

## 21 The Singular Value Decomposition; Clustering

### THE SINGULAR VALUE DECOMPOSITION (SVD) [and its Application to PCA]

Problems with PCA: Computing  $X^T X$  takes  $\Theta(nd^2)$  time.  
 $X^T X$  is poorly conditioned  $\rightarrow$  numerically inaccurate eigenvectors.  
 [The SVD improves both these problems.]

[Earlier this semester, we talked about the eigendecomposition of a square, symmetric matrix. Unfortunately, nonsymmetric matrices don't have nice eigendecompositions, and non-square matrices don't have eigenvectors at all. Happily, there is a similar decomposition that works for all matrices, even if they're not symmetric and not square.]

Fact: Every  $X \in \mathbb{R}^{n \times d}$  has a singular value decomposition  $X = UDV^T$  of the form

$$\begin{array}{c}
 X \\
 \begin{array}{|c|} \hline \\ \hline \end{array} \\
 n \times d
 \end{array}
 =
 \begin{array}{c}
 U \\
 \begin{array}{|c|} \hline \\ \hline \end{array} \\
 n \times n
 \end{array}
 \begin{array}{c}
 D \\
 \text{diagonal} \\
 \begin{array}{|c|} \hline \\ \hline \end{array} \\
 n \times d
 \end{array}
 \begin{array}{c}
 V^T \\
 \begin{array}{|c|} \hline \\ \hline \end{array} \\
 d \times d
 \end{array}
 =
 \sum_{i=1}^{\min\{n,d\}} \underbrace{\delta_i u_i v_i^T}_{\text{rank 1 outer product matrix}}$$

$U^T U = I = U U^T$   
 orthonormal  $u_i$ 's are left singular vectors of  $X$

$V^T V = I = V V^T$   
 orthonormal  $v_i$ 's are right singular vectors of  $X$

[Draw this by hand; write summation at the right last. [fullsvd.pdf](#)]

Diagonal entries  $\delta_1, \dots, \delta_{\min\{n,d\}}$  of  $D$  are nonnegative singular values of  $X$ .

[By convention, singular values are never negative, but some of them might be zero.]

Fact:  $v_i$  is an eigenvector of  $X^T X$  w/eigenvalue  $\delta_i^2$ , so the SVD solves PCA.

Proof:  $X^T X = V D^T U^T U D V^T = V D^T D V^T$   
 which is an eigendecomposition of  $X^T X$ , with each  $\delta_i^2$  on the diagonal of  $D^T D$ .

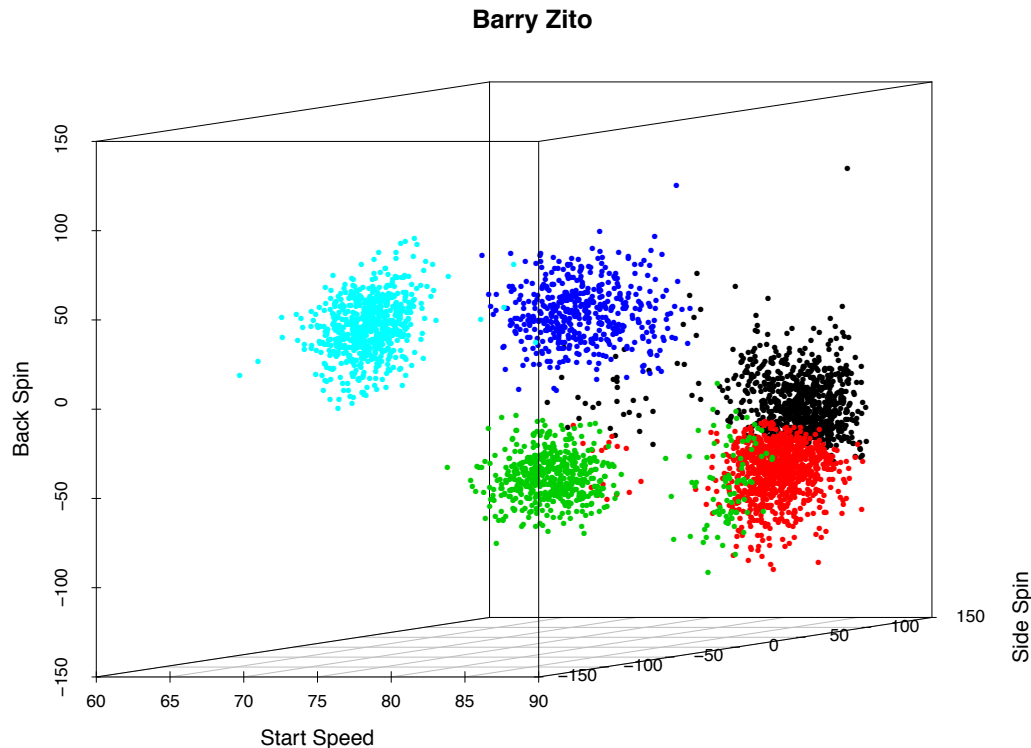
[The columns of  $V$  are the eigenvectors of  $X^T X$ , which are the principal components we need for PCA. The SVD also tells us their eigenvalues, which are the squares of the singular values of  $X$ . That's related to why the SVD is more numerically stable: the ratios between singular values are smaller than the ratios between eigenvalues, so it's easier to stably compute the singular values of  $X$  than the eigenvalues of  $X^T X$ .]

Important: Row  $i$  of  $UD$  gives the principal coordinates of sample point  $X_i$  (i.e.,  $\forall i, \forall j, X_i \cdot v_j = \delta_j U_{ij}$ ).

[So we don't need to explicitly compute the inner products  $X_i \cdot v_j$ ; the SVD has already done it for us.]

Proof:  $XV = UDV^T V = UD$ , so  $(XV)_{ij} = (UD)_{ij}$ .





4-Seam Fastball	2-Seam Fastball	Changeup	Slider	Curveball
<b>Black</b>	<b>Red</b>	<b>Green</b>	<b>Blue</b>	<b>Light Blue</b>

zito.pdf (from a talk by Michael Pane) [k-means clusters that classify Barry Zito's baseball pitches. Here we discover that there really are distinct classes of baseball pitches.]

### k-Means Clustering aka Lloyd's Algorithm (Stuart Lloyd, 1957)

Goal: Partition  $n$  points into  $k$  disjoint clusters.

Assign each sample point  $X_i$  a cluster label  $y_i \in [1, k]$ .

Cluster  $i$ 's mean is  $\mu_i = \frac{1}{n_i} \sum_{y_j=i} X_j$ , given  $n_i$  points in cluster  $i$ .

Find  $y$  that minimizes  $\sum_{i=1}^k \sum_{y_j=i} \|X_j - \mu_i\|^2$  [Sum of the squared distances from points to their cluster means.]

NP-hard. Solvable in  $O(nk^n)$  time. [Try every partition.]

*k*-means heuristic: Alternate between

(1)  $y_j$ 's are fixed; update  $\mu_i$ 's

(2)  $\mu_i$ 's are fixed; update  $y_j$ 's

Halt when step (2) changes no assignments.

[So, we have an assignment of points to clusters. We compute the cluster means. Then we reconsider the assignment. A point might change clusters if some other's cluster's mean is closer than its own cluster's mean. Then repeat.]

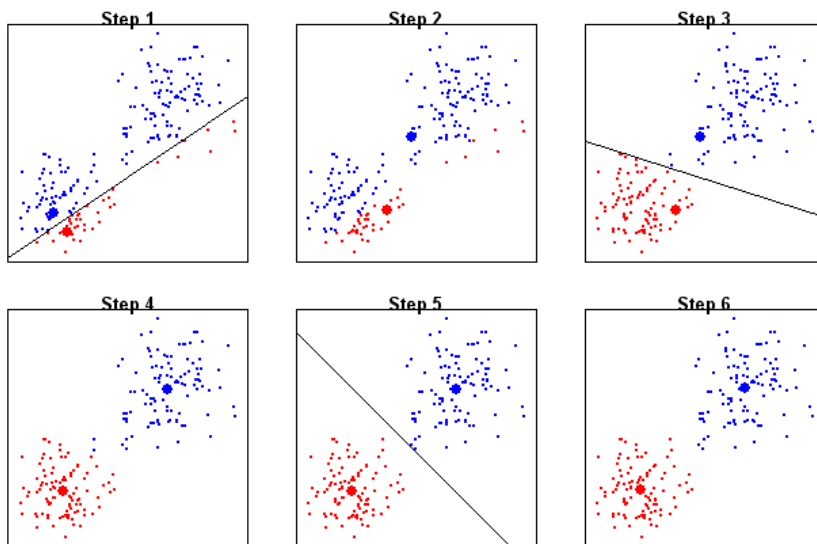
Step (1): One can show (calculus) the optimal  $\mu_i$  is the mean of the points in cluster  $i$ .

[This is easy calculus, so I leave it as a short exercise.]

Step (2): The optimal  $y$  assigns each point  $X_j$  to the closest mean  $\mu_i$ .

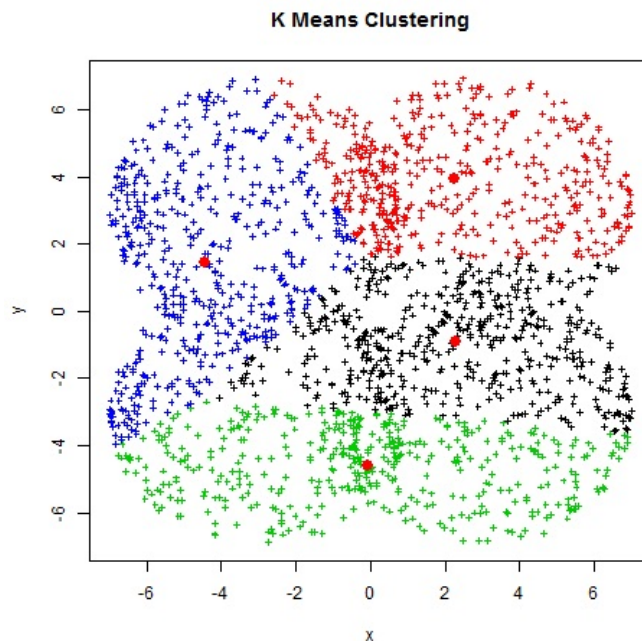
[If there's a tie, and one of the choices is for  $X_j$  to stay in the same cluster as the previous iteration, always take that choice.]

[...so both steps minimize the cost function, but they don't optimize all the variables at once.]



2means.png

[An example of 2-means. Odd-numbered steps reassign the data points. Even-numbered steps compute new means.]



4meansanimation.gif

[This is an animated GIF of 4-means with many points. Unfortunately, the animation doesn't work in the PDF lecture notes.]

Both steps decrease objective fn *unless* they change nothing.

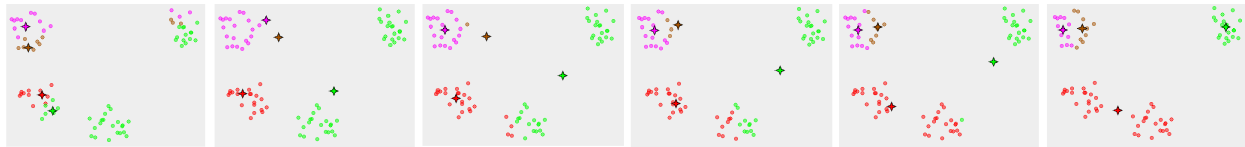
[Therefore, the algorithm never returns to a previous assignment.]

Hence alg. must terminate. [As there are only finitely many assignments.]

[This argument says that Lloyd's algorithm never loops forever. But it doesn't say anything optimistic about the running time, because we might see  $O(k^n)$  different assignments before we halt. In theory, one can actually construct point sets in the plane that take an exponential number of iterations, but those don't come up in practice.]

Usually very fast in practice. Finds a local minimum, often not global.

[... which is not surprising, as this problem is NP-hard.]



4meansbad.png [An example where 4-means clustering fails.]

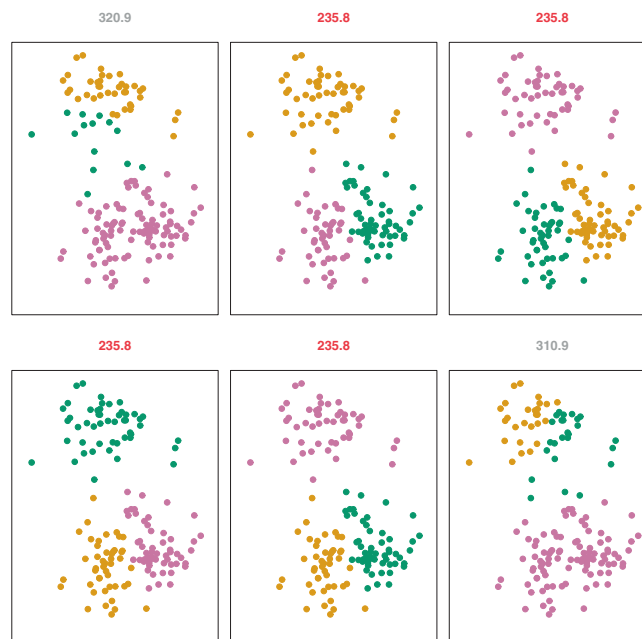
Getting started:

- Forgy method: choose  $k$  random sample points to be initial  $\mu_i$ 's; go to (2).
- Random partition: randomly assign each sample point to a cluster; go to (1).
- $k$ -means++: like Forgy, but biased distribution.

[Each  $\mu_i$  is chosen with a preference for points far from previous  $\mu_i$ 's.]

[ $k$ -means++ is a little more work, but it works well in practice and theory. Forgy seems to be better than random partition, but Wikipedia mentions some variants of  $k$ -means for which random partition is better.]

For best results, run  $k$ -means multiple times with random starts.



kmeans6times.pdf (ISL, Figure 10.7) [Clusters found by running 3-means 6 times on the same sample points, each time starting with a different random partition. The algorithm finds three different local minima.]

[Why did we choose that particular objective function to minimize? Partly because it is equivalent to minimizing the following function.]

Equivalent objective fn: the within-cluster variation

$$\text{Find } y \text{ that minimizes } \sum_{i=1}^k \frac{1}{n_i} \sum_{y_j=i} \sum_{y_m=i} \|X_j - X_m\|^2$$

[At the minimizer, this objective function is equal to twice the previous one. It's a worthwhile exercise to show that—it's harder than it looks. The nice thing about this expression is that it doesn't include the means; it's a function purely of the sample points and the clusters we assign them to. So it's more compelling.]

Normalize the data? [before applying  $k$ -means]

Same advice as for PCA. Sometimes yes, sometimes no.

[If some features are much larger than others, they will tend to dominate the Euclidean distance. So if you have features in different units of measurement, you probably should normalize them. If you have features in the same unit of measurement, you usually shouldn't, but it depends on context.]

[One difficulty with  $k$ -means is that you have to choose the number  $k$  of clusters before you start, and there isn't any reliable way to guess how many clusters will best fit the data. The next method, hierarchical clustering, has the advantage in that respect. By the way, there is a whole Wikipedia article on "Determining the number of clusters in a data set."]

## Hierarchical Clustering

Creates a tree; every subtree is a cluster.

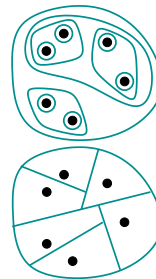
[So some clusters contain smaller clusters.]

Bottom-up, aka agglomerative clustering:

start with each point a cluster; repeatedly fuse pairs.

Top-down, aka divisive clustering:

start with all pts in one cluster; repeatedly split it.



[When the input is a point set, agglomerative clustering is used much more in practice than divisive clustering. But when the input is a graph, it's the other way around: divisive clustering is more common.]

We need a distance fn for clusters  $A, B$ :

complete linkage:  $d(A, B) = \max\{d(w, x) : w \in A, x \in B\}$

single linkage:  $d(A, B) = \min\{d(w, x) : w \in A, x \in B\}$

average linkage:  $d(A, B) = \frac{1}{|A||B|} \sum_{w \in A} \sum_{x \in B} d(w, x)$

centroid linkage:  $d(A, B) = d(\mu_A, \mu_B)$  where  $\mu_S$  is mean of  $S$

[The first three of these linkages work for any distance function, even if the input is just a matrix of distances between all pairs of sample points. The centroid linkage only really makes sense if we're using the Euclidean distance.]

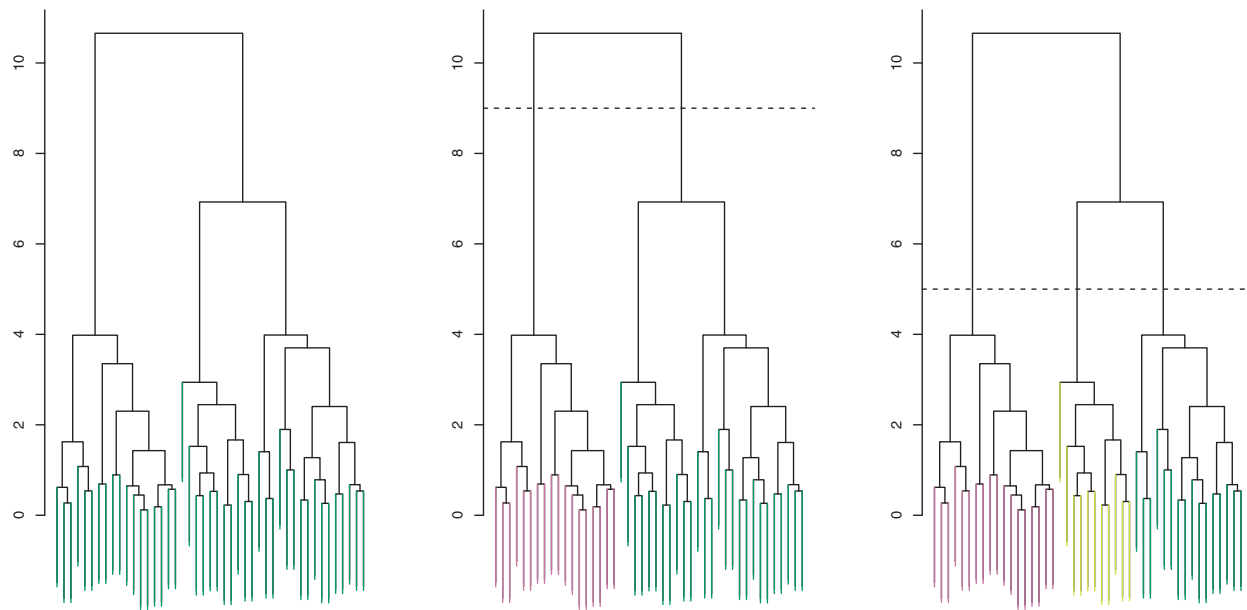
Greedy agglomerative alg.:

Repeatedly fuse the two clusters that minimize  $d(A, B)$

Naively takes  $O(n^3)$  time.

[But for complete and single linkage, there are more sophisticated algorithms called CLINK and SLINK, which run in  $O(n^2)$  time. A package called ELKI has publicly available implementations.]

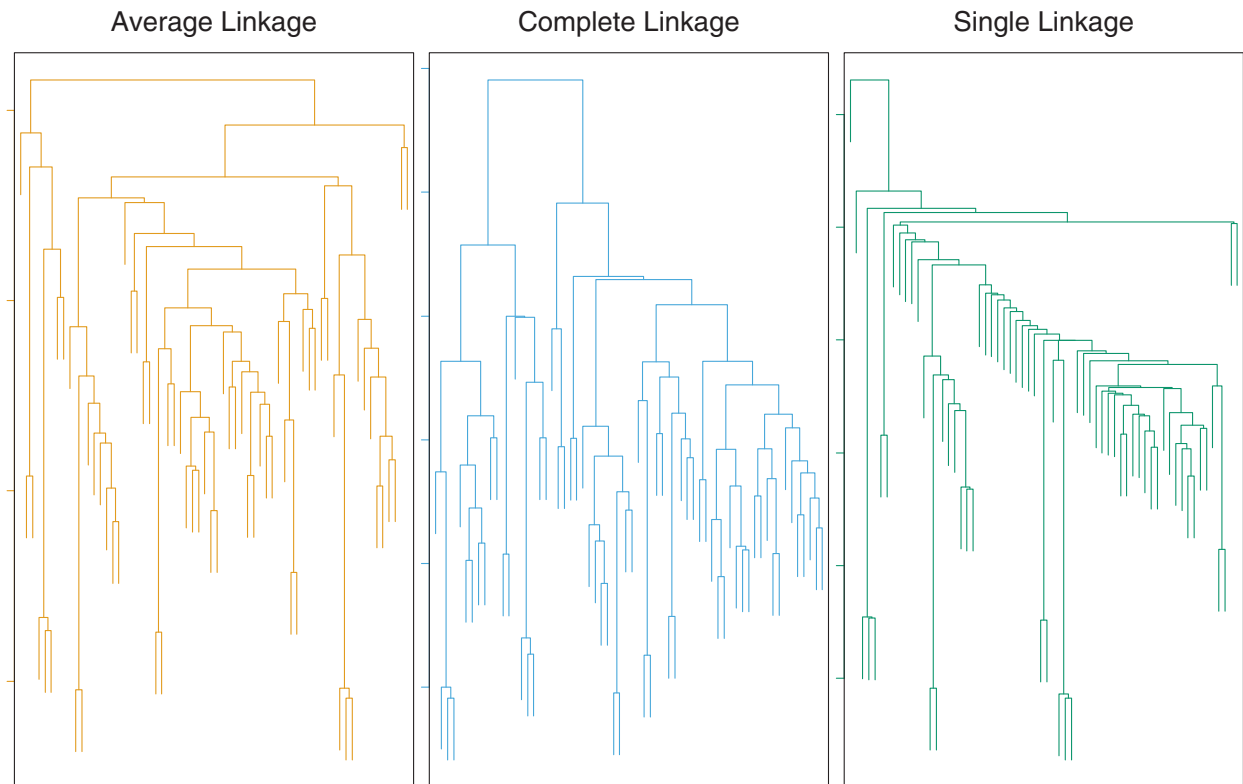
Dendrogram: Illustration of the cluster hierarchy (tree) in which the vertical axis encodes all the linkage distances.



dendrogram.pdf (ISL, Figure 10.9) [Example of a dendrogram cut into 1, 2, or 3 clusters.]

Cut dendrogram into clusters by horizontal line according to your choice of # of clusters OR intercluster distance.

[It's important to be aware that the horizontal axis of a dendrogram has no meaning. You could swap some treenode's left subtree and right subtree and it would still be the same dendrogram. It doesn't mean anything that two leaves happen to be next to each other.]



linkages.pdf (ISL, Figure 10.12) [Comparison of average, complete (max), and single (min) linkages. Observe that the complete linkage gives the best-balanced dendrogram, whereas the single linkage gives a very unbalanced dendrogram that is sensitive to outliers (especially near the top of the dendrogram).]

[Probably the worst of these is the single linkage, because it's very sensitive to outliers. Notice that if you cut this example into three clusters, two of them have only one sample point. It also tends to give you a very unbalanced tree.]

[The complete linkage tends to be the best balanced, because when a cluster gets large, the farthest point in the cluster is always far away. So large clusters are more resistant to growth than small ones. If balanced clusters are your goal, this is your best choice.]

[In most applications you probably want the average or complete linkage.]

Warning: centroid linkage can cause inversions where a parent cluster is fused at a lower height than its children.

[So statisticians don't like it, but nevertheless, centroid linkage is popular in genomics.]

[As a final note, all the clustering algorithms we've studied so far are unstable, in the sense that deleting a few sample points can sometimes give you very different results. But these unstable heuristics are still the most commonly used clustering algorithms. And it's not clear to me whether a truly stable clustering algorithm is even possible.]