

Bil 470 / YAP 470

Introduction to Machine Learning (Yapay Öğrenme)

Batuhan Bardak

Unsupervised Learning (PCA & K-means)

Date: 14.11.2022

Plan for today

- Dimensionality Reduction
 - PCA
- Unsupervised Learning
 - Hierarchical Clustering
 - K-Means

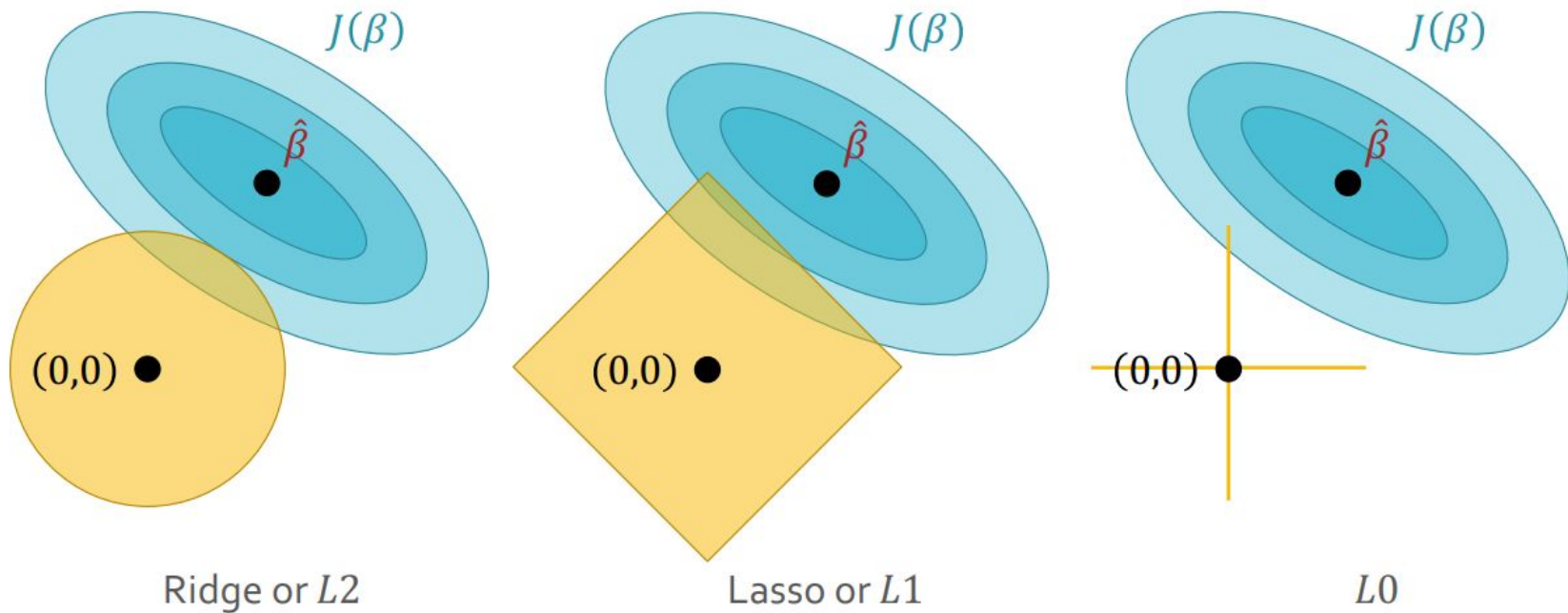
Learning Paradigms

- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$
 - Regression - $y^{(n)} \in \mathbb{R}$
 - Classification - $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$
 - Clustering
 - Dimensionality reduction
- Reinforcement learning - $\mathcal{D} = \{\mathbf{s}^{(n)}, \mathbf{a}^{(n)}, r^{(n)}\}_{n=1}^N$

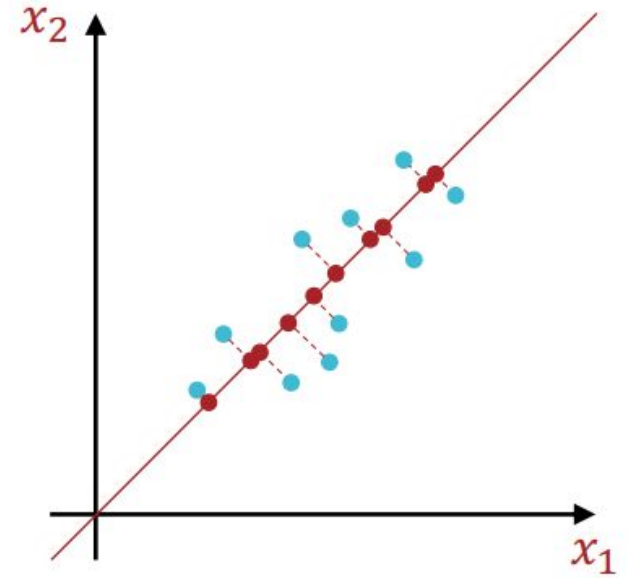
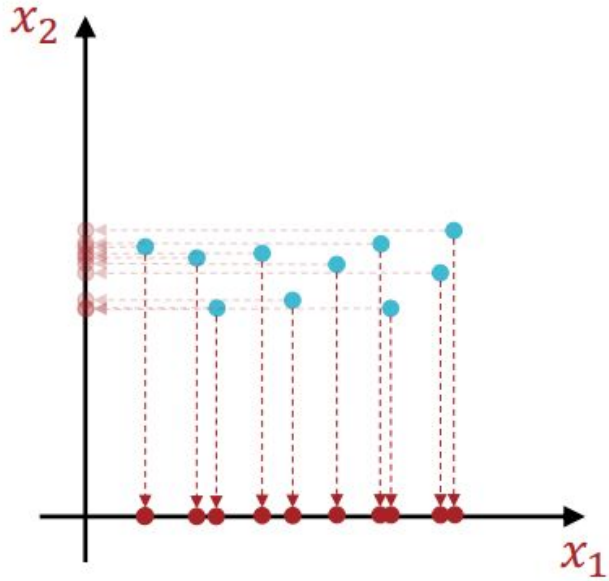
Dimensionality Reduction

- Goal: given some unlabeled data set, learn a latent (typically lower-dimensional) representation
- Use cases:
 - Reducing computational cost (runtime, storage, etc...)
 - Improving generalization
 - Visualizing data
- Applications:
 - High-resolution images/videos
 - Text data
 - Financial or transaction data

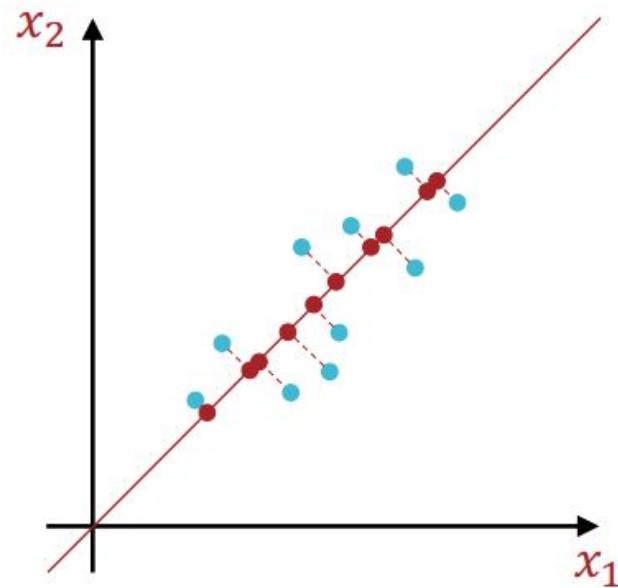
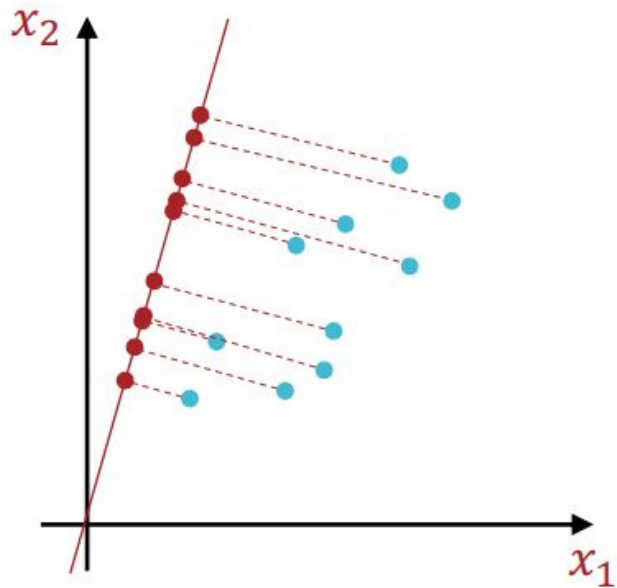
Recall: Regularization (L1)



Feature Elimination



Feature Reduction



Centering the Data

- To be consistent, we will constrain principal components to be *orthogonal unit vectors* that begin at the origin
- Preprocess data to be centered around the origin:

1. $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$

2. $\tilde{\mathbf{x}}^{(n)} = \mathbf{x}^{(n)} - \mu \forall n$

3. $X = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)T} \\ \tilde{\mathbf{x}}^{(2)T} \\ \vdots \\ \tilde{\mathbf{x}}^{(N)T} \end{bmatrix}$

Reconstruction Error

- The projection of $\tilde{\mathbf{x}}^{(n)}$ onto a **unit** vector \mathbf{v} is

$$\mathbf{z}^{(n)} = \left(\frac{\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}}{\|\mathbf{v}\|_2} \right) \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$$

Length of projection

Direction of projection

Reconstruction Error

- The projection of $\tilde{\mathbf{x}}^{(n)}$ onto a unit vector \mathbf{v} is

$$\mathbf{z}^{(n)} = (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}$$

$$\hat{\mathbf{v}} = \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2$$

$$\begin{aligned} & \|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2 \\ &= \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - 2(\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \tilde{\mathbf{x}}^{(n)} + (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \mathbf{v} \\ &= \tilde{\mathbf{x}}^{(n)T} \tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}^T \tilde{\mathbf{x}}^{(n)} \\ &= \|\tilde{\mathbf{x}}^{(n)}\|_2^2 - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \end{aligned}$$

Min. Reconstruction Error \Leftrightarrow Max. the Variance

$$\begin{aligned}\hat{\mathbf{v}} &= \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)} - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)}) \mathbf{v}\|_2^2 \\ &= \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N \|\tilde{\mathbf{x}}^{(n)}\|_2^2 - (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \\ &= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \sum_{n=1}^N (\mathbf{v}^T \tilde{\mathbf{x}}^{(n)})^2 \quad \leftarrow \begin{array}{l} \text{Variance of projections} \\ (\tilde{\mathbf{x}}^{(n)} \text{ are centered}) \end{array} \\ &= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T \left(\sum_{n=1}^N \tilde{\mathbf{x}}^{(n)} \tilde{\mathbf{x}}^{(n)T} \right) \mathbf{v} \\ &= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T (X^T X) \mathbf{v}\end{aligned}$$

Maximizing the Variance

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T (X^T X) \mathbf{v}$$

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \lambda) &= \mathbf{v}^T (X^T X) \mathbf{v} - \lambda (\|\mathbf{v}\|_2^2 - 1) \\ &= \mathbf{v}^T (X^T X) \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - 1) \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = (X^T X) \mathbf{v} - \lambda \mathbf{v}$$

$$\rightarrow (X^T X) \hat{\mathbf{v}} - \lambda \hat{\mathbf{v}} = 0 \rightarrow (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}$$

- $\hat{\mathbf{v}}$ is an eigenvector of $X^T X$ and λ is the corresponding eigenvalue! But which one?

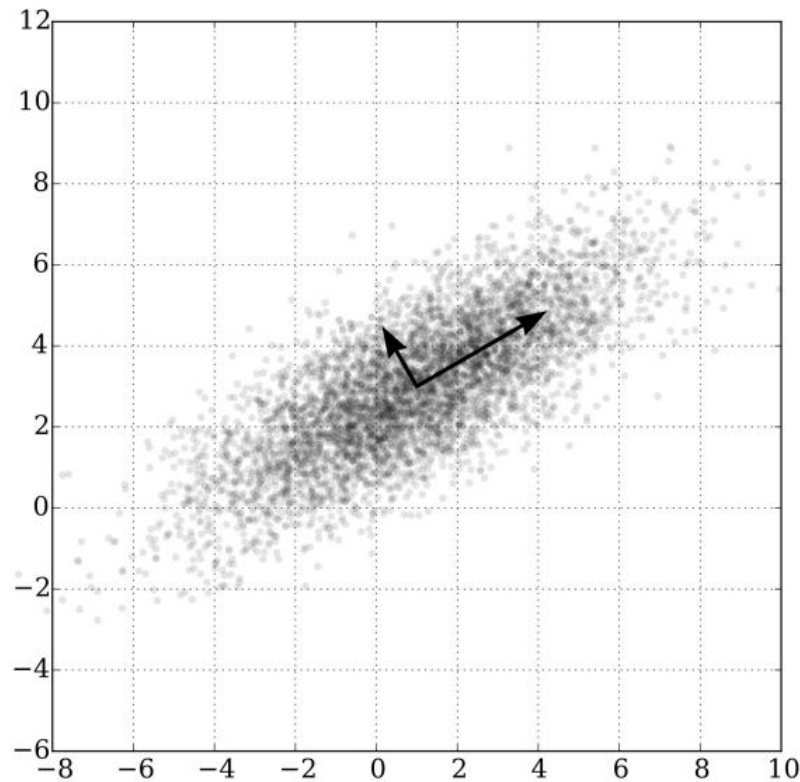
Maximizing the Variance

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|_2^2=1} \mathbf{v}^T (X^T X) \mathbf{v}$$

$$(X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}} \rightarrow \hat{\mathbf{v}}^T (X^T X) \hat{\mathbf{v}} = \lambda \hat{\mathbf{v}}^T \hat{\mathbf{v}} = \lambda$$

- The first principal component is the eigenvector $\hat{\mathbf{v}}_1$ that corresponds to the largest eigenvalue λ_1
- The second principal component is the eigenvector $\hat{\mathbf{v}}_2$ that corresponds to the second largest eigenvalue λ_2
 - $\hat{\mathbf{v}}_1$ and $\hat{\mathbf{v}}_2$ are orthogonal
- Etc ...
- λ_i is a measure of how much variance falls along $\hat{\mathbf{v}}_i$

How can we efficiently find principle components (eigenvectors)?



Singular Value Decomposition (SVD) for PCA

- Every real-valued matrix $X \in \mathbb{R}^{N \times D}$ can be expressed as

$$X = USV^T$$

where:

1. $U \in \mathbb{R}^{N \times N}$ - columns of U are eigenvectors of XX^T
2. $V \in \mathbb{R}^{D \times D}$ - columns of V are eigenvectors of X^TX
3. $S \in \mathbb{R}^{N \times D}$ - diagonal matrix whose entries are the eigenvalues of $X \rightarrow$ squared entries are the eigenvalues of XX^T and X^TX

PCA Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, \rho$
 1. Center the data
 2. Use SVD to compute the eigenvalues and eigenvectors of $\mathbf{X}^T \mathbf{X}$
 3. Collect the top ρ eigenvectors (corresponding to the ρ largest eigenvalues), $\mathbf{V}_\rho \in \mathbb{R}^{D \times \rho}$
 4. Project the data into the space defined by \mathbf{V}_ρ , $\mathbf{Z} = \mathbf{X} \mathbf{V}_\rho$
- Output: \mathbf{Z} , the transformed (potentially lower-dimensional) data

How many PCs should we use?

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, \rho$
 1. Center the data
 2. Use SVD to compute the eigenvalues and eigenvectors of $X^T X$
 3. Collect the top ρ eigenvectors (corresponding to the ρ largest eigenvalues), $V_\rho \in \mathbb{R}^{D \times \rho}$
 4. Project the data into the space defined by V_ρ , $Z = XV_\rho$
- Output: Z , the transformed (potentially lower-dimensional) data

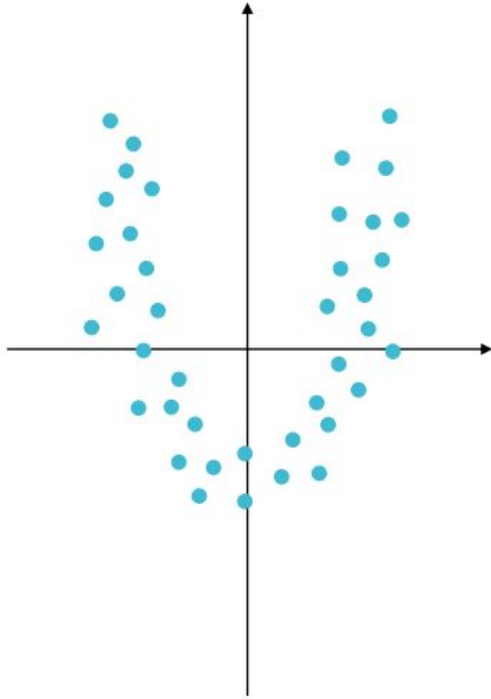
Choosing the number of PCs

- Define a percentage of explained variance for the i^{th} PC:

$$\lambda_i / \sum \lambda_j$$

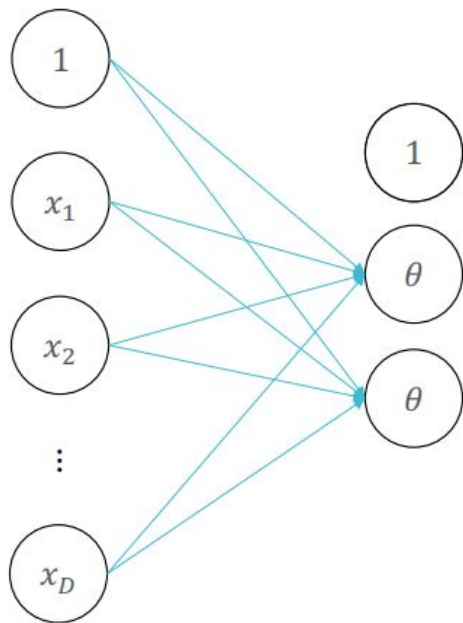
- Select all PCs above some threshold of explained variance, e.g., 5%
- Keep selecting PCs until the total explained variance exceeds some threshold, e.g., 90%
- Evaluate on some downstream metric

Shortcomings of PCA



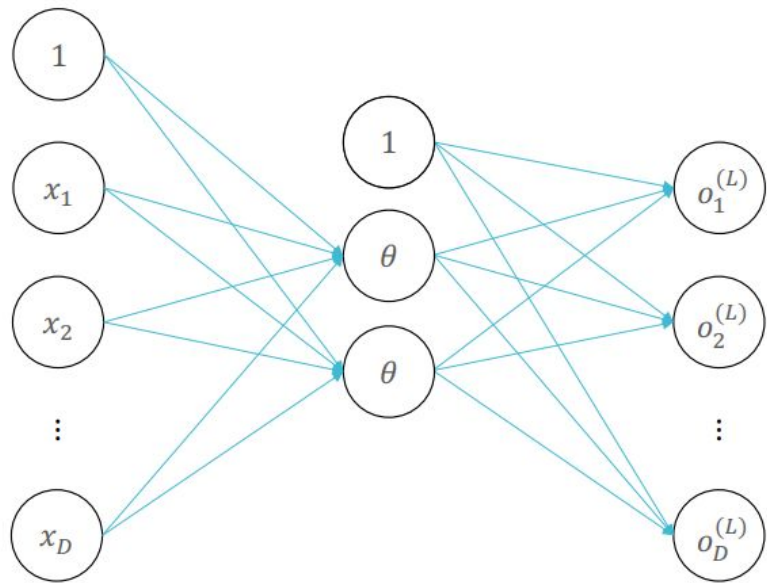
- Principal components are orthogonal (unit) vectors
- Principal components can be expressed as linear combinations of the data

Autoencoders



Insight: neural networks implicitly learn low-dimensional representations of inputs in hidden layers

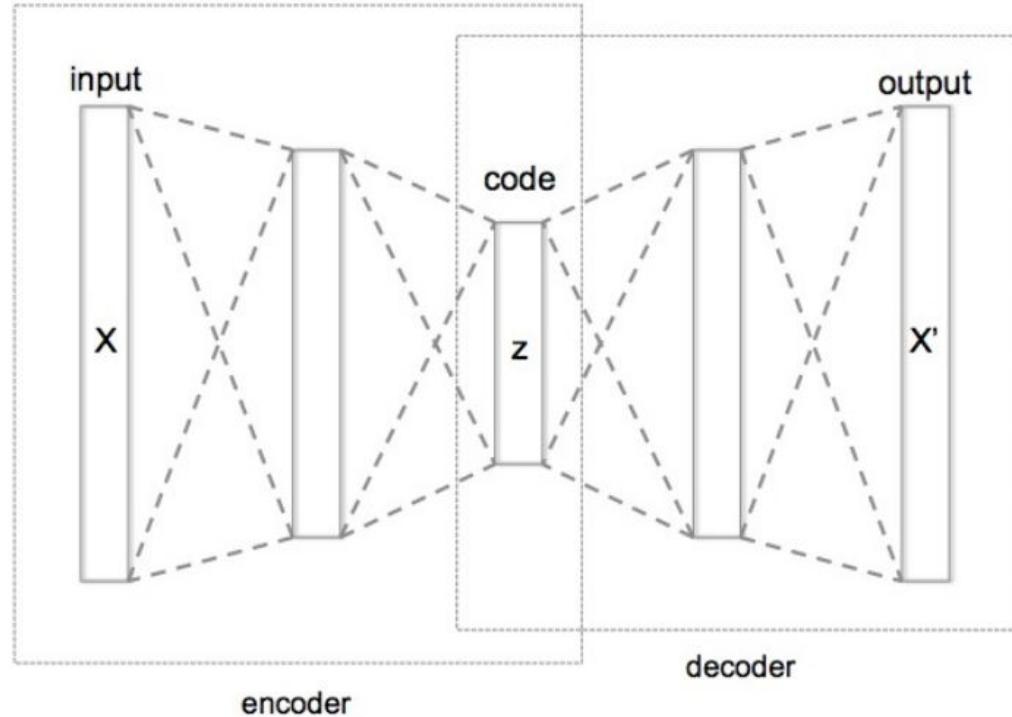
Autoencoders



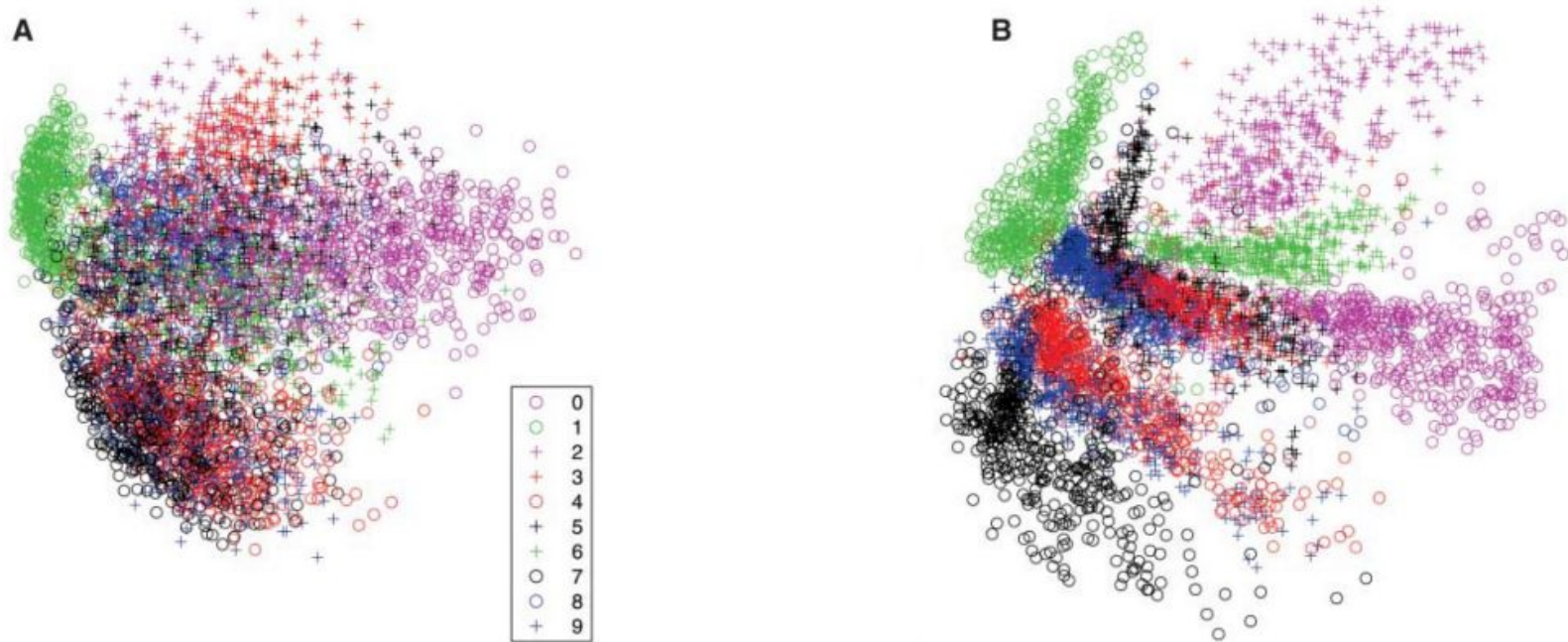
- Learn the weights by minimizing the reconstruction loss:

$$e(\mathbf{x}) = \|\mathbf{x} - \mathbf{o}^{(L)}\|_2^2$$

Deep Autoencoders



PCA (A) vs Autoencoder (B)



Takeaways

- PCA finds an orthonormal basis where the first principal component maximizes the variance \Leftrightarrow minimizes the reconstruction error
- Autoencoders use neural networks to automatically learn a latent representation that minimizes the reconstruction error

Clustering

- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$
 - Regression - $y^{(n)} \in \mathbb{R}$
 - Classification - $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$
 - Clustering
 - Dimensionality reduction
- Reinforcement learning - $\mathcal{D} = \{\mathbf{s}^{(n)}, \mathbf{a}^{(n)}, r^{(n)}\}_{n=1}^N$

Clustering

- Goal: split an unlabeled data set into groups or clusters of “similar” data points
- Use cases:
 - Organizing data
 - Discovering patterns or structure
 - Preprocessing for downstream machine learning tasks
- Applications:
 - Finding similar customer profiles for marketing
 - Reduce the number of data points via sampling from clusters
 - Use clustering to generate labels
 - Doing visualization

Similarity for k-nn

- Intuition: ~~predict the label of a data point to be the label of the “most similar” training point~~ two points are “similar” if the distance between them is small
- Euclidean distance: $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2$

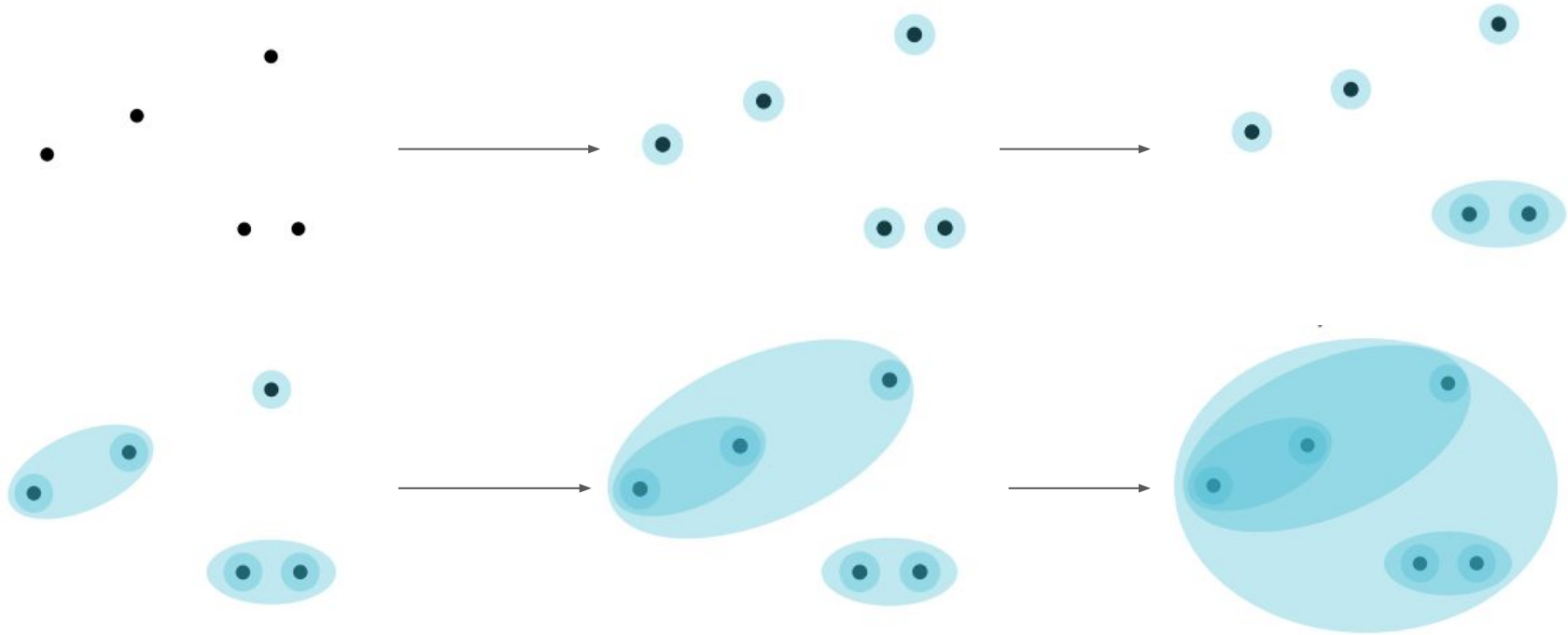
Clustering Algorithms

- Hierarchical
 - Top-down (divisive)
 - Bottom-up (agglomerative)
- Partitioning
 - K-means

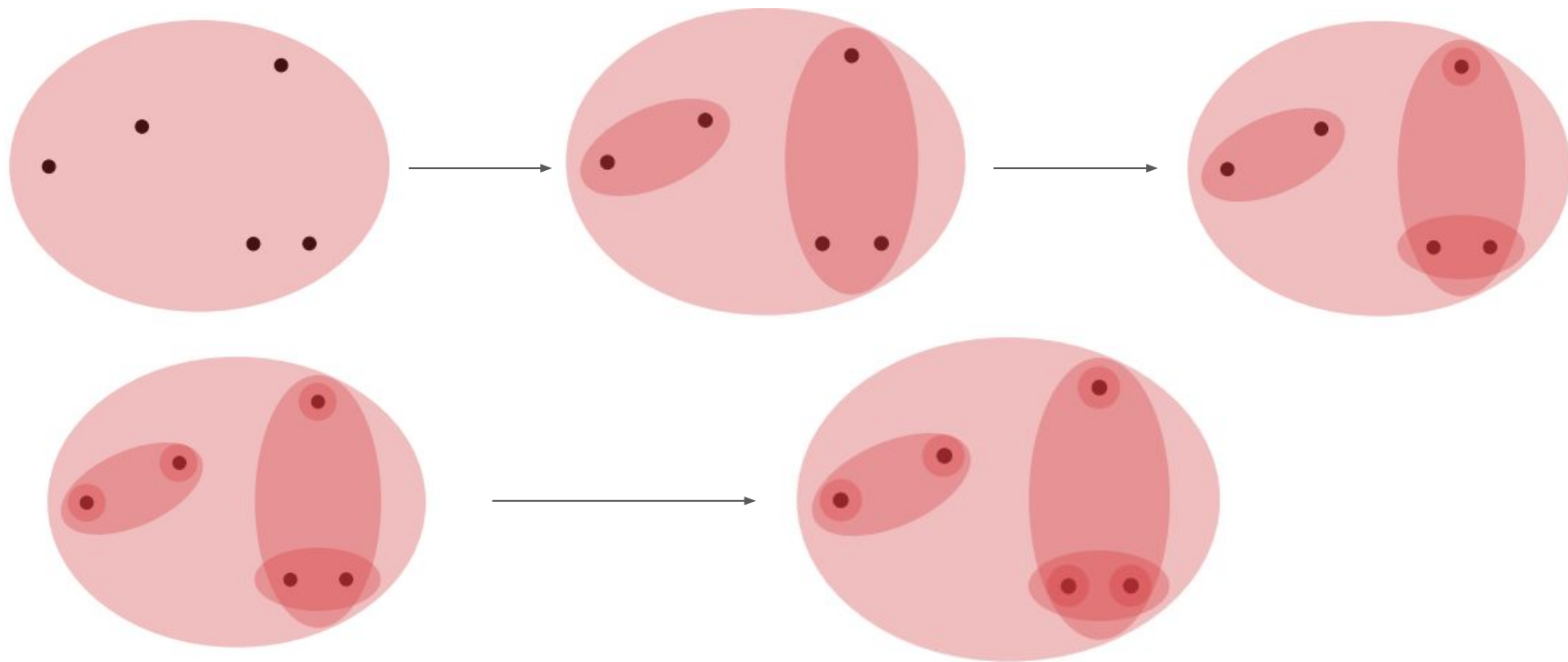
Hierarchical Clustering

- Bottom-up (agglomerative)
 - Start with each data point in its own cluster
 - Iteratively combine the most similar clusters
 - Stop when all data points are in a single cluster
- Top-down (divisive)
 - Start with all data points in one cluster
 - Iteratively split the largest cluster into two clusters
 - Stop when all clusters are single data points

Bottom-up (agglomerative)



Top-down (divisive)



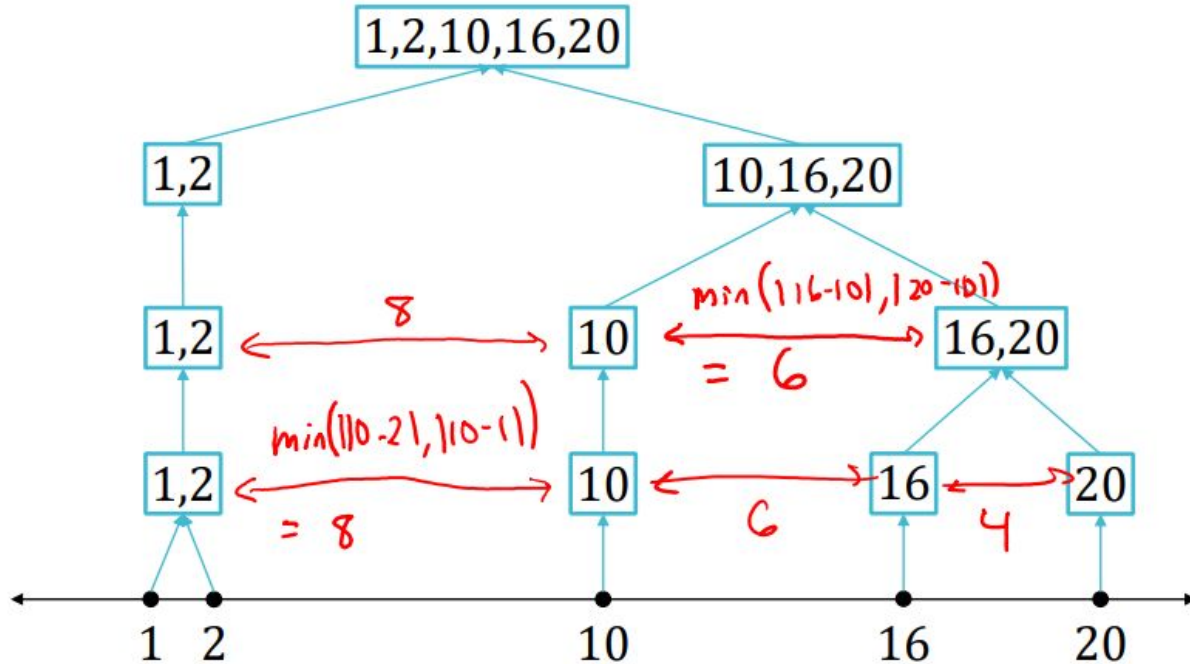
Bottom-up Hierarchical Clustering

- Bottom-up (agglomerative)
 - Start with each data point in its own cluster
 - Iteratively combine the **most similar** clusters
 - Stop when all data points are in a single cluster
- Key question: how do we define similarity between clusters?
 - Single-linkage: consider the closest data points
- Complete-linkage: consider the farthest data points

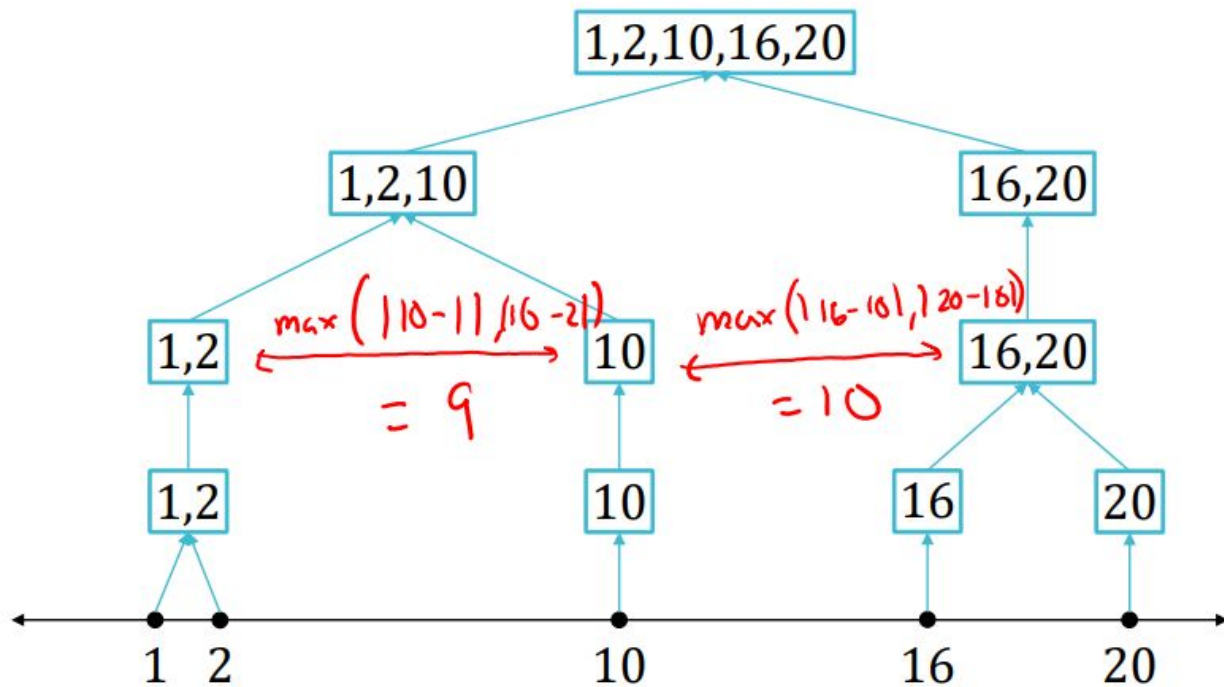
$$d_{SL}(C_i, C_j) = \min_{x \in C_i, y \in C_j} d(x, y)$$

$$d_{CL}(C_i, C_j) = \max_{x \in C_i, y \in C_j} d(x, y)$$

Single Linkage Dendrogram

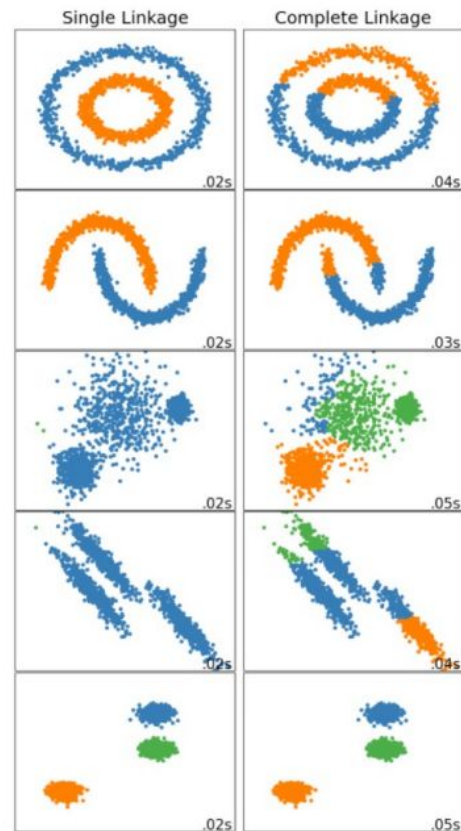


Complete Linkage Dendrogram



Single vs Complete Linkage

- Single-linkage prioritizes local behaviour, can lead to long, snakelike (i.e., nonconvex) clusters.
- Complete-linkage tries to keep the diameter of clusters small; clusters tend to be more spherical (i.e., convex)



Top-down Hierarchical Clustering

- Top-down (divisive)
 - Start with all data points in one cluster
 - Iteratively split the largest cluster into two clusters
 - Stop when all clusters are single data points
- **Key question:** how can we partition a cluster?
 - K-Means with $K = 2$

Partition-Based Clustering

- Given a desired number of clusters, K , return a partition of the data set into K groups or clusters, $\{C_1, \dots, C_K\}$, that optimize some objective function
 - Can be used as a subroutine for top-down hierarchical clustering when $K = 2$
1. What objective function should we optimize?
 2. How can we perform optimization in this setting?

Which partition is better?



Option A



Option B



Recipe for K-means

- Define a model and model parameters
 - Assume K clusters and use the Euclidean distance
 - Parameters: μ_1, \dots, μ_K and $z^{(1)}, \dots, z^{(N)}$

- Write down an objective function

$$\sum_{n=1}^N \|\mathbf{x}^{(n)} - \mu_{z^{(n)}}\|_2$$

- Optimize the objective w.r.t. the model parameters
 - Use (block) coordinate descent

Optimizing the K-means objective

$$\hat{\mu}_1, \dots, \hat{\mu}_K, z^{(1)}, \dots, z^{(N)} = \operatorname{argmin} \sum_{n=1}^N \|\mathbf{x}^{(n)} - \mu_{z^{(n)}}\|_2$$

- If μ_1, \dots, μ_K are fixed

$$\hat{z}^{(n)} = \operatorname{argmin}_{k \in \{1, \dots, K\}} \|\mathbf{x}^{(n)} - \mu_k\|_2$$

- If $z^{(1)}, \dots, z^{(N)}$ are fixed

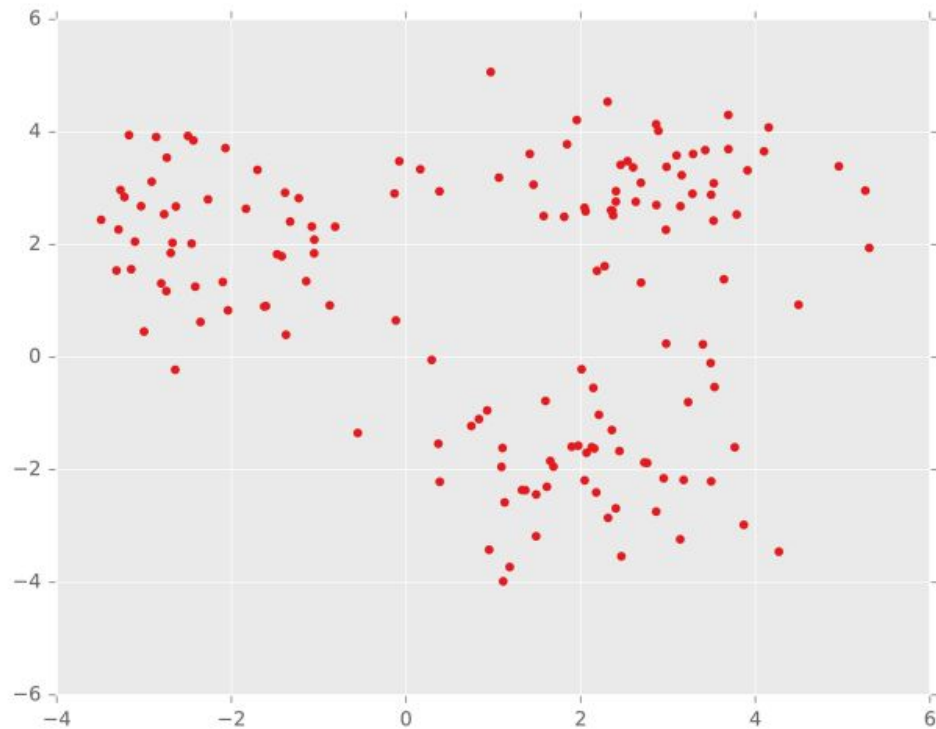
$$\hat{\mu}_k = \operatorname{argmin}_{\mu} \sum_{n: z^{(n)} = k} \|\mathbf{x}^{(n)} - \mu\|_2$$

$$= \frac{1}{N_k} \sum_{n: z^{(n)} = k} \mathbf{x}^{(n)}$$

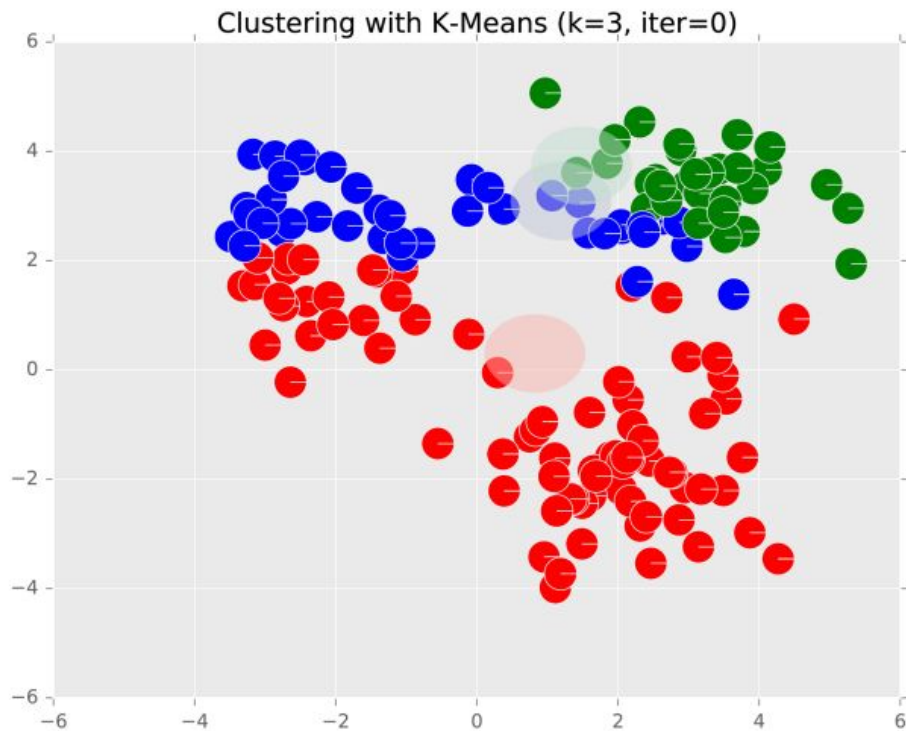
K-means Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, K$
 1. Initialize cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$
 2. While NOT CONVERGED
 - a. Assign each data point to the cluster with the nearest cluster center:
$$z^{(n)} = \underset{k}{\operatorname{argmin}} \|\mathbf{x}^{(n)} - \boldsymbol{\mu}_k\|_2$$
 - b. Recompute the cluster centers:
$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n: z^{(n)}=k} \mathbf{x}^{(n)}$$
where N_k is the number of data points in cluster k
- Output: cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ and cluster assignments $z^{(1)}, \dots, z^{(N)}$

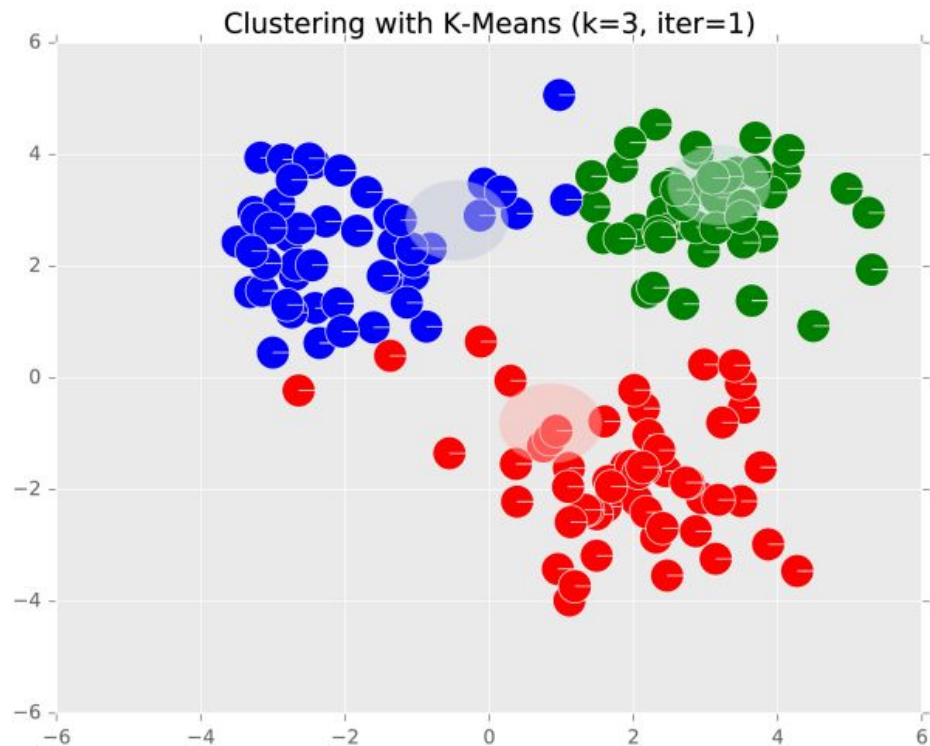
K-means: Example ($K = 3$)



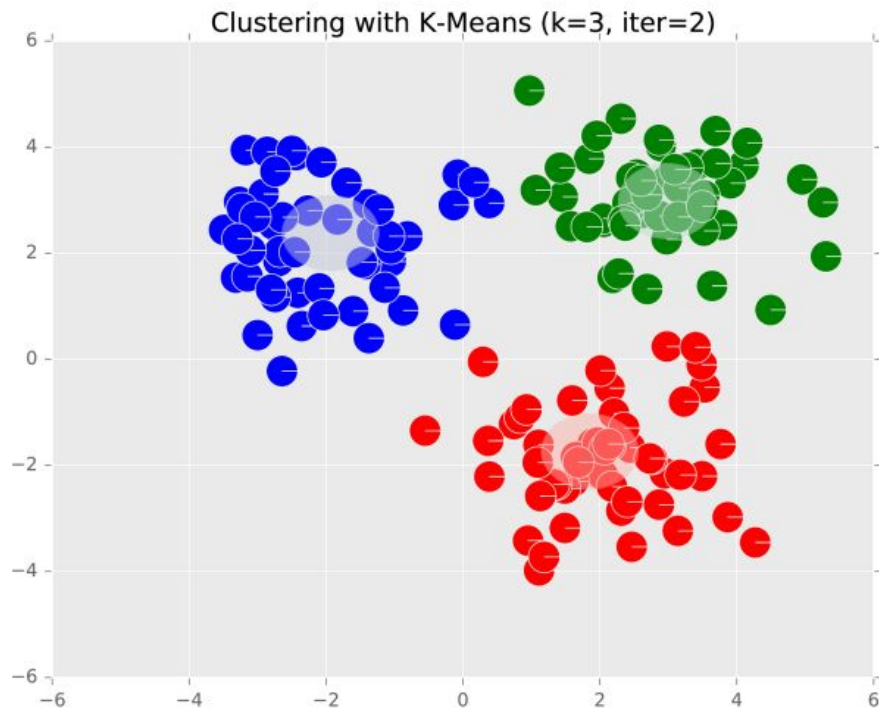
K-means: Example ($K = 3$)



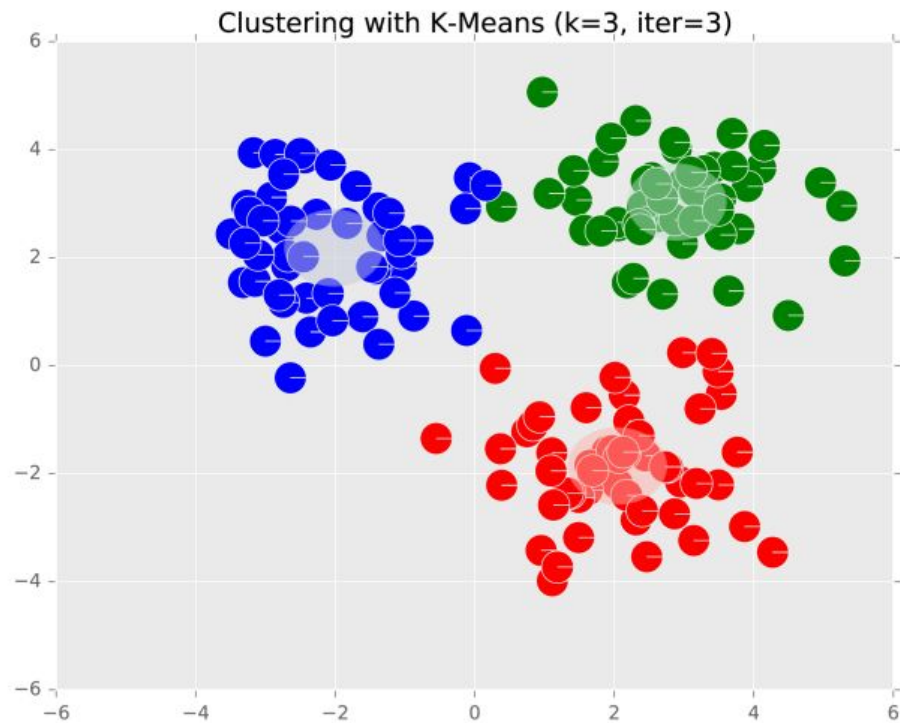
K-means: Example ($K = 3$)



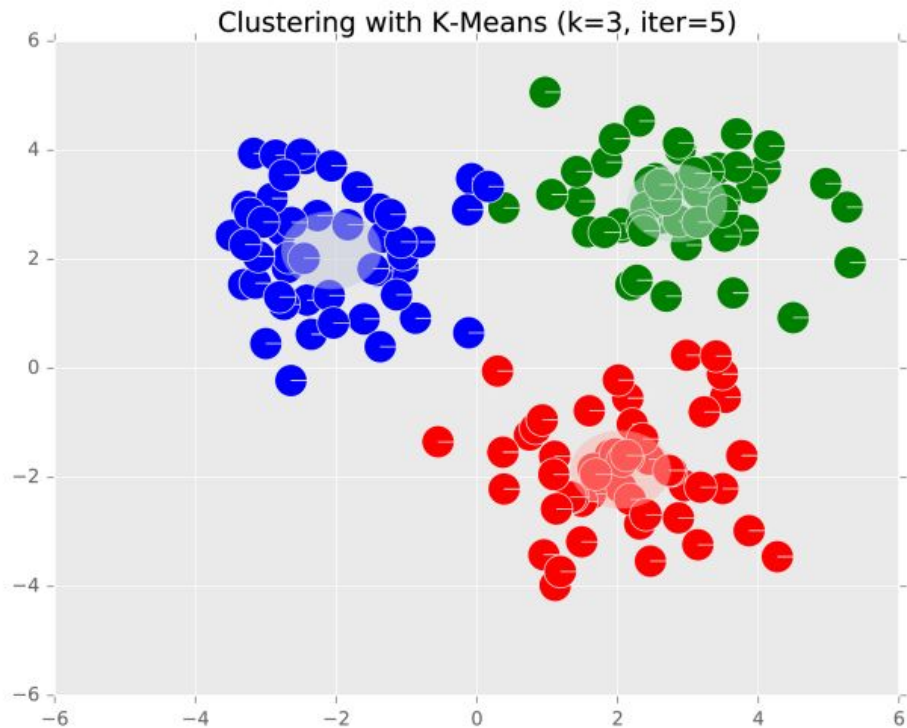
K-means: Example ($K = 3$)



K-means: Example ($K = 3$)

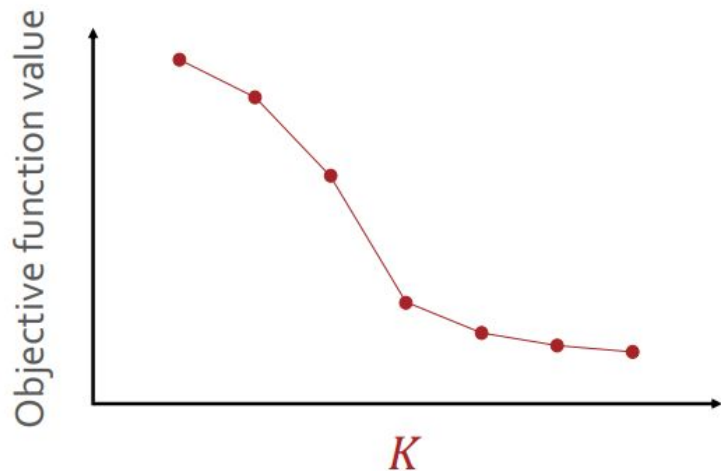


K-means: Example ($K = 3$)



Setting K

- Idea: choose the value of K that minimizes the objective function



- Look for the characteristic “elbow” or largest decrease when going from $K - 1$ to K

Next Class:

Machine Learning System Design