

Clustering: K -means and Kernel K -means

Piyush Rai

Machine Learning (CS771A)

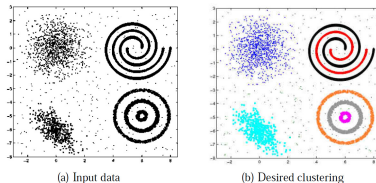
Aug 31, 2016

Clustering

- Usually an **unsupervised learning** problem
- Given: N **unlabeled** examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; no. of desired partitions K
- Goal: Group the examples into K “homogeneous” partitions

Clustering

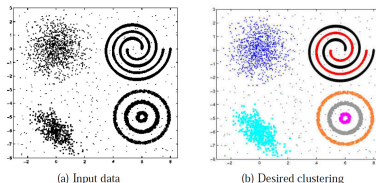
- Usually an **unsupervised learning** problem
- Given: N **unlabeled** examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; no. of desired partitions K
- Goal: Group the examples into K “homogeneous” partitions



Picture courtesy: “Data Clustering: 50 Years Beyond K-Means”, A.K. Jain (2008)

Clustering

- Usually an **unsupervised learning** problem
- Given: N **unlabeled** examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; no. of desired partitions K
- Goal: Group the examples into K “homogeneous” partitions

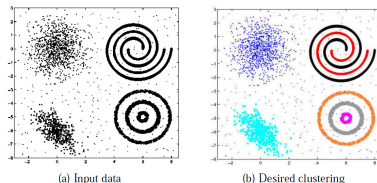


Picture courtesy: “Data Clustering: 50 Years Beyond K-Means”, A.K. Jain (2008)

- Loosely speaking, it is classification without ground truth labels

Clustering

- Usually an **unsupervised learning** problem
- Given: N **unlabeled** examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; no. of desired partitions K
- Goal: Group the examples into K “homogeneous” partitions



Picture courtesy: “Data Clustering: 50 Years Beyond K-Means”, A.K. Jain (2008)

- Loosely speaking, it is classification without ground truth labels
- A good clustering is one that achieves:
 - **High within-cluster similarity**
 - **Low inter-cluster similarity**

Similarity can be Subjective

- Clustering only looks at similarities, no labels are given
- **Without labels**, similarity can be hard to define



Similarity can be Subjective

- Clustering only looks at similarities, no labels are given
- **Without labels**, similarity can be hard to define



- Thus using the right distance/similarity is very important in clustering

Similarity can be Subjective

- Clustering only looks at similarities, no labels are given
- Without labels, similarity can be hard to define



- Thus using the right distance/similarity is very important in clustering
- Also important to define/ask: “Clustering based on what”?

Similarity can be Subjective

- Clustering only looks at similarities, no labels are given
- Without labels, similarity can be hard to define



- Thus using the right distance/similarity is very important in clustering
- Also important to define/ask: “Clustering based on what”?

Clustering: Some Examples

- Document/Image/Webpage Clustering
- Image Segmentation (clustering pixels)



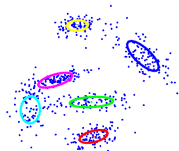
- Clustering web-search results
- Clustering (people) nodes in (social) networks/graphs
- .. and many more..

Picture courtesy: <http://people.cs.uchicago.edu/~pff/segment/>

Types of Clustering

1 Flat or Partitional clustering

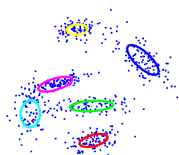
- Partitions are independent of each other



Types of Clustering

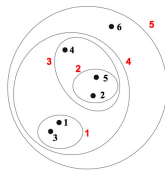
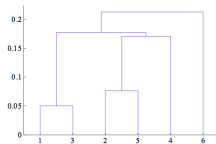
1 Flat or Partitional clustering

- Partitions are independent of each other



2 Hierarchical clustering

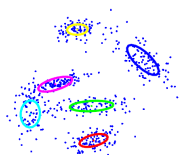
- Partitions can be visualized using a tree structure (a dendrogram)



Types of Clustering

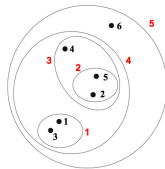
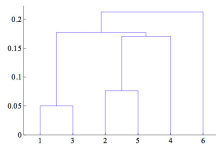
1 Flat or Partitional clustering

- Partitions are **independent of each other**



2 Hierarchical clustering

- Partitions can be visualized using a tree structure (a **dendrogram**)



- Possible to view partitions at **different levels of granularities** (i.e., can refine/coarsen clusters) using different K

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**
 - (Re)-Assign each example \mathbf{x}_n to its closest cluster center (based on the smallest Euclidean distance)

$$\mathcal{C}_k = \{n : k = \arg \min_k \|\mathbf{x}_n - \mu_k\|^2\}$$

(\mathcal{C}_k is the set of examples assigned to cluster k with center μ_k)

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**
 - (Re)-Assign each example \mathbf{x}_n to its closest cluster center (based on the smallest Euclidean distance)

$$\mathcal{C}_k = \{n : k = \arg \min_k \|\mathbf{x}_n - \mu_k\|^2\}$$

(\mathcal{C}_k is the set of examples assigned to cluster k with center μ_k)

- Update the cluster means

$$\mu_k = \text{mean}(\mathcal{C}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**
 - (Re)-Assign each example \mathbf{x}_n to its closest cluster center (based on the smallest Euclidean distance)

$$\mathcal{C}_k = \{n : k = \arg \min_k \|\mathbf{x}_n - \mu_k\|^2\}$$

(\mathcal{C}_k is the set of examples assigned to cluster k with center μ_k)

- **Update** the cluster means

$$\mu_k = \text{mean}(\mathcal{C}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

- **Repeat** while not converged

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**
 - (Re)-Assign each example \mathbf{x}_n to its closest cluster center (based on the smallest Euclidean distance)

$$\mathcal{C}_k = \{n : k = \arg \min_k \|\mathbf{x}_n - \mu_k\|^2\}$$

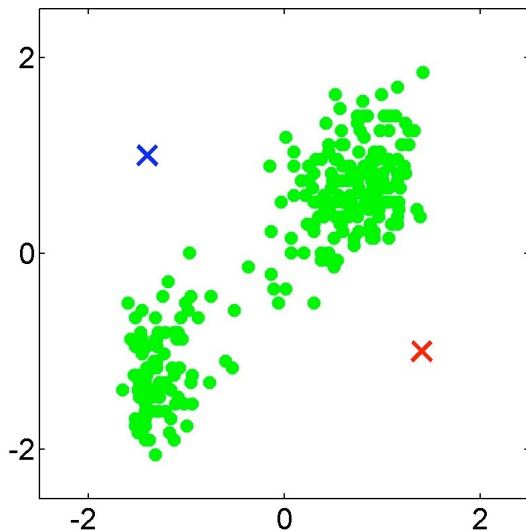
(\mathcal{C}_k is the set of examples assigned to cluster k with center μ_k)

- Update the cluster means

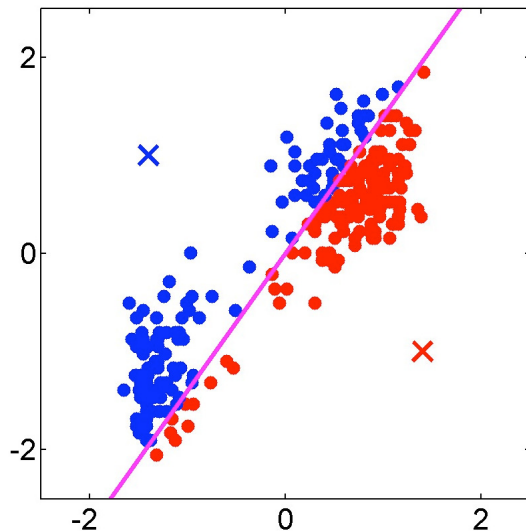
$$\mu_k = \text{mean}(\mathcal{C}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

- Repeat while not converged
- Stop when cluster means or the “loss” does not change by much

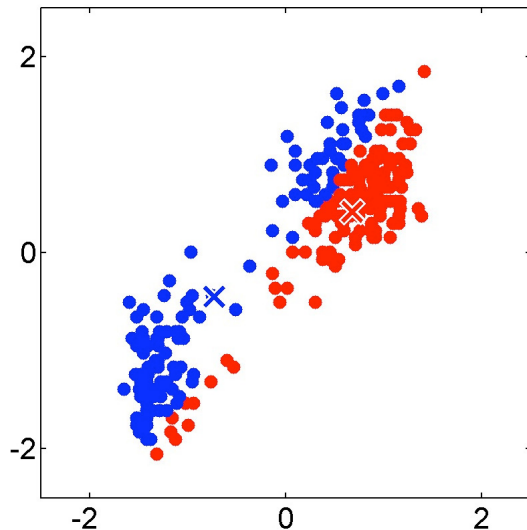
K -means: Initialization (assume $K = 2$)



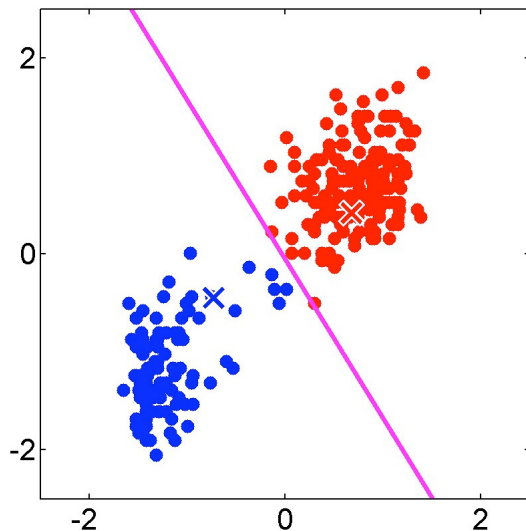
K -means iteration 1: Assigning points



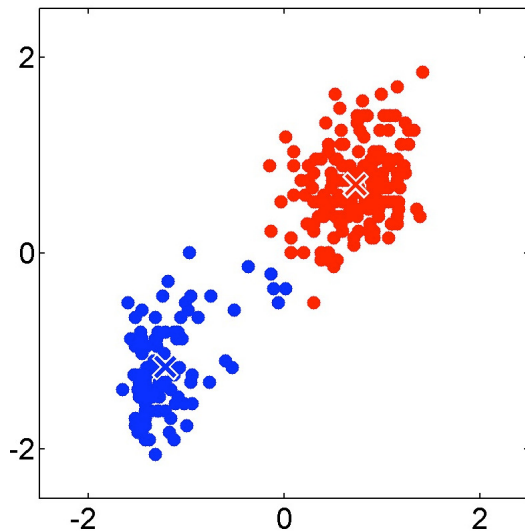
K -means iteration 1: Recomputing the centers



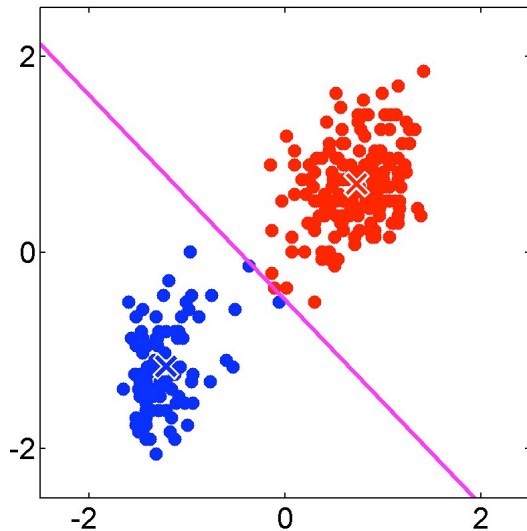
K -means iteration 2: Assigning points



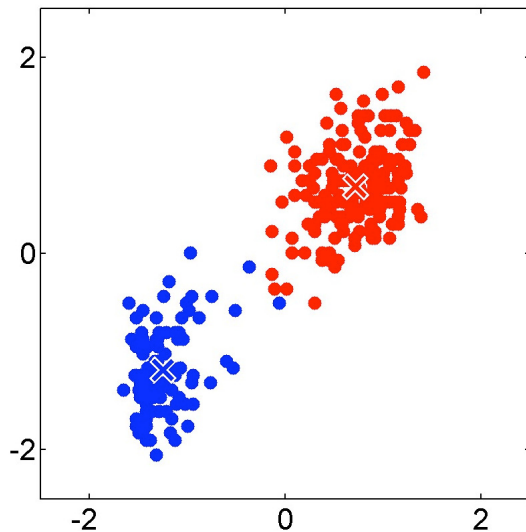
K -means iteration 2: Recomputing the centers



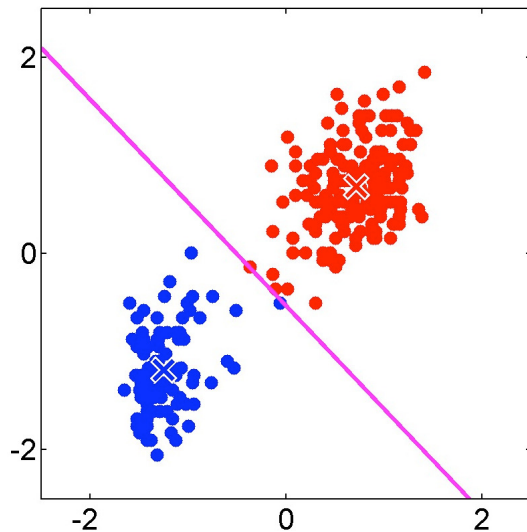
K -means iteration 3: Assigning points



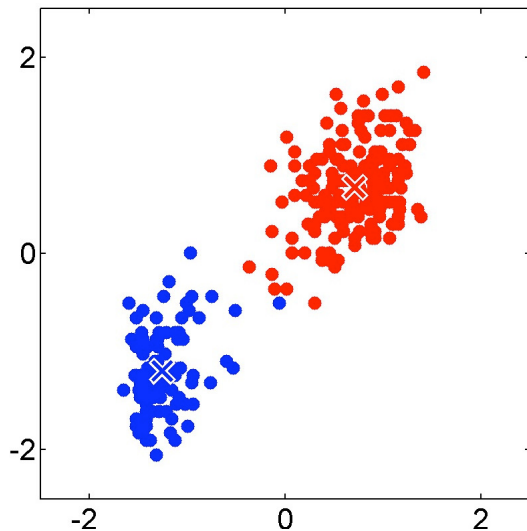
K -means iteration 3: Recomputing the centers



K -means iteration 4: Assigning points



K -means iteration 4: Recomputing the centers



What Loss Function is K -means Optimizing?

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K **one-hot encoding** of \mathbf{x}_n

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K one-hot encoding of \mathbf{x}_n
- Define the distortion or “loss” for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K one-hot encoding of \mathbf{x}_n
- Define the distortion or “loss” for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

- Total distortion over all points defines the K -means “loss function”

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K **one-hot encoding** of \mathbf{x}_n
- Define the **distortion** or **“loss”** for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

- Total distortion over all points defines the K -means “loss function”

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2 = \|\mathbf{X} - \mathbf{Z}\mu\|^2$$

where \mathbf{Z} is $N \times K$ (row n is \mathbf{z}_n) and μ is $K \times D$ (row k is μ_k)

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K **one-hot encoding** of \mathbf{x}_n
- Define the **distortion** or **“loss”** for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

- Total distortion over all points defines the K -means “loss function”

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2 = \|\mathbf{X} - \mathbf{Z}\mu\|^2$$

where \mathbf{Z} is $N \times K$ (row n is \mathbf{z}_n) and μ is $K \times D$ (row k is μ_k)

- The K -means **problem** is to minimize this objective w.r.t. μ and \mathbf{Z}

What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K **one-hot encoding** of \mathbf{x}_n
- Define the **distortion** or **“loss”** for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

- Total distortion over all points defines the K -means “loss function”

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2 = \|\mathbf{X} - \mathbf{Z}\mu\|^2$$

where \mathbf{Z} is $N \times K$ (row n is \mathbf{z}_n) and μ is $K \times D$ (row k is μ_k)

- The K -means **problem** is to minimize this objective w.r.t. μ and \mathbf{Z}
 - Note that the objective only minimizes **within-cluster distortions**

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)
- The **K -means algorithm** we saw is a heuristic to optimize this function

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)
- The **K -means algorithm** we saw is a heuristic to optimize this function
- K -means algorithm alternated between the following two steps

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)
- The **K -means algorithm** we saw is a heuristic to optimize this function
- K -means algorithm alternated between the following two steps
 - Fix $\boldsymbol{\mu}$, minimize w.r.t. \mathbf{Z} (**assign points to closest centers**)

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)
- The **K -means algorithm** we saw is a heuristic to optimize this function
- K -means algorithm alternated between the following two steps
 - Fix $\boldsymbol{\mu}$, minimize w.r.t. \mathbf{Z} (assign points to closest centers)
 - Fix \mathbf{Z} , minimize w.r.t. $\boldsymbol{\mu}$ (recompute the center means)

K-means Objective

- Consider the K -means objective function

$$L(\boldsymbol{\mu}, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- It is a **non-convex** objective function
 - Many local minima possible
- Also **NP-hard** to minimize in general (note that \mathbf{Z} is discrete)
- The **K -means algorithm** we saw is a heuristic to optimize this function
- K -means algorithm alternated between the following two steps
 - Fix $\boldsymbol{\mu}$, minimize w.r.t. \mathbf{Z} (**assign points to closest centers**)
 - Fix \mathbf{Z} , minimize w.r.t. $\boldsymbol{\mu}$ (**recompute the center means**)
- Note: The algorithm usually converges to a local minima (though may not always, and it may just converge “somewhere”). Multiple runs with different initializations can be tried to find a good solution.

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or μ) can never increase the objective

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or $\boldsymbol{\mu}$) can never increase the objective
- When we update \mathbf{Z} from $\mathbf{Z}^{(t-1)}$ to $\mathbf{Z}^{(t)}$

$$L(\boldsymbol{\mu}^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\boldsymbol{\mu}^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t-1)})$$

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or $\boldsymbol{\mu}$) can never increase the objective
- When we update \mathbf{Z} from $\mathbf{Z}^{(t-1)}$ to $\mathbf{Z}^{(t)}$

$$L(\boldsymbol{\mu}^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\boldsymbol{\mu}^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t-1)})$$

because the new $\mathbf{Z}^{(t)} = \arg \min_{\mathbf{Z}} L(\boldsymbol{\mu}^{(t-1)}, \mathbf{X}, \mathbf{Z})$

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or μ) can never increase the objective
- When we update \mathbf{Z} from $\mathbf{Z}^{(t-1)}$ to $\mathbf{Z}^{(t)}$

$$L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t-1)})$$

because the new $\mathbf{Z}^{(t)} = \arg \min_{\mathbf{Z}} L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z})$

- When we update μ from $\mu^{(t-1)}$ to $\mu^{(t)}$

$$L(\mu^{(t)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)})$$

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or μ) can never increase the objective
- When we update \mathbf{Z} from $\mathbf{Z}^{(t-1)}$ to $\mathbf{Z}^{(t)}$

$$L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t-1)})$$

because the new $\mathbf{Z}^{(t)} = \arg \min_{\mathbf{Z}} L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z})$

- When we update μ from $\mu^{(t-1)}$ to $\mu^{(t)}$

$$L(\mu^{(t)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)})$$

because the new $\mu^{(t)} = \arg \min_{\mu} L(\mu, \mathbf{X}, \mathbf{Z}^{(t)})$

Convergence of K -means Algorithm

- Each step (updating \mathbf{Z} or μ) can never increase the objective
- When we update \mathbf{Z} from $\mathbf{Z}^{(t-1)}$ to $\mathbf{Z}^{(t)}$

$$L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t-1)})$$

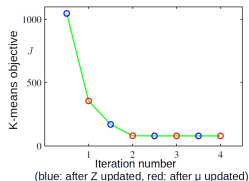
because the new $\mathbf{Z}^{(t)} = \arg \min_{\mathbf{Z}} L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z})$

- When we update μ from $\mu^{(t-1)}$ to $\mu^{(t)}$

$$L(\mu^{(t)}, \mathbf{X}, \mathbf{Z}^{(t)}) \leq L(\mu^{(t-1)}, \mathbf{X}, \mathbf{Z}^{(t)})$$

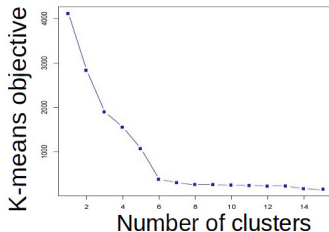
because the new $\mu^{(t)} = \arg \min_{\mu} L(\mu, \mathbf{X}, \mathbf{Z}^{(t)})$

- Thus the K -means algorithm monotonically decreases the objective



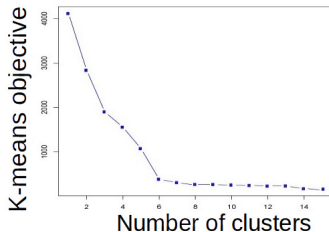
K-means: Choosing K

- One way to select K for the K -means algorithm is to try different values of K , plot the K -means objective versus K , and look at the “elbow-point”



K-means: Choosing K

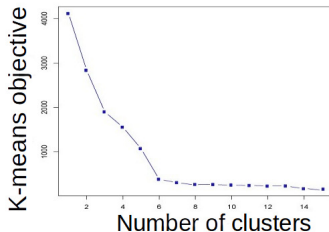
- One way to select K for the K -means algorithm is to try different values of K , plot the K -means objective versus K , and look at the “elbow-point”



- For the above plot, $K = 6$ is the elbow point

K-means: Choosing K

- One way to select K for the K -means algorithm is to try different values of K , plot the K -means objective versus K , and look at the “elbow-point”



- For the above plot, $K = 6$ is the elbow point
- Can also use information criterion such as AIC (Akaike Information Criterion)

$$AIC = 2L(\hat{\mu}, \mathbf{X}, \hat{\mathbf{Z}}) + K \log D$$

.. and choose the K that has the smallest AIC (discourages large K)

K -means: Some Limitations

- Makes **hard assignments** of points to clusters

K -means: Some Limitations

- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all

K-means: Some Limitations

- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all
 - No notion of a **soft assignment** (i.e., **probability** of being assigned to each cluster: say $K = 3$ and for some point x_n , $p_1 = 0.7$, $p_2 = 0.2$, $p_3 = 0.1$)

K-means: Some Limitations

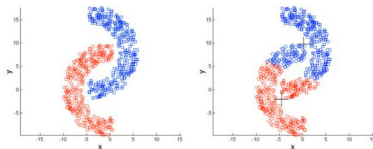
- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all
 - No notion of a **soft assignment** (i.e., **probability** of being assigned to each cluster: say $K = 3$ and for some point x_n , $p_1 = 0.7$, $p_2 = 0.2$, $p_3 = 0.1$)
- Works well only if the clusters are **roughly of equal sizes**

K-means: Some Limitations

- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all
 - No notion of a **soft assignment** (i.e., **probability** of being assigned to each cluster: say $K = 3$ and for some point x_n , $p_1 = 0.7$, $p_2 = 0.2$, $p_3 = 0.1$)
- Works well only if the clusters are **roughly of equal sizes**
- Probabilistic clustering methods such as **Gaussian mixture models** can handle both these issues (model each cluster using a Gaussian distribution)

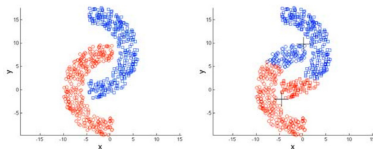
K-means: Some Limitations

- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all
 - No notion of a **soft assignment** (i.e., **probability** of being assigned to each cluster: say $K = 3$ and for some point x_n , $p_1 = 0.7, p_2 = 0.2, p_3 = 0.1$)
- Works well only if the clusters are **roughly of equal sizes**
- Probabilistic clustering methods such as **Gaussian mixture models** can handle both these issues (model each cluster using a Gaussian distribution)
- K-means also works well only when the clusters are **round-shaped** and does badly if the clusters have **non-convex shapes**



K-means: Some Limitations

- Makes **hard assignments** of points to clusters
 - A point either completely belongs to a cluster or doesn't belong at all
 - No notion of a **soft assignment** (i.e., **probability** of being assigned to each cluster: say $K = 3$ and for some point x_n , $p_1 = 0.7$, $p_2 = 0.2$, $p_3 = 0.1$)
- Works well only if the clusters are **roughly of equal sizes**
- Probabilistic clustering methods such as **Gaussian mixture models** can handle both these issues (model each cluster using a Gaussian distribution)
- K-means also works well only when the clusters are **round-shaped** and does badly if the clusters have **non-convex shapes**



- **Kernel K-means** or **Spectral clustering** can handle non-convex

Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2$$

Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2 = \|\phi(\mathbf{x}_n)\|^2 + \|\phi(\boldsymbol{\mu}_k)\|^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k)$$

Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\begin{aligned}\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2 &= \|\phi(\mathbf{x}_n)\|^2 + \|\phi(\boldsymbol{\mu}_k)\|^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k) \\ &= k(\mathbf{x}_n, \mathbf{x}_n) + k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k) - 2k(\mathbf{x}_n, \boldsymbol{\mu}_k)\end{aligned}$$

- Here $k(.,.)$ denotes the kernel function and ϕ is its (implicit) feature map

Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\begin{aligned}\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2 &= \|\phi(\mathbf{x}_n)\|^2 + \|\phi(\boldsymbol{\mu}_k)\|^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k) \\ &= k(\mathbf{x}_n, \mathbf{x}_n) + k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k) - 2k(\mathbf{x}_n, \boldsymbol{\mu}_k)\end{aligned}$$

- Here $k(.,.)$ denotes the kernel function and ϕ is its (implicit) feature map
- Note: ϕ doesn't have to be computed/stored for data $\{\mathbf{x}_n\}_{n=1}^N$ or the cluster means $\{\boldsymbol{\mu}_k\}_{k=1}^K$ because computations only depend on kernel evaluations

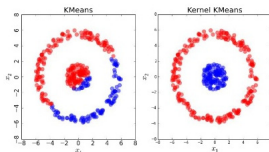
Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\begin{aligned}\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2 &= \|\phi(\mathbf{x}_n)\|^2 + \|\phi(\boldsymbol{\mu}_k)\|^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k) \\ &= k(\mathbf{x}_n, \mathbf{x}_n) + k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k) - 2k(\mathbf{x}_n, \boldsymbol{\mu}_k)\end{aligned}$$

- Here $k(.,.)$ denotes the kernel function and ϕ is its (implicit) feature map
- Note: ϕ doesn't have to be computed/stored for data $\{\mathbf{x}_n\}_{n=1}^N$ or the cluster means $\{\boldsymbol{\mu}_k\}_{k=1}^K$ because computations only depend on kernel evaluations

Kernel K-means vs. K-means



Pyclus: Open Source Data Clustering Package

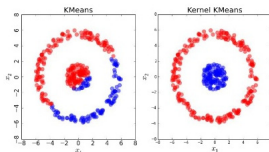
Kernel K -means

- **Basic idea:** Replace the Euclidean distance/similarity computations in K -means by the **kernelized versions**. E.g., $d(\mathbf{x}_n, \boldsymbol{\mu}_k) = \|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|$ by

$$\begin{aligned}\|\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)\|^2 &= \|\phi(\mathbf{x}_n)\|^2 + \|\phi(\boldsymbol{\mu}_k)\|^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k) \\ &= k(\mathbf{x}_n, \mathbf{x}_n) + k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k) - 2k(\mathbf{x}_n, \boldsymbol{\mu}_k)\end{aligned}$$

- Here $k(.,.)$ denotes the kernel function and ϕ is its (implicit) feature map
- Note: ϕ doesn't have to be computed/stored for data $\{\mathbf{x}_n\}_{n=1}^N$ or the cluster means $\{\boldsymbol{\mu}_k\}_{k=1}^K$ because computations only depend on kernel evaluations

Kernel K-means vs. K-means



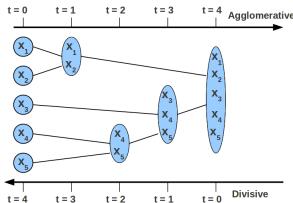
Pycluster: Open Source Data Clustering Package

- **A small technical note:** When computing $k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k)$ and $k(\mathbf{x}_n, \boldsymbol{\mu}_k)$, remember that $\phi(\boldsymbol{\mu}_k)$ is the average of ϕ 's the data points assigned to cluster k

Extra Slides:

Hierarchical Clustering

Hierarchical Clustering



- **Agglomerative (bottom-up) Clustering**

- 1 Start with each example in its own **singleton cluster**
- 2 At each time-step, greedily **merge** 2 most similar clusters
- 3 Stop when there is a single cluster of all examples, else go to 2

- **Divisive (top-down) Clustering**

- 1 Start with all examples in the same cluster
- 2 At each time-step, remove the “outsiders” from the **least cohesive cluster**
- 3 Stop when each example is in its own singleton cluster, else go to 2

- Agglomerative is more popular and simpler than divisive (but less accurate)

(Dis)similarity between clusters

- We know how to compute the dissimilarity $d(\mathbf{x}_i, \mathbf{x}_j)$ between two examples
- How to compute the dissimilarity between two clusters R and S ?
- **Min-link or single-link:** results in chaining (clusters can get very large)

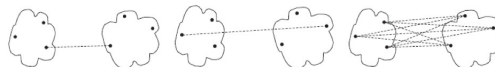
$$d(R, S) = \min_{\mathbf{x}_R \in R, \mathbf{x}_S \in S} d(\mathbf{x}_R, \mathbf{x}_S)$$

- **Max-link or complete-link:** results in small, round shaped clusters

$$d(R, S) = \max_{\mathbf{x}_R \in R, \mathbf{x}_S \in S} d(\mathbf{x}_R, \mathbf{x}_S)$$

- **Average-link:** compromise between single and complete linkage

$$d(R, S) = \frac{1}{|R||S|} \sum_{\mathbf{x}_R \in R, \mathbf{x}_S \in S} d(\mathbf{x}_R, \mathbf{x}_S)$$



(a) MIN (single link.)

(b) MAX (complete link.)

(c) Group average.

Flat vs Hierarchical Clustering

- Flat clustering produces a single partitioning
- Hierarchical Clustering can give different partitionings depending on the level-of-resolution we are looking at
- Flat clustering is usually more efficient run-time wise
- Hierarchical clustering can be slow (has to make several merge/split decisions)