

SVD

Recall from PCA: Need to compute $X^T X$, and its (i, j) pairs

Problem: - $\Theta(nd^2)$

- $X^T X$ might be poorly conditioned

↳ numerically inaccurate eigenvectors

SVD improves both
vs

SVD Works for all Matrices!

SVD: can always find $X = UDV^T$

$$X = UDV^T$$

$$= \left[\vec{u}_1 \dots \vec{u}_d \right] \begin{bmatrix} r_1 & & 0 \\ & \ddots & \\ 0 & & r_d \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

left Sing. Values

Outerproduct form: $X = \sum_{i=1}^d g_i u_i v_i^T$

If $n \geq d$, then \mathbf{U} is $n \times d$, \mathbf{D} is $d \times d$, and \mathbf{V}^T is $d \times d$. If $n < d$, then \mathbf{U} is $n \times n$, \mathbf{D} is $n \times n$, and \mathbf{V}^T is $n \times d$.

Number of $\hat{\sigma}_i \neq 0 = \text{rk}(X)$

↳ if X is centered Matrix for pts. that all lie on a line, then only $\sigma_1 \neq 0$

Theorem: Consider $X = UDV^T$, then

- \vec{v}_i is an eigenvector of $X^T X$ w/ eigenvalue σ_i^2

$$\text{proof: } x^T x = V D^T D v^T = V D^2 V^T \rightarrow f_i = \sigma_i^2$$

PCA needs eigenvectors of $X^T X$ (rows of V^T)

- SVD also gives the eigenvalue λ_i of \widehat{V} :

↳ SVD more numerically stable since ratio between σ_i & σ_j smaller than λ_i & λ_j

* if $n > d$, then V omits corresponding (σ_i, \vec{v}_i) pairs

- in $O(ndk)$ time, we can find the k greatest singular values & their corresponding eigenvectors

↳ save time by computing only some of the vectors

Important: row i of UD gives principal coordinates of X_i , i.e. $x_i \cdot \vec{v}_j = \sigma_j u_{ij} \forall j$
↳ SVD gives inner product $x_i \cdot \vec{v}_j$, don't need to compute it ourselves

Clustering: goal: partition data into clusters by similarity

examples of where we'd apply clustering:

- Discovery: find songs similar to songs you like, determine market segments
- Hierarchy: find good taxonomy of species from genes
- Quantization: compress dataset by reducing choices
- Graph Partitioning: image segmentation, finding groups in social networks

k -means clustering: goal - to partition n points into k disjoint clusters

- assign each input x_i a cluster label $y_i \in \{1, k\}$,

- cluster i 's mean: $\bar{M}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_j$, n_i pts. in cluster i

formal optimization problem: $\bar{Y}^* = \underset{Y}{\operatorname{argmin}} \sum_{i=1}^k \sum_{j:y_j=i} \|x_j - M_i\|_2^2$

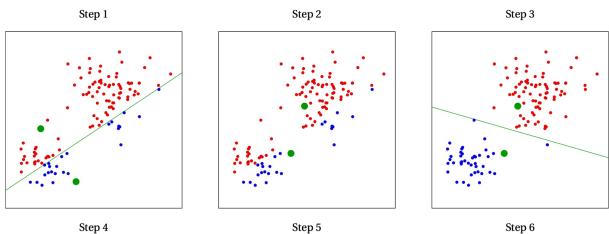
- NP-hard

- $O(nk^N)$ time by trying every partition

- We can use heuristics to get an approximation
- Lloyd's algo:
 - alternate between:
 - 1) fix \vec{Y}_j 's, update M_i 's
 - 2) fix M_i , update \vec{Y}_j
 - halt when (2) makes no changes

- One can show through calculus that optimal M_i in (1) is mean of all pts. in cluster;
- one can also show that the optimal \vec{Y} in (2) assigns each pt. X_j to the closest mean M_i
- if there's a tie but X_j has option to stay in old cluster, choose that

odd Steps: relabel datapoints →
 even steps: compute new means



- both steps decrease value of obj fn.
- ↳ no previous state is ever revisited
- ↳ Lloyd's algo. never runs forever
- ↳ doesn't say anything optimistic about runtime
- ↳ might see $O(K^M)$ assignments before halting
- ↳ Lloyd's algo. usually finds local minima, not global minima
- expected since problem is NP-hard

Figure 21.1: An example of 2-means clustering

Ways to initialize the means:

better than random partitioning ↘

- 1) **Forgy Method:** Choose k random sample pts. to be the means M_i
- 2) **Random Partition:** randomly assign each sample pt. to a cluster, then update means
- 3) **k -means++:** Forgy but w/ biased distribution

- each center is chosen w/ preference towards points far from previous centers

best results: run k-means multiple times w/ random starts

Equivalent Objective fn.: minimize Within-cluster Variation:

$$\sum_{i=1}^k \frac{1}{n_i} \sum_{j: v_j=i} \sum_{m=j} \|x_j - x_m\|^2$$

- * at minimizes, this fn. = Z-value of prev. objective function
- this expression doesn't include the means
 - only depends on inputs & clusters we assign them to

normalization before k-means?

- Same w/ PCA: Yes and no

Yes: Some features much larger than others \rightarrow they dominate Euclidean distance

No: If features are in same unit of measurement

k-medoids clustering: generalizes k-means beyond Euclidean distance

\hookrightarrow means not optimal for other distance metrics

- Specify distance fn. $d(\vec{x}, \vec{y})$ between \vec{x} & \vec{y} , known as dissimilarity
 - $d(\vec{x}, \vec{y})$ is arbitrary
 - ideally satisfies triangle inequality $d(\vec{x}, \vec{y}) \leq d(\vec{x}, \vec{z}) + d(\vec{z}, \vec{y})$
 - Some use l_1 or l_∞ or matrix of pairwise distances between inputs

e.g. When Euclidean dist. Undesired: market analysis

- Want to cluster customers who buy similar products
- Makes more sense to use angle of vectors as dissimilarity rather than euclidean dist.

Medoid: Sample pt. that minimizes total dissimilarity

\hookrightarrow Medoid always one of the sample points

problem w/ k-means or k-medoids: need to choose k before starting
↳ no reliable way of guessing best k

Hierarchical Clustering: create a tree in which every subtree is a cluster
↳ some clusters contain smaller clusters

agglomerative clustering: each point is its own cluster, repeatedly fusing them
(bottom-up approach) use when input is a point set $\mathcal{O}(n^3)$

divisive clustering: start w/ all points in 1 cluster, repeatedly split them
(top-down) use when input is a graph

- hierarchical clustering needs a distance function/linkage for clusters A & B

Complete linkage: $d(A, B) = \max_{\vec{w} \in A, \vec{x} \in B} d(\vec{w}, \vec{x})$ } work for any distance fn.

Single linkage: $d(A, B) = \min_{\vec{w} \in A, \vec{x} \in B} d(\vec{w}, \vec{x})$

Average Linkage: $d(A, B) = \frac{1}{|A| \cdot |B|} \sum_{\vec{w} \in A} \sum_{\vec{x} \in B} d(\vec{w}, \vec{x})$

Centroid linkage $d(A, B) = d(\vec{\mu}_A, \vec{\mu}_B)$ (only if using Euclidean distance)

↳ can use medoids instead of means

* medoids more robust to outliers than means

All Clustering algorithms are unstable

continue below

Dendograms: Illustration of cluster hierarchy (tree)

- Vertical axis: encodes linkage distances
- horiz. axis has no meaning

draw horizontal line to split tree into clusters

complete linkage is the best balanced: When cluster is large, outside points are always far away \rightarrow large clusters more resistant to growth than small ones

* in most applications: use avg. or complete linkage

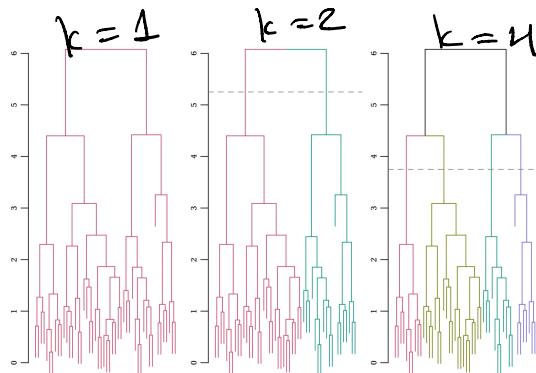


Figure 21.2: Partitioning a dendrogram into clusters

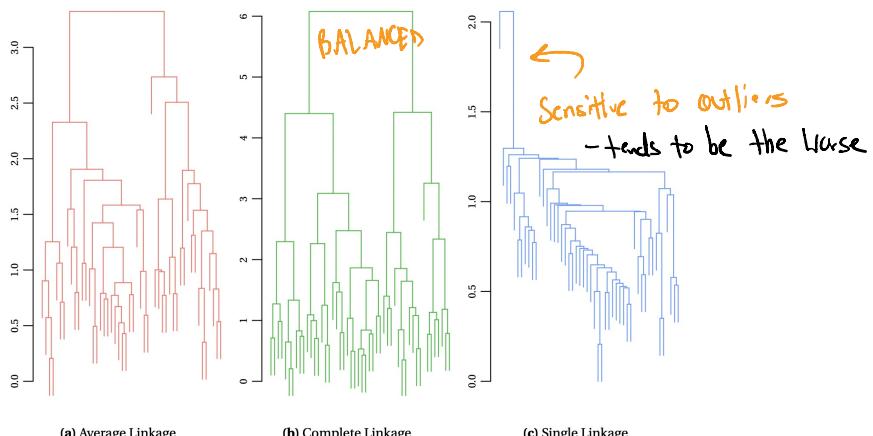


Figure 21.3: Comparison between dendograms of various linkage functions

