means

# Measure of distortion



Charlie Brown and Linus VQ Objective

- The objective gives a measure of how distorted the compressed picture is relative to the original picture

# Measure of distortion



Charlie Brown and Linus VQ Objective

- The objective gives a measure of how distorted the compressed picture is relative to the original picture
- For more clusters, the picture is less distorted.

- In many practical settings, Euclidean distance is not appropriate. When?

# *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,

# *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,
  - Discrete multivariate data, such as purchase histories

## *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,
    - Discrete multivariate data, such as purchase histories
    - Positive data, such as time spent on a web-page

# *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,
    - Discrete multivariate data, such as purchase histories
    - Positive data, such as time spent on a web-page
- *k*-medoids is an algorithm that only requires knowing distances between data points, $d_{n,m} = d(x_n, x_{m_k})$.

## *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,
    - Discrete multivariate data, such as purchase histories
    - Positive data, such as time spent on a web-page
- *k*-medoids is an algorithm that only requires knowing distances between data points, $d_{n,m} = d(x_n, x_{m_k})$.
- *No need to define the mean.*

# *k*-medoids

- In many practical settings, Euclidean distance is not appropriate. When?
- For example,
  - Discrete multivariate data, such as purchase histories
  - Positive data, such as time spent on a web-page
- *k*-medoids is an algorithm that only requires knowing distances between data points, $d_{n,m} = d(x_n, x_{m_k})$.
- *No need to define the mean.*
- Each of the clusters is associated with its most typical example

# $k$-medoids algorithm

1. Initialization

# $k$-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$

# $k$-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster identities $\mathbf{m}_{1:k}$

# $k$-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster identities $\mathbf{m}_{1:k}$
2. Repeat

# *k*-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster identities $\mathbf{m}_{1:k}$
2. Repeat
   1. Assign each data point to its closest center

$$z_n = \arg \min_{i \in \{1,\dots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

# *k*-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster identities $\mathbf{m}_{1:k}$
2. Repeat
   1. Assign each data point to its closest center

      $$z_n = \arg\min_{i \in \{1,\dots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

   2. For each cluster, find the data point in that cluster that is closest to the other points in that cluster

      $$i_k = \arg\min_{\{n\,:\,z_n=k\}} \sum_{\{m\,:\,z_m=k\}} d(\mathbf{x}_n, \mathbf{x}_m)$$

# *k*-medoids algorithm

1. Initialization
   - Data are $\mathbf{x}_{1:N}$
   - Choose initial cluster identities $\mathbf{m}_{1:k}$
2. Repeat
   1. Assign each data point to its closest center

   $$z_n = \arg \min_{i \in \{1,\dots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

   2. For each cluster, find the data point in that cluster that is closest to the other points in that cluster

   $$i_k = \arg \min_{\{n : z_n = k\}} \sum_{\{m : z_m = k\}} d(\mathbf{x}_n, \mathbf{x}_m)$$

   3. Set each cluster center equal to their closest data points

   $$m_k = \mathbf{x}_{i_k}$$

# *k*-medoids algorithm

**1** Initialization
- Data are $\mathbf{x}_{1:N}$
- Choose initial cluster identities $\mathbf{m}_{1:k}$

**2** Repeat

  **1** Assign each data point to its closest center

  $$z_n = \arg \min_{i \in \{1,\dots,k\}} d(\mathbf{x}_n, \mathbf{m}_i)$$

  **2** For each cluster, find the data point in that cluster that is closest to the other points in that cluster

  $$i_k = \arg\min_{\{n : z_n = k\}} \sum_{\{m : z_m = k\}} d(\mathbf{x}_n, \mathbf{x}_m)$$

  **3** Set each cluster center equal to their closest data points

  $$m_k = \mathbf{x}_{i_k}$$

**3** Until assignments $\mathbf{z}_{1:N}$ do not change

- Choosing $k$ is a nagging problem in cluster analysis

- Choosing $k$ is a nagging problem in cluster analysis
- Sometimes, the problem determines $k$

- Choosing $k$ is a nagging problem in cluster analysis
- Sometimes, the problem determines $k$
  - A certain required compression in VQ

- Choosing *k* is a nagging problem in cluster analysis
- Sometimes, the problem determines *k*
    - A certain required compression in VQ
    - Clustering customers for *k* salespeople in a business

## Choosing $k$

- Choosing $k$ is a nagging problem in cluster analysis
- Sometimes, the problem determines $k$
    - A certain required compression in VQ
    - Clustering customers for $k$ salespeople in a business
- Usually, we seek the "natural" clustering, but what does this mean?

## Choosing $k$

- Choosing $k$ is a nagging problem in cluster analysis
- Sometimes, the problem determines $k$
    - A certain required compression in VQ
    - Clustering customers for $k$ salespeople in a business
- Usually, we seek the "natural" clustering, but what does this mean?
- It is not well-defined.

## What happens as $k$ increases?

# What happens as $k$ increases?

# What happens as $k$ increases?

# What happens as $k$ increases?



OBJ=9.97e+00

# What happens as $k$ increases?

# What happens as $k$ increases?

# What happens as $k$ increases?

# What happens as $k$ increases?

# Heuristic: A kink in the objective



- Notice the "kink" in the objective between 3 and 5.

- Notice the "kink" in the objective between 3 and 5.
- This suggests that 4 is the right number of clusters.

# Heuristic: A kink in the objective



- Notice the "kink" in the objective between 3 and 5.
- This suggests that 4 is the right number of clusters.
- Tibshirani (2001) presents a method for finding this kink.

# Archeology

- Spatial and Statistical Inference of Late Bronze Age Polities in the Southern Levant (Savage and Falconer)

## Archeology

- Spatial and Statistical Inference of Late Bronze Age Polities in the Southern Levant (Savage and Falconer)
- Cluster the location of archeological sites in Israel

## Archeology

- Spatial and Statistical Inference of Late Bronze Age Polities in the Southern Levant (Savage and Falconer)
- Cluster the location of archeological sites in Israel
- Make inferences about political history based on the clusters

## Archeology

- Spatial and Statistical Inference of Late Bronze Age Polities in the Southern Levant (Savage and Falconer)
- Cluster the location of archeological sites in Israel
- Make inferences about political history based on the clusters
- Choose $k$ very carefully, with a complicated computational technique.

Legend:
- Late Bronze polities based on textual evidence.
- Late Bronze "city-states"
- ③ Site clusters from k-means analysis
- ■ Key sites
- · Other sites

Lebanon

Sea of Galilee

Jordan River

Dead Sea

Negev Desert

Labeled sites: Acco, Achshaf, Shim'on, Dor, Megiddo, Gath-Carmel, Gath-Padalla, Shechem, Aphek, Gezer, Gath, Ashkelon, Lachish, Debir, Yurza, Hazor, Anaharath, Ta'anach, Rehov, Jerusalem

Cluster numbers: 23, 12, 8, 5, 15, 18, 3, 11, 24, 21, 20, 14, 22, 19, 7, 13, 2, 10, 16, 17, 9, 1

LBA - 24 cluster solution

0    20    40 km

N
W      E
S

- Coping with cold: An integrative, multitissue analysis of the transciptome of a poikilothermic vertebrate (Gracey et al., 2004)

## Computational Biology

- Coping with cold: An integrative, multitissue analysis of the transciptome of a poikilothermic vertebrate (Gracey et al., 2004)
- Exposed carp to different levels of cold

- Coping with cold: An integrative, multitissue analysis of the transciptome of a poikilothermic vertebrate (Gracey et al., 2004)
- Exposed carp to different levels of cold
- Clustered genes based on their response in different tissues

## Computational Biology

- Coping with cold: An integrative, multitissue analysis of the transciptome of a poikilothermic vertebrate (Gracey et al., 2004)
- Exposed carp to different levels of cold
- Clustered genes based on their response in different tissues
- (No mention of how $k = 23$ was chosen.)

## Education

- Teachers as Sources of Middle School Students' Motivational Identity: Variable-Centered and Person-Centered Analytic Approaches (Murdock and Miller, 2003)

- Teachers as Sources of Middle School Students' Motivational Identity: Variable-Centered and Person-Centered Analytic Approaches (Murdock and Miller, 2003)
- Clustered survey results of 206 students

# Education

- Teachers as Sources of Middle School Students' Motivational Identity: Variable-Centered and Person-Centered Analytic Approaches (Murdock and Miller, 2003)

- Clustered survey results of 206 students

- Used the clusters to identify groups to buttress an analysis of what affects motivation.

- Teachers as Sources of Middle School Students' Motivational Identity: Variable-Centered and Person-Centered Analytic Approaches (Murdock and Miller, 2003)

- Clustered survey results of 206 students

- Used the clusters to identify groups to buttress an analysis of what affects motivation.

- I.e., the levels of encouragement are corrected for

## Education

- Teachers as Sources of Middle School Students' Motivational Identity: Variable-Centered and Person-Centered Analytic Approaches (Murdock and Miller, 2003)
- Clustered survey results of 206 students
- Used the clusters to identify groups to buttress an analysis of what affects motivation.
- I.e., the levels of encouragement are corrected for
- Chose the number of clusters to get nice results

TABLE 3. Five-Cluster Solution: Z scores on Each Clustering Variable

| | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 | Cluster 5 |
|---|---|---|---|---|---|
| Teacher caring | −.5 | −.5 to .5 | −.5 to .5 | −.5 | 1.0 |
| Peers' academic support | 1.0 | −.5 | 1.0 | −.5 | −.5 to .5 |
| Parents' academic support | .5 | −1.0 | −.5 to .5 | −.5 to .5 | 1.0 |

TABLE 4. Means and Standard Deviations for Each Cluster on Grade 8 Motivational Variables

| Cluster | Academic Self-Efficacy | | Intrinsic Valuing of Education | | Teacher-Rated Effort | |
|---|---|---|---|---|---|---|
| | M | SD | M | SD | M | SD |
| 1. All positive | 3.59 | .48[a] | 2.99 | .55[a] | 3.74 | .26[a] |
| 2. Peer negative, parents very negative | 2.44 | .66[b] | 2.16 | .51[b] | 3.05 | .61[b] |
| 3. Peer positive | 3.01 | .73[c] | 2.43 | .66[b] | 3.26 | .66[b] |
| 4. Negative teacher and peer | 2.47 | .63[b] | 2.24 | .51[b] | 3.17 | .59[b] |
| 5. Positive teacher and parents | 3.19 | .65[c] | 2.89 | .62[a] | 3.54 | .47[a] |

- Implications of Racial and Gender Differences in Patterns of Adolescent Risk Behavior for HIV and other Sexually Transmitted Diseases (Halpert et al., 2004)

# Sociology

- Implications of Racial and Gender Differences in Patterns of Adolescent Risk Behavior for HIV and other Sexually Transmitted Diseases (Halpert et al., 2004)
- Clustered survey results of 13,998 students to understand patterns of drug abuse and sexual activity

# Sociology

- Implications of Racial and Gender Differences in Patterns of Adolescent Risk Behavior for HIV and other Sexually Transmitted Diseases (Halpert et al., 2004)
- Clustered survey results of 13,998 students to understand patterns of drug abuse and sexual activity
- *K* chosen for interpretability and "stability," which means that they could interpret multiple *k*-means runs on different data in the same way.

- Implications of Racial and Gender Differences in Patterns of Adolescent Risk Behavior for HIV and other Sexually Transmitted Diseases (Halpert et al., 2004)
- Clustered survey results of 13,998 students to understand patterns of drug abuse and sexual activity
- $K$ chosen for interpretability and "stability," which means that they could interpret multiple $k$-means runs on different data in the same way.
- Draw the conclusion that patterns exist. What's wrong with this?

## Sociology

- Implications of Racial and Gender Differences in Patterns of Adolescent Risk Behavior for HIV and other Sexually Transmitted Diseases (Halpert et al., 2004)
- Clustered survey results of 13,998 students to understand patterns of drug abuse and sexual activity
- *K* chosen for interpretability and "stability," which means that they could interpret multiple *k*-means runs on different data in the same way.
- Draw the conclusion that patterns exist. What's wrong with this?
- *k*-means will find patterns everywhere!

**TABLE 2. Percentage distribution of participants, by cluster, and behavioral patterns defining each cluster**

| Cluster type and behavioral patterns | % |
|---|---|
| **Light substance dabblers**—infrequent or no current use of substances†<br>None have had sex | 24.4 |
| **Abstainers**—none have ever used substances† or had sex | 22.7 |
| **Sex dabblers**—all have had sex<br>Median no. of partners=1<br>60% used a condom at last sex<br>Infrequent use of substances† | 14.5 |
| **Drinkers**—all consumed alcohol in past 12 mos.<br>49% report binge drinking<br>Infrequent or no illicit drug use<br>None have had sex | 7.4 |
| **Smokers**—all smoke cigarettes daily<br>Infrequent use of alcohol/illicit drugs<br>62% have had sex | 7.3 |
| **Alcohol-and-sex dabblers**—all drink occasionally; all have had sex<br>Infrequent tobacco/illicit drug use | 5.4 |
| **Binge drinkers**—all binge frequently<br>Infrequent cigarette, marijuana and other drug use<br>60% binge ≥1 time/wk.<br>45% have had sex | 4.4 |
| **Heavy dabblers**—all smoke, drink and binge drink with moderate frequency<br>45% use marijuana; few use other illicit drugs<br>91% have had sex | 3.6 |
| **Combination sex and drug use**—all have had sex; all used alcohol/illicit drug at last sex | 3.4 |
| **Marijuana users**—all use marijuana frequently; few have used other illicit drugs<br>94% use alcohol<br>79% smoke cigarettes<br>74% have had sex | 1.7 |
| **Multiple partners**—all report ≥14 sexual partners<br>75% report low or moderate use of substances† | 1.3 |
| **Sex for drugs or money**—all have had sex for drugs or money<br>50% report low or moderate use of substances†<br>Median no. of partners=3 | 1.2 |
| **High marijuana use and sex**—all use marijuana frequently; all have had sex<br>All used alcohol/other drug at last sex<br>82% have had >1 partner (median=6) | 1.1 |
| **Marijuana and other drug users**—95% report heavy marijuana use; all use other illicit drugs<br>68% have had sex<br>28% used alcohol/other drug at last sex | 0.6 |
| **Injection-drug users**—all have injected drugs<br>82% have had sex<br>Median no. of partners=4 | 0.6 |
| **Males who have sex with males**—all are males who have sex with another male<br>78% have had multiple partners (median=5)<br>40% used marijuana in past 30 days<br>50% use alcohol ≥1 time/mo.<br>17% have had sex for drugs or money | 0.3 |

# Summary