

K-Means

Machine Learning 1 — Lecture 12

13th June 2023

Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology

Supervised Learning

- Linear regression, logistic regression
- Neural networks
- K-nearest neighbours
- Decision trees, random forests
- Support vector machines (SVMs)
- Linear discriminant analysis

Unsupervised Learning

- Principal component analysis
- K-means (clustering, today)
- Gaussian mixture models (density estimation, soft clustering, next time)

Clustering

Organize N objects in K clusters, i.e., groups of similar objects.

Relation to classification: “classes” are not given, but should be found in unsupervised way.

Clustering is not well-defined and there are multiple questions:

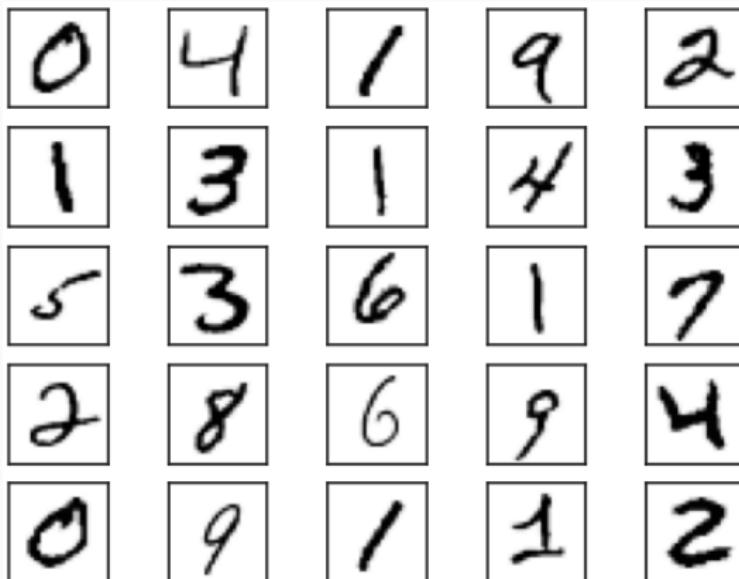
- We want to group similar objects together, but what does “similar” mean?
- What is K , i.e., the number of clusters? Provided by the user or found automatically?
- Can an object belong to more than one cluster? (**multi-view clustering**)
- Does an object need to belong to a cluster, i.e., do we allow **outliers**?
- Is membership of objects to clusters “hard” (0 or 1) or “soft” (probabilistic, fuzzy)?

How do we cluster a deck of cards?



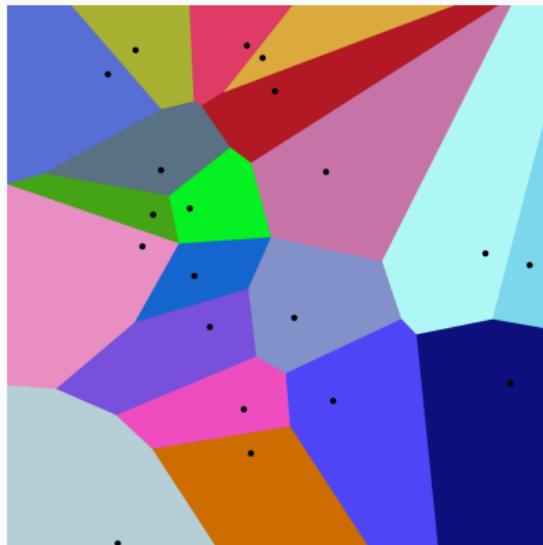
- suits? (hearts, spades, clubs, diamonds)
- value? (2, 3, . . . , King, Ace)
- court cards vs. numerical?

How to cluster handwritten digits (MNIST data)?



- Objects (data points) are represented with D -dimensional vectors $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
- **Cluster centers** are represented as D -dimensional vectors $\mu^{(1)}, \dots, \mu^{(K)}$
- Each data point is assigned to cluster whose center is closest
- Hence, similarity is defined via **Euclidean distance**
- Data points which are closest to $\mu^{(k)}$ form the k^{th} **cluster**
$$\mathcal{D}_k = \{\mathbf{x} \mid \mathbf{x} \in \mathcal{D} \wedge k = \arg \min_l \|\mathbf{x} - \mu^{(l)}\|_2\}$$
- $\mu^{(k)}$ is a **representative** or **prototype** of cluster \mathcal{D}_k
- \mathcal{D}_k is the set of all data points which fall in the **Voronoi cell** specified by $\mu^{(k)}$

Voronoi cells: Given a set of vectors $\mu^{(1)}, \dots, \mu^{(K)}$, the k^{th} Voronoi cell is the set of all points which are closest $\mu^{(k)}$. Voronoi cells form a partition of the Euclidean space.



- Cluster centers: $\mu^{(1)}, \dots, \mu^{(K)}$
- Clusters: $\mathcal{D}_k = \{\mathbf{x} \mid \mathbf{x} \in \mathcal{D} \wedge k = \arg \min_l \|\mathbf{x} - \mu^{(l)}\|_2\}$
- The **k-means problem**: Minimize the sum of squared distances between data points and their assigned cluster center:

$$\min_{\mu^1, \dots, \mu^K} \mathcal{L}_{WCSS} = \min_{\mu^1, \dots, \mu^K} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{D}_k} \|\mathbf{x} - \mu^{(k)}\|_2^2$$

- The objective \mathcal{L}_{WCSS} is called the **within-cluster sum-of-squares** (WCSS), **inertia**, or simply **k-means objective**
- For a proper optimization problem, we should explicitly encode the cluster assignments \mathcal{D}_k

- Let $z_{i,k}$ be **binary indicator variables** defined as

$$z_{i,k} = \begin{cases} 1 & \text{if data point } \mathbf{x}^{(i)} \text{ is assigned to } \boldsymbol{\mu}^{(k)} \\ 0 & \text{otherwise} \end{cases}$$

- Let \mathbf{Z} be the $N \times K$ matrix collecting all indicators:

$$\mathbf{Z} = \begin{pmatrix} z_{1,1} & z_{1,2} & \dots & z_{1,K} \\ z_{2,1} & z_{2,2} & \dots & z_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ z_{N,1} & z_{N,2} & \dots & z_{N,K} \end{pmatrix}$$

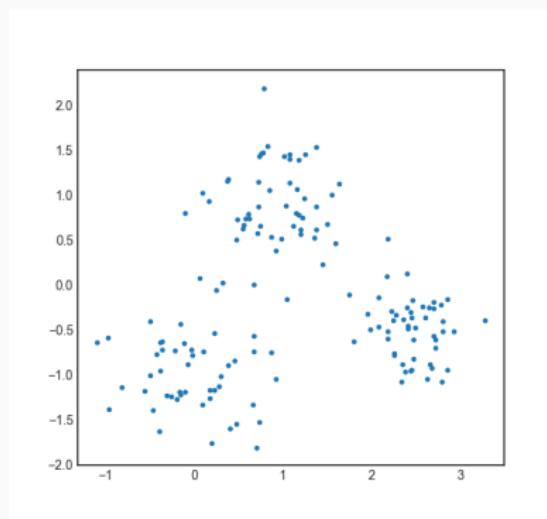
- Rows are a **one-hot encoding**

- We can now write the k-means problem as

$$\begin{aligned} \min_{\mu^1, \dots, \mu^K, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|^2 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \quad (\text{integer constraints}) \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

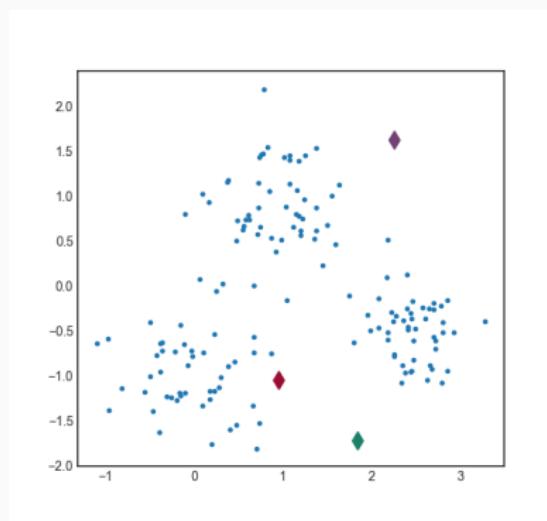
- Mixed Integer Quadratic Program (MIQP)**
- Constraints for $z_{i,k}$ ensure that each data point is assigned to exactly one cluster center (one hot encoding)
- Multiplication with indicator $z_{i,k}$ in the objective ensures that only the distance to the assigned cluster center counts
- We can now decompose the problem and alternate between Z and $\{\mu^1, \dots, \mu^K\}$, leading to the **k-means algorithm**

Assume the data set below. It seems to consist of roughly three clusters, so we set $K = 3$.



Let's "guess" a first solution for $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$, e.g. by drawing them uniformly between minimal and maximal range in each dimension.

Given $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$, how do we chose Z ?

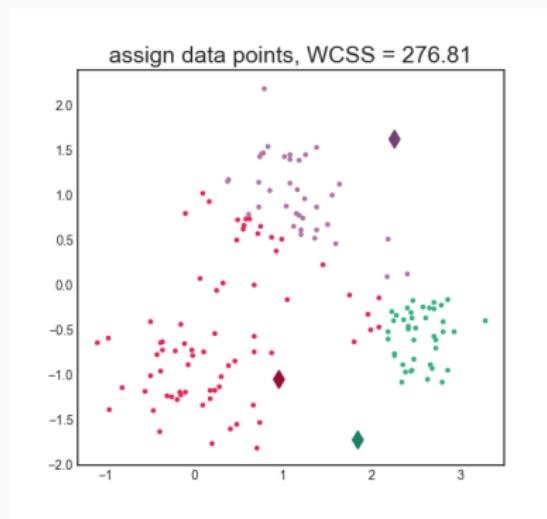


K-Means Algorithm

Example

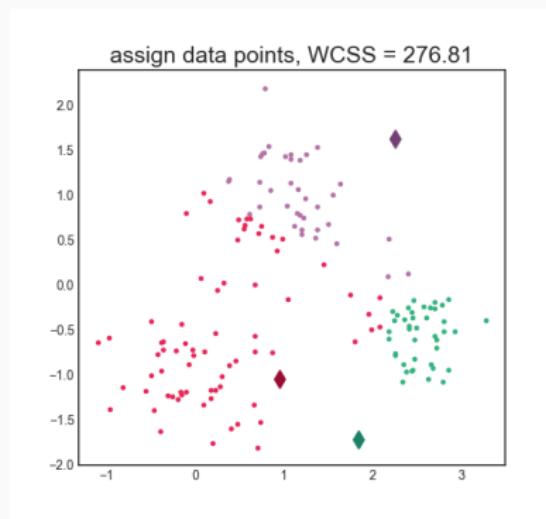
Given $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$, we can easily find optimal Z , by setting

$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$



Now we have a new cluster assignment Z and the cluster centers seem “outdated.”

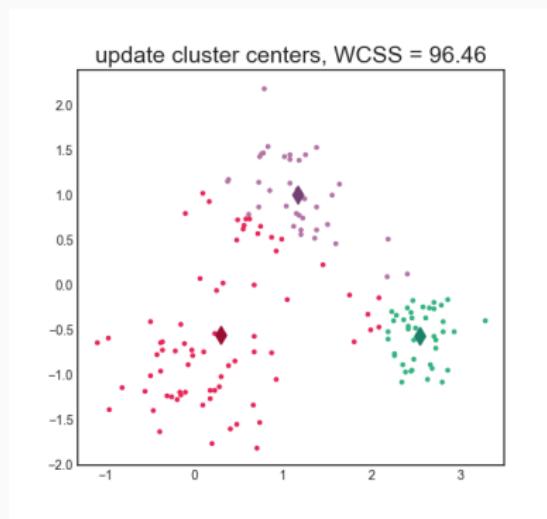
How do we chose new $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$?



Intuitively, we would set $\mu^{(k)}$ to the **mean vector** of the k^{th} cluster (hence **k-means**):

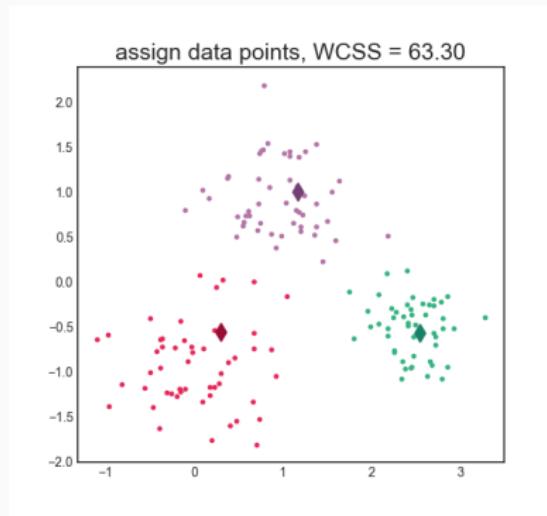
$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$

Now Z seems outdated.



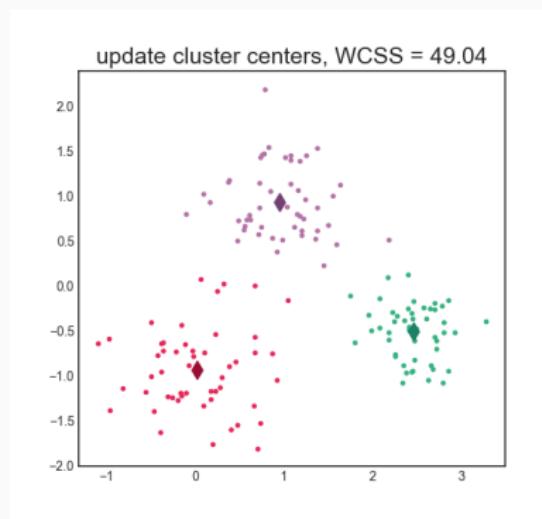
Note that we have effectively the **same situation as after initialization**, i.e. we set again

$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$



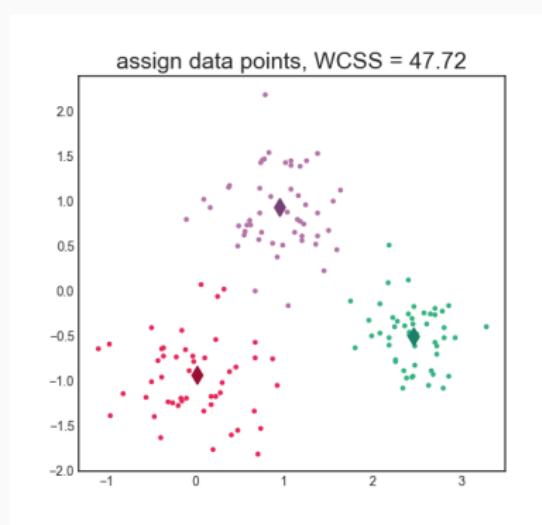
Again, update mean vectors:

$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$



Re-assign the data points to clusters:

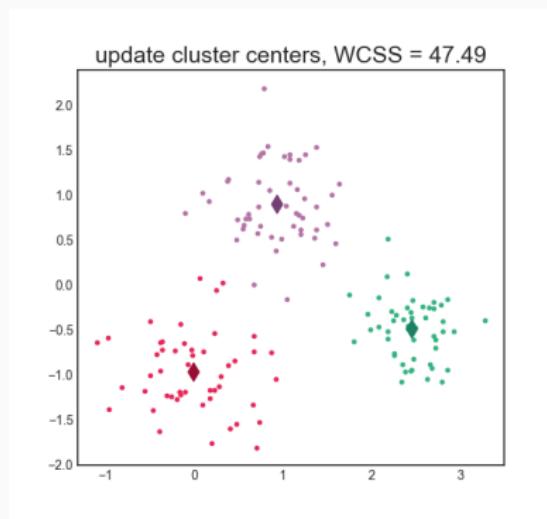
$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$



Update means:

$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$

After that, the assignments Z won't change any more. Hence, also $\mu^{(k)}$'s won't change any more. The algorithm has **converged**.



K-means Algorithm

input: Data $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

output: Cluster centers $\mu^{(1)}, \dots, \mu^{(K)}$, assignment matrix Z

Initialize $\mu^{(1)}, \dots, \mu^{(K)}$ randomly

for counter = 1 ... max number of iterations **do**

for $i = 1 \dots N$ **do**

 let $k^* = \arg \min_k \|x^{(i)} - \mu^{(k)}\|_2$

 set $z_{i,k} \leftarrow \begin{cases} 1 & \text{if } k = k^* \\ 0 & \text{otherwise} \end{cases}$

end

for $k = 1 \dots K$ **do**

 let $\mathcal{D}_k = \{x^{(i)} | z_{i,k} = 1\}$

$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$

end

if WCSS = WCSS of previous iteration **then**

 | break

end

end

K-Means Problem vs. K-Means Algorithm

- Nitpick: difference between **k-means problem** vs. **k-means algorithm**
- **k-means problem**: the optimization problem

$$\begin{aligned} \min_{\mu^1, \dots, \mu^K, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|_2^2 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \quad (\text{integer constraints}) \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

- **k-means algorithm**: the iterative algorithm alternating between updating Z and updating $\mu^{(1)}, \dots, \mu^{(K)}$. Also called **Lloyd's algorithm**.
- Is the k-means algorithm solving the k-means problem?

Convergence of K-Means Algorithm

Does the k-means algorithm always converge?

Or can it “oscillate”?

Convergence of K-Means Algorithm

Recall the k-means problem:

$$\begin{aligned} \min_{\mu^{(1)}, \dots, \mu^{(K)}, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|_2^2 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

The k-means algorithm really alternates the following two steps:

- (globally) optimize the k-means problem w.r.t. Z , keeping $\mu^{(1)}, \dots, \mu^{(K)}$ fixed
- (globally) optimize the k-means problem w.r.t. μ^1, \dots, μ^K , keeping Z fixed

(I) Optimality of Cluster Assignments

For (currently) fixed $\mu^{(1)}, \dots, \mu^{(K)}$, k-means computes Z as:

```
for i = 1 ... N do
    let  $k^* = \arg \min_k \|x^{(i)} - \mu^{(k)}\|$ 
    set  $z_{i,k} \leftarrow \begin{cases} 1 & \text{if } k = k^* \\ 0 & \text{otherwise} \end{cases}$ 
end
```

This is clearly optimal, since any other choice for $z_{i,k}$ would lead to a larger objective:

$$\sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|_2^2$$

(II) Optimality of Means

For (currently) fixed \mathbf{Z} , k-means computes μ^1, \dots, μ^K as:

```
for k = 1 ... K do
    let  $\mathcal{D}_k = \{\mathbf{x}^{(i)} \mid z_{i,k} = 1\}$ 
     $\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{\mathbf{x} \in \mathcal{D}_k} \mathbf{x}$ 
end
```

Is this minimizing the k-means objective?

$$\min_{\mu^1, \dots, \mu^K} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{D}_k} \|\mathbf{x} - \mu^{(k)}\|_2^2$$

Observation: the objective has the form $\sum_{k=1}^K \mathcal{L}_k$, where

$$\mathcal{L}_k = \sum_{\mathbf{x} \in \mathcal{D}_k} \|\mathbf{x} - \mu^{(k)}\|_2^2$$

depends only on μ_k and we can minimize \mathcal{L}_k **independently** for each k .

(II) Optimality of Means cont'd

Since $\mathcal{L}_k = \sum_{\mathbf{x} \in \mathcal{D}_k} \sum_{d=1}^D \left(x_d - \mu_d^{(k)} \right)^2$, the gradient w.r.t. $\boldsymbol{\mu}^{(k)}$ is

$$\frac{\partial \mathcal{L}_k}{\partial \mu_d^{(k)}} = \sum_{\mathbf{x} \in \mathcal{D}_k} 2(x_d - \mu_d^{(k)}) \quad \nabla_{\boldsymbol{\mu}^{(k)}} \mathcal{L}_k = \sum_{\mathbf{x} \in \mathcal{D}_k} 2(\mathbf{x} - \boldsymbol{\mu}^{(k)})$$

Setting the gradient to zero yields indeed the mean:

$$\sum_{\mathbf{x} \in \mathcal{D}_k} 2(\mathbf{x} - \boldsymbol{\mu}^{(k)}) = 0$$

$$\sum_{\mathbf{x} \in \mathcal{D}_k} \boldsymbol{\mu}^{(k)} = \sum_{\mathbf{x} \in \mathcal{D}_k} \mathbf{x}$$

$$|\mathcal{D}_k| \boldsymbol{\mu}^{(k)} = \sum_{\mathbf{x} \in \mathcal{D}_k} \mathbf{x}$$

$$\boldsymbol{\mu}^{(k)} = \frac{1}{|\mathcal{D}_k|} \sum_{\mathbf{x} \in \mathcal{D}_k} \mathbf{x}$$

Convergence of K-Means Algorithm cont'd

$$\begin{aligned} \min_{\mu^{(1)}, \dots, \mu^{(K)}, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|^2 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

- We have shown that the k-means algorithm alternates between:
 - globally optimizing w.r.t. Z , $\mu^{(1)}, \dots, \mu^{(K)}$ fixed
 - globally optimizing w.r.t. $\mu^{(1)}, \dots, \mu^{(K)}$, Z fixed
- The k-means objective **can only decrease** after each step
- Moreover, there are only **finitely many** indicator matrices Z
- We can **strictly** decrease at most as often as there are different Z
- **Thus, the k-means algorithm converges in terms of the WCSS objective**

Global Optimum?

So, k-means performs **globally optimal sub-steps**,

- minimizing Z
- minimizing $\{\mu^{(1)}, \dots, \mu^{(K)}\}$

in alternating fashion. This is known as **block-coordinate descent** in “blocks” Z and $\{\mu^{(1)}, \dots, \mu^{(K)}\}$. Thus, does the k-means algorithm globally solve the k-means problem?

Global Optimum?

So, k-means performs **globally optimal sub-steps**,

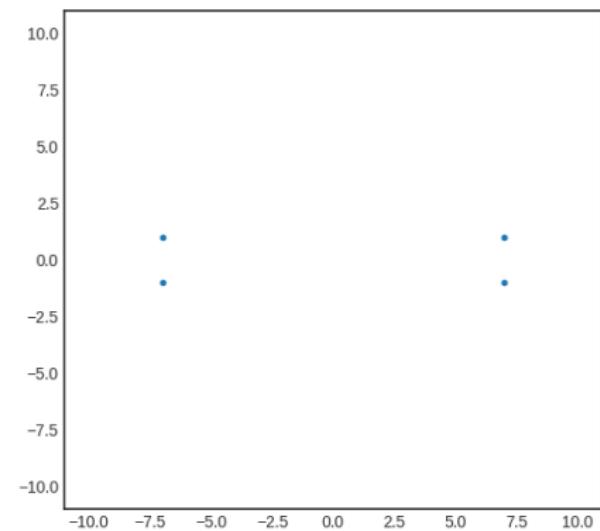
- minimizing Z
- minimizing $\{\mu^{(1)}, \dots, \mu^{(K)}\}$

in alternating fashion. This is known as **block-coordinate descent** in “blocks” Z and $\{\mu^{(1)}, \dots, \mu^{(K)}\}$. Thus, does the k-means algorithm globally solve the k-means problem?

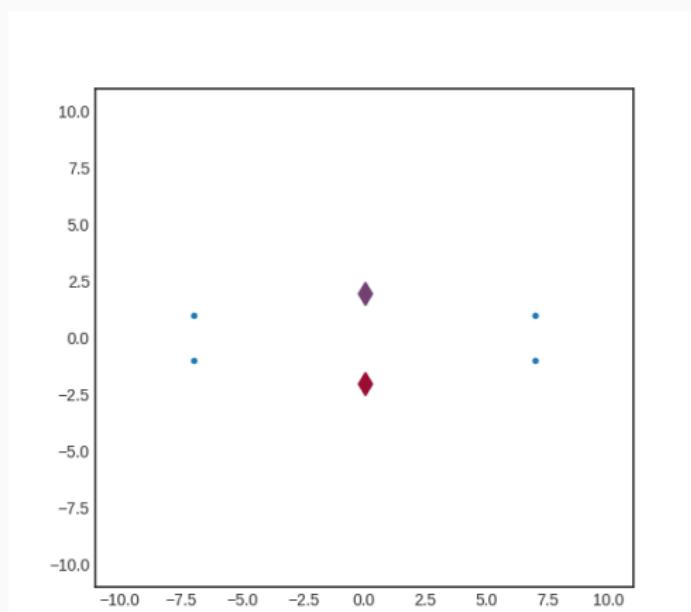
No!

- The k-means algorithm typically converges to a sub-optimal solution!
- The actual solution depends on the initialization and can in theory be arbitrarily bad
- Solving the k-means problem can be shown to be **NP-hard!**

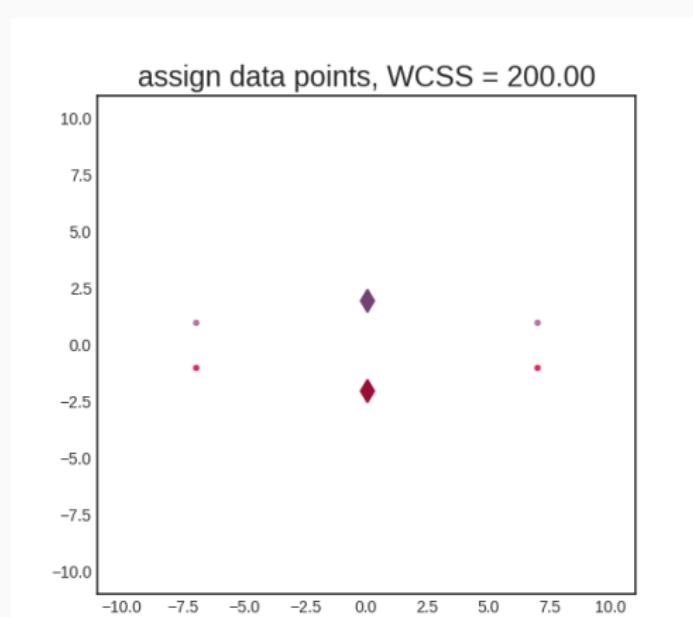
Assume the data set below and assume we set $K = 2$.



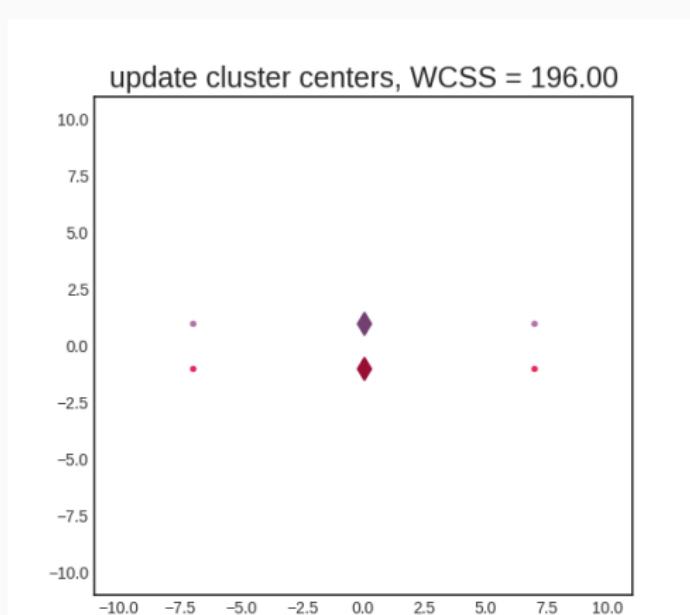
Say the cluster centers are initialized to lie on the y -axis.



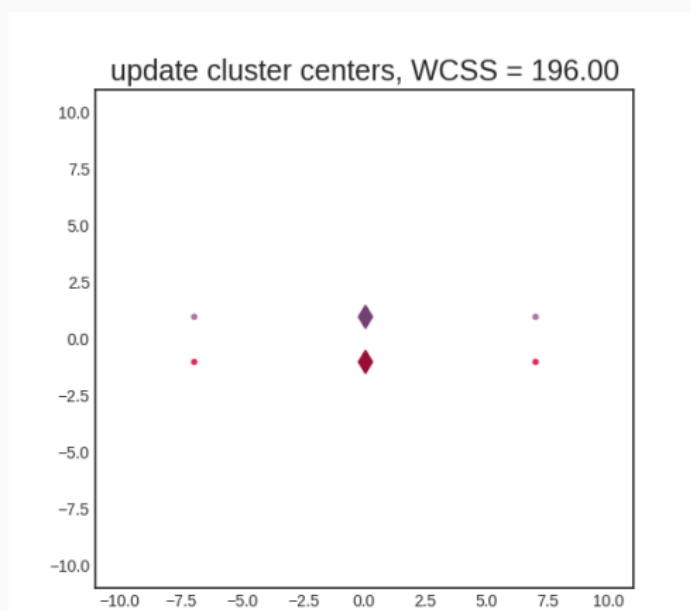
Assign data points to clusters...



Update cluster centers . . .



K-means has already converged, to an evidently a sub-optimal result. By making the rectangle described by the data points wider, the result can be made arbitrarily bad.



Avoiding Bad Clustering

- Random re-starts, i.e., run k-means multiple times with different random initializations
- Pick the clustering of the best run, i.e., with the lowest final objective
- Use clever initializations, like k-means++

K-Means++

- K-Means++ is standard k-means with a clever initialization, proposed by Arthur & Vassilvitskii, 2007
- Basic idea: use data points as initial cluster seeds $\mu^{(1)}, \dots, \mu^{(K)}$, spreaded over the whole dataset
- $\mu^{(1)}$ is selected uniform at random from \mathcal{D}
- For $k > 1$:
 - compute the distance d_i between any $x^{(i)}$ and the closest $\mu^{(l)}$, for $l < k$
 - select $\mu^{(k)} = x^{(i)}$ with probability proportional to d_i^2
- Thus, samples which are far from already selected seeds are preferred

K-Means++ cont'd

- The selection of initial seeds adds to the overall computational time
- However, the authors observed that with their initialization the main phase converges much faster, so that in total k-means++ is usually faster than vanilla k-means
- They typically observed much better clustering results than for vanilla k-means
- Moreover, in expectation over the randomness used in its initialization, k-means++ is guaranteed to achieve an objective which is only a multiplicative factor of $\mathcal{O}(\log K)$ worse than the global optimal solution!
- K-means++ is the de-facto standard implementation of k-means in many libraries

Extending K-Means to Other Distances

- We can adapt k-means to distances other than Euclidean distance
- For example, we could formulate an ℓ_1 -version of k-means:

$$\begin{aligned} \min_{\mu^{(1)}, \dots, \mu^{(K)}, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}^{(k)} \right\|_1 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

- We can easily adapt the k-means algorithm to this case
 - optimizing Z : $\mathbf{x}^{(i)}$ is assigned to closest $\boldsymbol{\mu}^{(k)}$ in ℓ_1 -sense
 - $\boldsymbol{\mu}^{(k)}$ are updated by using the **median** instead of the mean:

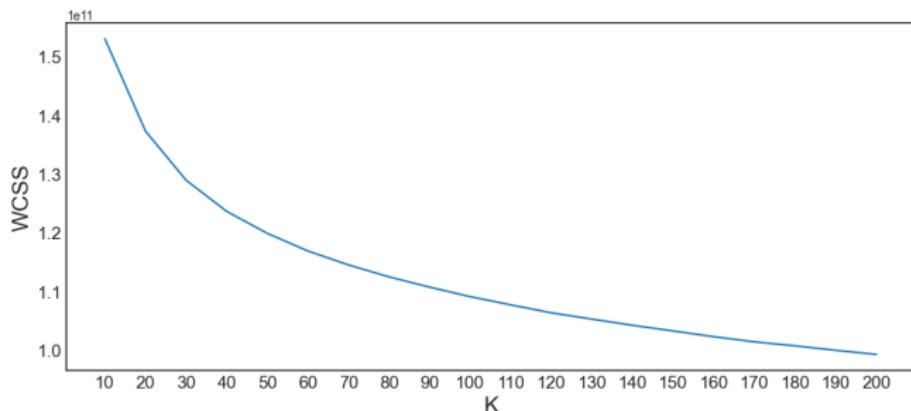
$$\boldsymbol{\mu}_d^{(k)} \leftarrow median(\mathcal{D}_k)$$

where median acts dimension-wise

- This method is known as **k-median algorithm**

Selecting K

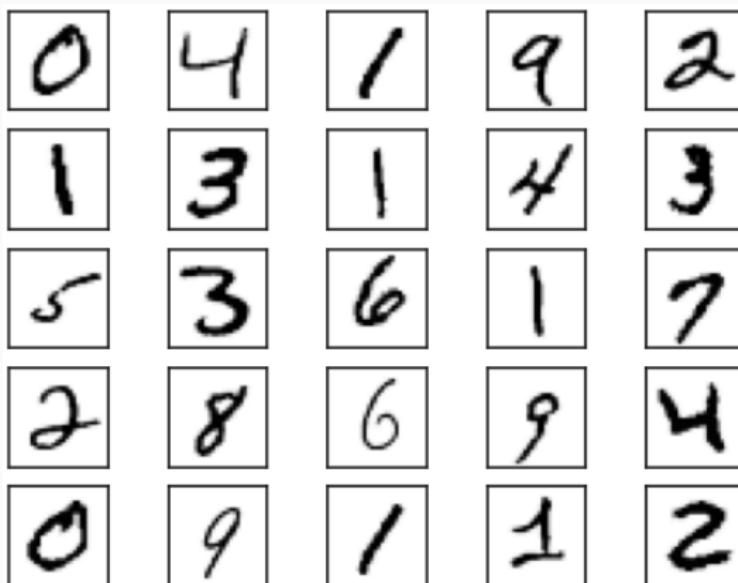
- The k-means algorithm assumes that K is given
- A simple method to find a suitable K is to run k-means for a range of different K 's
- For each K , we record the final objective WCSS
- Plotting the objective against the range of K allows us to “visually select” a suitable K



Application: Data Exploration, Data Summary

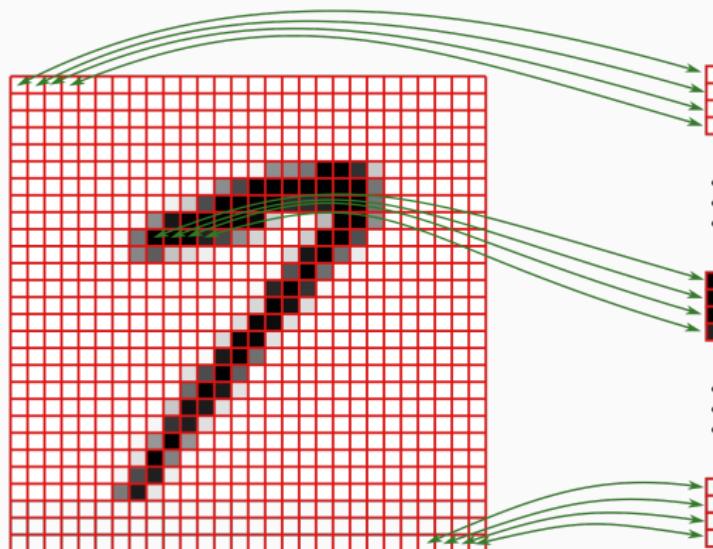
- Data sets can be massive nowadays
- Clustering algorithms like k-means can extract a few exemplars or prototypes, representing whole clusters in the data

- collection of 60,000 images of handwritten digits
- images are 28×28 pixels, with 256 gray scales
- white = 0, black = 255

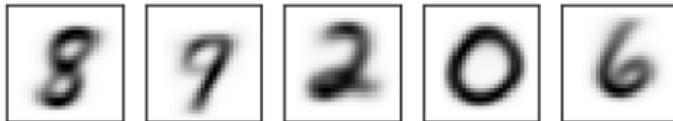
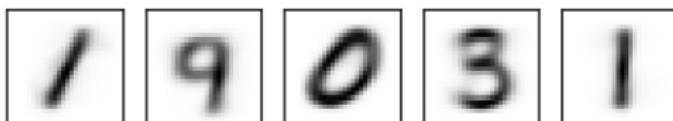


28 × 28 image

784-dim. vector



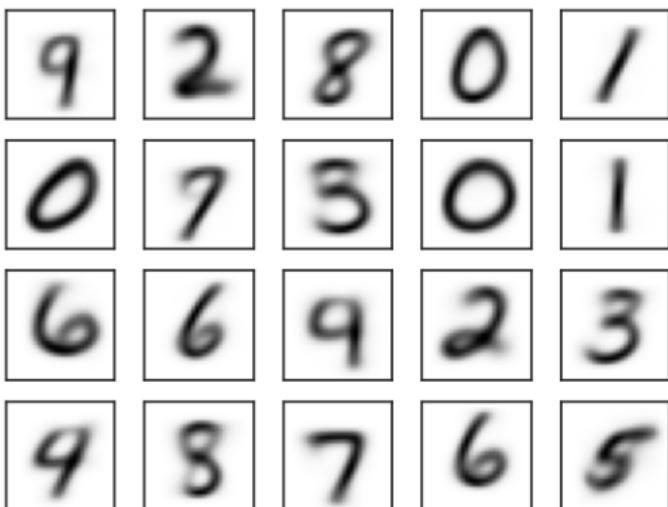
Clustering with $K = 10$, using k-means++ and 10 restarts, delivers following centers $\mu^{(1)}, \dots, \mu^{(10)}$, which summarize the whole dataset. We see that not all “natural classes” are represented.



Plotting some examples from each cluster. First column corresponds to $\mu^{(k)}$, columns 2 – 8 to examples from the cluster

1	1	5	1	7	5	4	1
9	4	4	4	7	7	4	4
0	0	0	0	0	0	0	0
3	3	3	3	9	3	3	2
1	1	1	1	1	3	1	1
8	5	3	8	8	8	8	5
9	9	4	7	9	9	2	9
2	2	2	2	2	2	2	2
0	0	0	0	0	0	0	0
6	6	6	6	6	6	6	4

Same as before, but with $K = 20$. Now all “natural classes” are represented (with duplicates).



Clustering MNIST cont'd

Example

Plotting some examples from each cluster.



Application: Image Compression

Example

Consider following color image, of size 683×1024 pixels:





- Each pixel has 3 channels, red, green, and blue
- Each channel is encoded with 1 byte, yielding $16.7M$ colors
- In uncompressed format we require $683 \times 1024 \times 3 = 2MB$ memory to store this image
- We can apply k-means as a simple lossy compression algorithm
- This technique is also known as **vector quantization**

Image Compression with K-Means (Vector Quantization)

- Interpret each pixel as 3-dimensional data point
- For the previous example image, this yields a “data set” of $683 \times 1024 = 699392$ data points
- Apply k-means with some K , yielding $\mu^{(1)}, \dots, \mu^{(K)}$
- Each $\mu^{(k)}$ is a 3-dimensional vector representing a color, thus $\mu^{(1)}, \dots, \mu^{(K)}$ can be seen as a “learned color palette” for the input image
- Now, each pixel value is replaced by its closest center $\mu^{(k)}$
- Crucially, for each pixel, we just need to store the index k of the closest $\mu^{(k)}$, which can be done with $\log(K)$ bits
- Additionally, we also need to store $\mu^{(1)}, \dots, \mu^{(K)}$, which can be done with $K \times 3 \times 8$ bits

Image Compression $K = 2$

Example



Compression ratio = 24, signal-to-noise-ratio = 10.99 dB

Image Compression $K = 16$

Example



Compression ratio = 6, signal-to-noise-ratio = 19.93 dB

Image Compression $K = 64$

Example



Compression ratio = 4, signal-to-noise-ratio = 25.12 dB



Compression ratio = 3.42, signal-to-noise-ratio = 27.35 dB

- **Clustering**: group similar objects
- Not well-defined, many choices to define clustering
- Data exploration, data summary, data compression
- **K-means problem**: find K mean vectors, which specify clusters as set of points which are closest to them
- Objective: **Within-Cluster Sum-of-Squared distances (WCSS)**
- NP-hard to optimize
- **K-mean algorithm**: alternate between
 - optimal assignment of points to cluster centers
 - taking means of cluster
- Reduces the objective in each step, but converges to a sub-optimal solution
- Re-starts, k-means++

Please Evaluate

ML1 (VO) is nominated for an Award for Excellence in Teaching

