# STAT 432: Basics of Statistical Learning

Clustering and PCA

Shiwei Lan, Ph.D. <shiwei@illinois.edu>

http://shiwei.stat.illinois.edu/stat432.html

January 22, 2019

University of Illinois at Urbana-Champaign

# Unsupervised Learning

## Unsupervised Learning

- No response variable $Y$, only $\{x_i\}_{i=1}^n$.
- Goal: learn patterns in $X$.
- Examples
    - Estimate the density, covariance, graph (network), etc. of $X$ — could be difficult in high-dimensional settings
    - Cluster analysis: identify multiple regions of the feature space that contains modes of density.
    - Dimension reduction: identify low-dimensional manifold within the feature space $\mathcal{X}$ that represents high data density.
- Often times, there is no clear measure of success.

# Cluster Analysis

## Cluster Analysis

- Group the dataset into subsets so that those within each subset are more closely related (similar) to each other than those objects assigned to other subsets. Each subset is called a cluster

- Flat clustering vs. hierarchical clustering: flat clustering divides the dataset into $k$ cluster, and hierarchical clustering arranges the clusters into a natural hierarchy.

- Clustering results are crucially dependent on the measure of similarity (or distance) between the "points" to be clustered.

# Distance Metric

- A distance metric or a distance function is a function that defines the similarity of two elements (points, sets, etc.)

- For the distance of two points (with continuous entries), the most commonly used measurement is the Euclidian distance:

$$d(u,v) = \|u - v\|_2$$
$$= \sqrt{\sum_{j=1}^{p}(u_j - v_j)^2}$$

- For categorical entries, the Hamming distance is usually used

$$d(u,v) = \sum_{j=1}^{p} \mathbf{1}\{u_j \neq v_j\}$$

- Distance measures should be defined based on the application. There is no universally best approach.

# Clustering

- Suppose we have a set of $n$ data points
- We want to form $K \ll n$ clusters, indexed by $k \in \{1, \ldots, K\}$.
- Let $C(\cdot)$ be a cluster index function that assign th $i$th observation or cluster $C(i)$.
- Consider: search for a function $C : \{1, \ldots, n\} \to \{1, \ldots, K\}$ to minimize the within cluster distance:

$$W(C) = \frac{1}{2} \sum_{k=1}^{K} \sum_{C(i), C(i')=k} d(x_i, x_{i'}).$$

## Clustering

- This is equivalent to maximizing the between cluster distance

$$B(C) = \frac{1}{2} \sum_{k=1}^{K} \sum_{C(i)=k} \sum_{C(i')\neq k} d_{ii'}$$

- Note that the total distance can be broke down into

$$T = \frac{1}{2} \sum_{i=1}^{n} \sum_{i'=1}^{n} d_{ii'} = \frac{1}{2} \sum_{k=1}^{K} \sum_{C(i)=k} \left[ \sum_{C(i')=k} d_{ii'} + \sum_{C(i')\neq k} d_{ii'} \right]$$

$$= W(C) + B(C)$$

- The total distance is fixed for a given set of data, hence

$$\text{minimize } W(C) \iff \text{ maximize } B(C)$$

# Clustering

- Given a specific distance measure $d(\cdot, \cdot)$, several algorithms can be used to find the clusters
  - Combinatorial algorithm
  - $K$-means clustering
  - Hierarchical clustering

# Combinatorial Algorithm

## Combinatorial Algorithms

- For small $n$ and $K$, we could minimize $W$ by brute-force search.
- However, the number of "tries" needed to complete the search is

$$S(n,K) = \frac{1}{K!} \sum_{k=1}^{K} (-1)^{K-k} \binom{K}{k} k^n$$

- For example $S(10,4) = 34,105$; $S(19,4) \approx 10^{10}$.
- This is not feasible for large $n$ and $K$, since the number of distinct assignments can be extremely large.
- It calls for more efficient algorithms: may not be optimal but a reasonably good suboptimal partition.

# $K$-means Clustering

# $K$-means Clustering

- Consider an enlarged optimization problem:

$$\min_{C, \{m_k\}_{k=1}^K} \sum_{k=1}^K \sum_{C(i)=k} \|x_i - m_k\|^2$$

- Hence, we are solving both
  - the cluster index function $C(\cdot)$,
  - and also the cluster centers $m_k$, $k = 1 \ldots K$.
- This problem is NP-hard for $\geq 2$ dimensions.

# $K$-means Clustering

- Combinatorial algorithm is too expansive.
- Instead, consider an algorithm that alternatively updates the two components:
  - $C$, the cluster assignments
  - $\{m_k\}_{k=1}^{K}$: the cluster means
- We will do <span style="color:orange">an iterative update</span> by:
  1) Fixing $C$, find the best $\{m_k\}_{k=1}^{K}$
  2) Fixing $\{m_k\}_{k=1}^{K}$, find the best $C$

# $K$-means Clustering

- Fixing $C$, we know the cluster label of each subject. For any set $\{i : C(i) = k\}$, finding the mean is

$$m_k = \arg\min_m \sum_{C(i)=k} \|x_i - m\|^2.$$

- This is simply finding the mean within cluster $k$, i.e.

$$m_k = \frac{\sum_{C(i)=k} x_i}{\sum_i \mathbf{1}\{C(i) = k\}}$$

# $K$-means Clustering

- Fixing the cluster means $\{m_k\}_{k=1}^K$, to find the new cluster assignments, we simply recalculate the distance from an observation to each of the cluster mean.

$$C(i) = \underset{k \in \{1,...,K\}}{\arg\min} \ d(x_i, m_k)$$

- Hence each point will be assigned to the closest cluster mean

# $K$-means Clustering

- A $K$-means Clustering algorithm:
  1) Randomly split the dataset into $K$ different subsets. Assign each subsets a cluster label. Then iterate between 2) and 3).
  2) Given the cluster assignment $C$, calculate the cluster mean vectors $m_1, \ldots, m_K$.
  3) Given the current set of means $\{m_1, \ldots, m_K\}$, assign each observation to the closest current cluster mean.
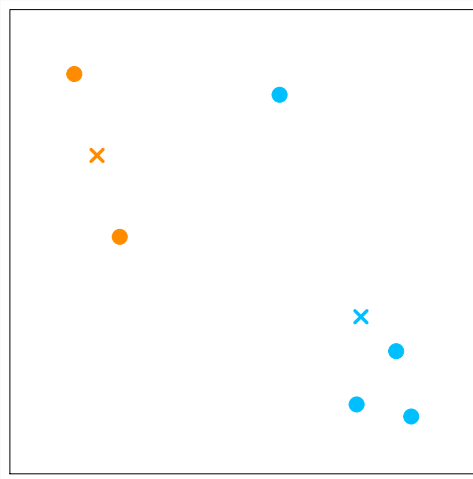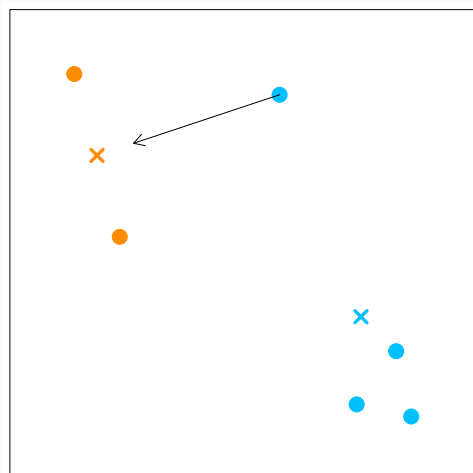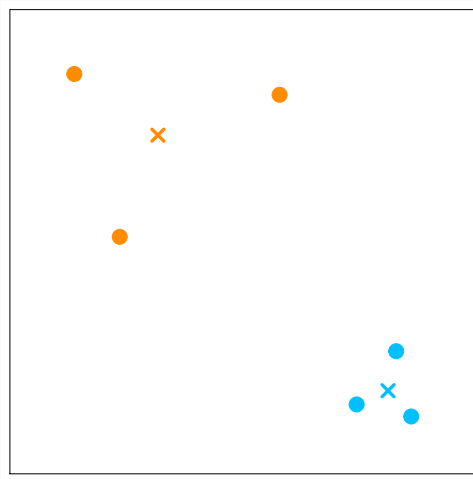- Stop the algorithm when $C$ does not change

- Note: We usually initiate the cluster labels randomly. However, this algorithm does not guarantee to global minimum. Example?
- The algorithm still has a descent property, which leads to a local minimizer.

- $K$-medoids is an alternative version of $K$-means:
- Replace the second step by searching for the one observation that minimizes the distance to all others in the cluster

$$i_k^* = \underset{i:C(i)=k}{\arg\min} \sum_{C(i')=k} D(x_i, x_{i'})$$

- Use $x_{i_k^*}$ as the "center" of cluster $k$.

- See the supplementary R file
- Example 1: the iris data
- Example 2: cluster pixels in a picture

- Why do clustering when we have the true label?
- How to choose $K$
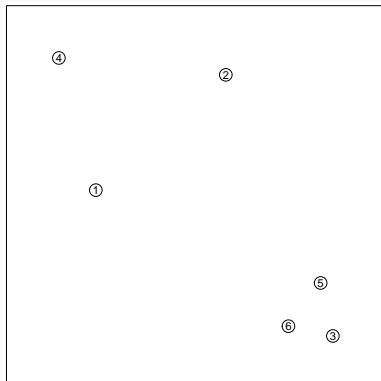- Which variables to include?

# Hierarchical Clustering

# Hierarchical Clustering

- Choosing the number of clusters $K$ can be difficult
- A hierarchical representation which
    - at the lowest level, each cluster contains a single observation.
    - at the highest level there is only one cluster containing all observations.
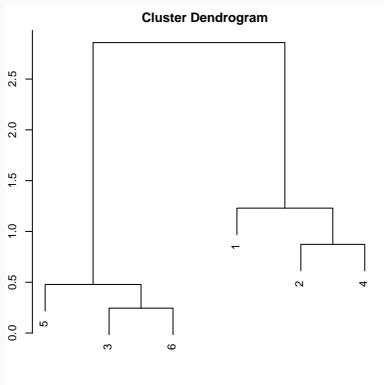- Use dendrogram to display the clustering result.
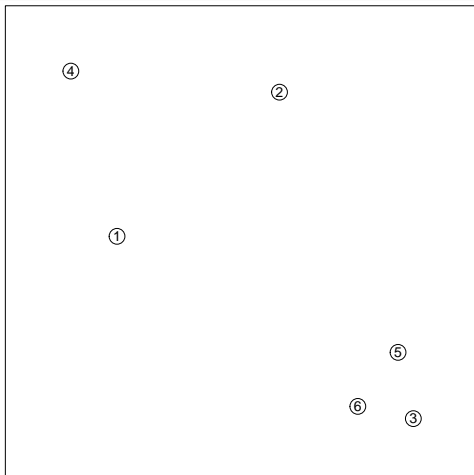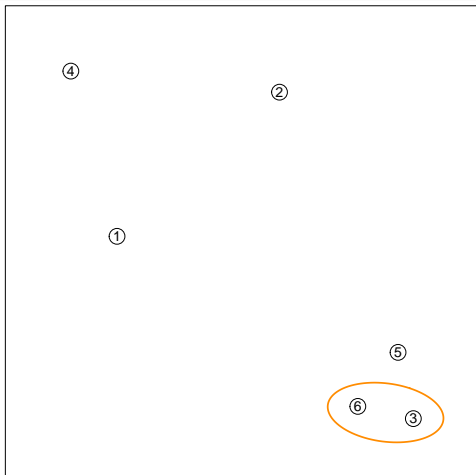
- Suppose we have a set of 6 observations

## Dendrogram

- A typical dendrogram from hierarchical clustering
- How to construct this?



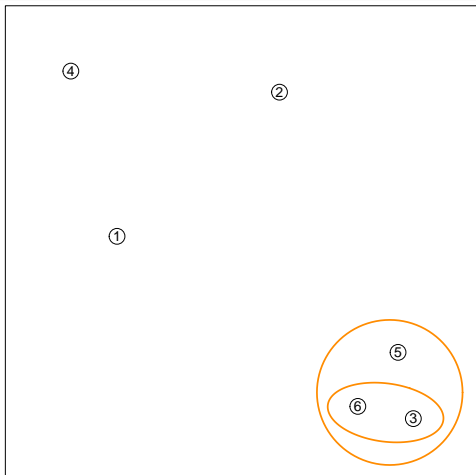**Cluster Dendrogram**

## Algorithm (agglomerative)

- An agglomerative algorithm is a "bottom up" approach:
  - Begin with every observation representing a singleton cluster.
  - At each step, merge two "closest" clusters into one cluster and reduce the number of clusters by one.
  - Stop when all observations are in the same cluster
- How to choose which two clusters to merge?
- This requires:
  - A distance measure between any two observations $d(x_i, x_{j'})$
  - A distance measure between any two sets of observations $d(G, H)$

## Distance Measures

- Distance $d(G, H)$ between two clusters $G$ and $H$ can be defined in different ways:
  - Complete linkage (default of $\text{hclust()}$): the furthest pair

  $$d(G, H) = \max_{i \in G, i' \in H} d_{ii'}$$

  - Single linkage: the closest pair

  $$d(G, H) = \min_{i \in G, i' \in H} d_{ii'}$$

  - Average linkage: average dissimilarity

  $$d(G, H) = \frac{1}{n_G n_H} \sum_{i \in G} \sum_{i' \in H} d_{ii'}$$

- Different choice may results in different hierarchical structures

- To perform a hierarchical clustering, a matrix of all the pair-wise distances is sufficient
- We don't have to know the values of the original observations
- This is an $n \times n$ matrix: the $(i, i')$'s element represents the distance between $x_i$ and $x_{i'}$
- This matrix is also called a dissimilarity matrix.
  - Symmetric
  - diagonal elements are zero

## Applications

- See the supplementary $R$ file
- Example 1: the iris data
- Example 2: RNA expression data

# Principle Component Analysis

## Principle Component Analysis

- Principle Component Analysis (PCA) is an old but very useful technique invented by Karl Pearson in 1901
- The main purpose is data visualization (mostly in 2d)
- It also serves as a dimension reduction method
- Unsupervised method, can be used for preprocessing the data.

# Principle Component Analysis

- Given that we have a $n \times p$ design matrix $\mathbf{X}$, there are many equivalent approaches (motivations):
    - Explain the most variation: Produce a derived set of uncorrelated variables $\mathbf{Z}_k = \mathbf{X}\alpha_k$, $k = 1, \ldots, q < p$ that are linear combinations of the original variables, and explain most of the variation in the original set
    - Approximate the design matrix: Approximate the design matrix $\mathbf{X}$ by the best (using Frobenius norm) rank-$q$ matrix, which can be performed through SVD

## Principle Component Analysis

- Suppose we have an $n \times p$ design matrix $\mathbf{X}$
- The first step is to center each variable, i.e., subtract the column means from each column respectively.
- In the following, we assume that $\mathbf{X}$ is already centered.
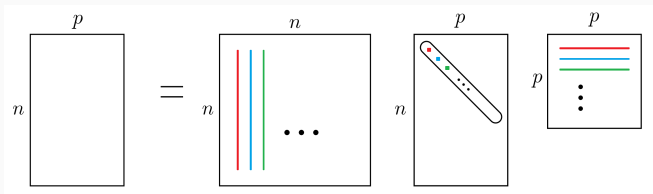- Centering $\mathbf{X}$ does nothing but re-positioning the axis

# Singular Value Decomposition

- One way to understand the PCA is using a singular value decomposition (SVD)

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times n} \mathbf{D}_{n \times p} \mathbf{V}_{p \times p}^{\mathsf{T}},$$

where both $\mathbf{U}$ and $\mathbf{V}$ are orthogonal matrices, i.e. $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{U}\mathbf{U}^{\mathsf{T}} = \mathbf{I}$, and same for $\mathbf{V}$; and $\mathbf{D}_{n \times p}$ is a diagonal matrix.



- The diagonal elements of $\mathbf{D}_{n \times p}$ are of a decreasing order.

## Low Rank Approximation

- If we have to choose a rank-1 matrix $\mathbf{A}$ to approximate $\mathbf{X}_{n \times p}$, what would we do?

- Turns out that the best choice is

$$\mathbf{U}_1 d_{11} \mathbf{V}_1^\mathsf{T},$$

where $\mathbf{U}_1$ is the first column of $\mathbf{U}$, $\mathbf{V}_1$ is the first column of $\mathbf{V}$, and $d_{11}$ is the first diagonal element of $\mathbf{D}$

- In other words, $\|\mathbf{X} - \mathbf{U}_1 d_{11} \mathbf{V}_1^\mathsf{T}\|_2^2$ is minimized with this choice.

- Hence, $\mathbf{U}_1 d_{11} \mathbf{V}_1^\mathsf{T}$ is a rank-1 matrix that explained the variations in $\mathbf{X}$ as much as possible.

## Principle Component Analysis

- Let's consider the sample covariance matrix $\widehat{\Sigma} = \mathbf{X}^\mathsf{T}\mathbf{X}/(n-1)$, since $\mathbf{X}$ is already centered.
- $\widehat{\Sigma}$ can be diagonalize into

$$\widehat{\Sigma} = \mathbf{V}\mathbf{D}^*\mathbf{V}^\mathsf{T},$$

where columns of $\mathbf{V}$ are principle directions (loadings) and projecting $\mathbf{X}$ on these loadings gives the principal components

- On the other hand, based on SVD,

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\mathsf{T},$$

we can rewrite $\widehat{\Sigma}$ as

$$\widehat{\Sigma} = \mathbf{V}\mathbf{D}^\mathsf{T}\mathbf{U}^\mathsf{T}\mathbf{U}\mathbf{D}\mathbf{V}^\mathsf{T}/(n-1) = \mathbf{V}\frac{\mathbf{D}^2}{n-1}\mathbf{V}^\mathsf{T}$$
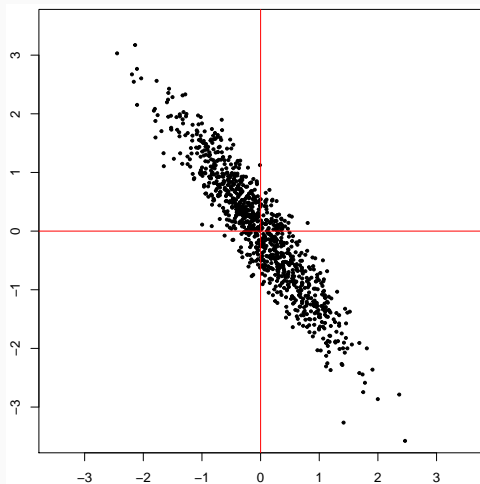
## Principle Component Analysis

- So the right singular vectors $\mathbf{V}$ of $\mathbf{X}$ are just the principle directions, and the principal components are basically projecting each row (observation) of $\mathbf{X}$ onto those directions:

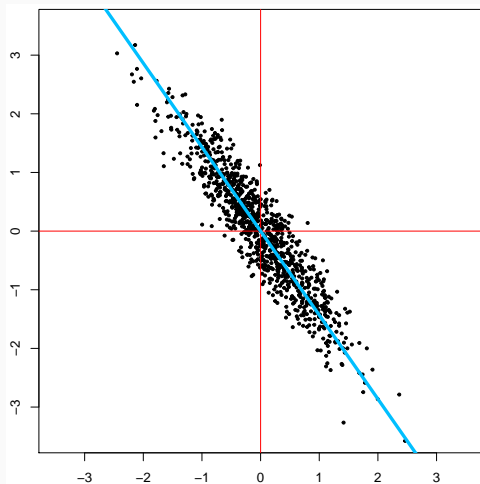$$\mathbf{XV} = \mathbf{UDV}^\mathsf{T}\mathbf{V} = \mathbf{UD}$$

- The first column of $\mathbf{U}$ is the first PC, and $d_{11}$ is the variation along that direction, which is also the squared eigenvalue from SVD.

- PCA should be performed by centering $\mathbf{X}$ first (column-wise, i.e., by each variable).

# Demonstration

# Demonstration

# Principle Component Analysis

- PCA is a dimension reduction tool, often used for visualization
- The leading PCs may display interesting information of the underlying (unobserved) clusters/manifold
- PCA is unsupervised, i.e., the directions does not necessarily reflect the relationship with the response (if there is any)