

Machine Learning

An overview of unsupervised methods

William La Cava

Postdoctoral Researcher

Computational Genetics Laboratory

1acava@upenn.edu

October 10, 2018

Outline

1 Unsupervised Learning

Outline

- 1 Unsupervised Learning
- 2 Examples

Outline

- 1 Unsupervised Learning
- 2 Examples
- 3 K-Means

Outline

- 1 Unsupervised Learning
- 2 Examples
- 3 K-Means
- 4 Heirarchical Agglomerative Clustering

Outline

- 1 Unsupervised Learning
- 2 Examples
- 3 K-Means
- 4 Heirarchical Agglomerative Clustering
- 5 PCA

Outline

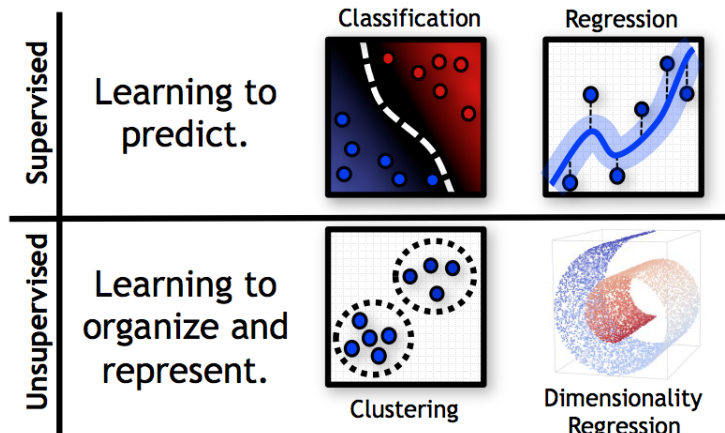
- 1 Unsupervised Learning
- 2 Examples
- 3 K-Means
- 4 Heirarchical Agglomerative Clustering
- 5 PCA
- 6 Examples

Outline

- 1 Unsupervised Learning
- 2 Examples
- 3 K-Means
- 4 Heirarchical Agglomerative Clustering
- 5 PCA
- 6 Examples
- 7 Conclusions

Machine Learning

Tasks



Unsupervised Learning

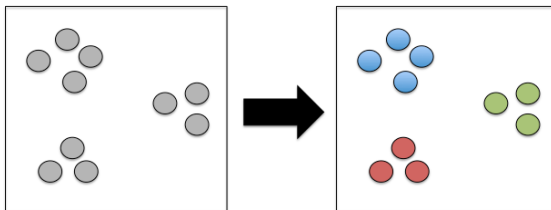
- Set of data: $\{\mathbf{x}_i, i = 1 \dots N\}$ with d features

Unsupervised Learning

- Set of data: $\{\mathbf{x}_i, i = 1 \dots N\}$ with d features

Definition (Clustering)

Given a set of data \mathbf{x} , partition the data cases into groups such that the data cases within each *partition* are more *similar* to each other than they are to data cases in other partitions.

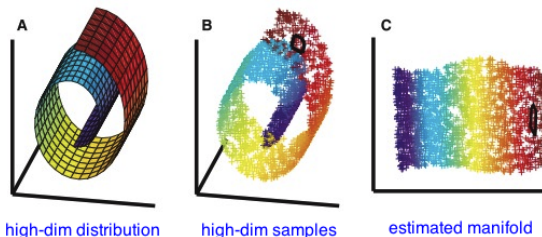


Unsupervised Learning

- Set of data: $\{\mathbf{x}_i, i = 1 \dots N\}$ with d features

Definition (Dimensionality Reduction)

Given a set of data $\mathbf{x} \in \mathbb{R}^d$, map the feature vectors into a lower dimensional space \mathbb{R}^k where $k < d$ while preserving certain properties of the data.



Examples

Supervised learning questions:

- *Clinical* What patient health characteristics are predictive of response to this treatment?
- *Genetics* For a cohort of patients, I have measured genotypes and the effective therapeutic dose of a drug. In new patients where I also measured genotypes, what dose should I use?

Examples

Unsupervised learning questions:

- *Clinical* Are there identifiable sub-groups of patients in my data (e.g., patients with similar demographics or that respond similarly to different treatments?)
- *Genetics* Are there patterns of gene expression in biopsies that I collected that suggest patients could be more precisely characterized in different molecular groups?

K-Means

- Attempts to group data into K clusters.
- Begins with randomly initialized centroids, $\mu_1 \dots \mu_K$
- Two step process:
 - 1 Calculate distance from every data case to each centroid, and assign clusters accordingly.
 - 2 Update cluster centers μ to the mean of the data cases assigned to them.
- Keep going until cluster positions stop changing.

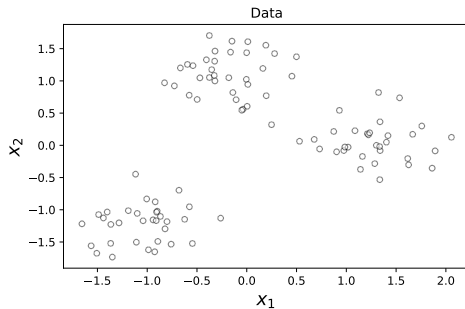
K-Means

- Minimizes the *within-cluster* variation:

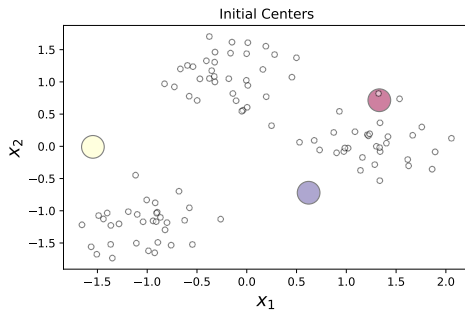
$$\mathcal{C}^* = \arg \min_{\mathcal{C}} \sum_{k=1}^K \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$

- K-Means converges to the local optima of its initial centroid positions.

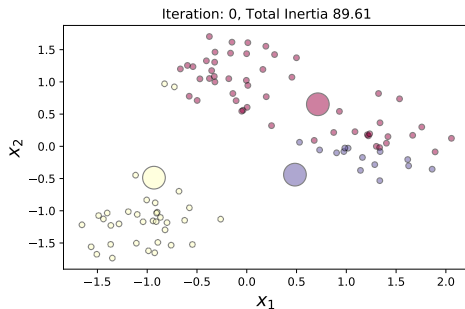
K-Means



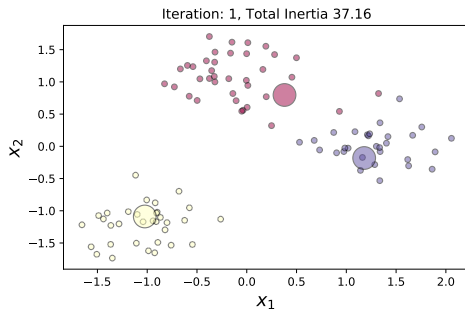
K-Means



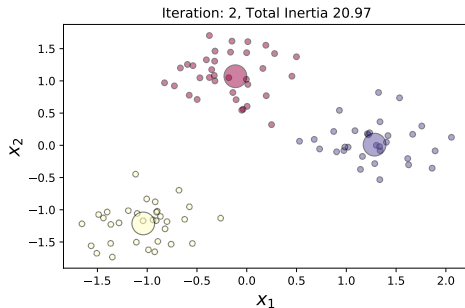
K-Means



K-Means

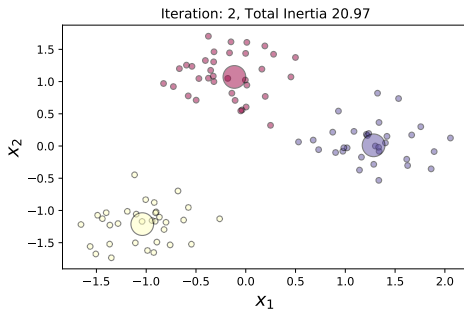


K-Means



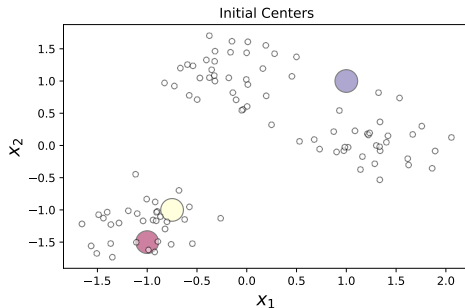
K-Means

Done!



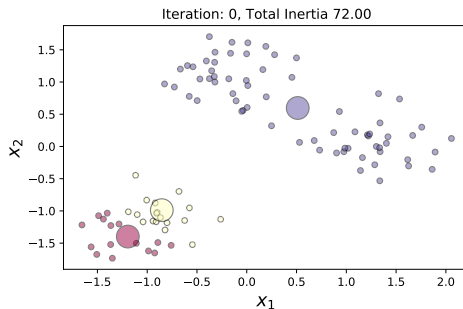
K-Means

Bad Initialization



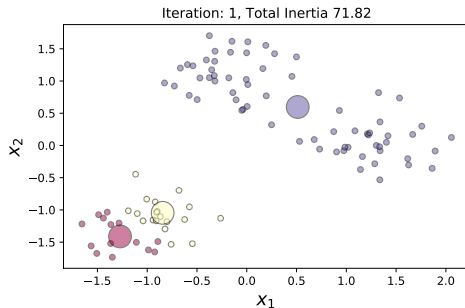
K-Means

Bad Initialization



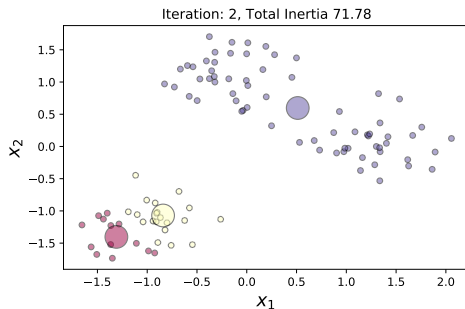
K-Means

Bad Initialization



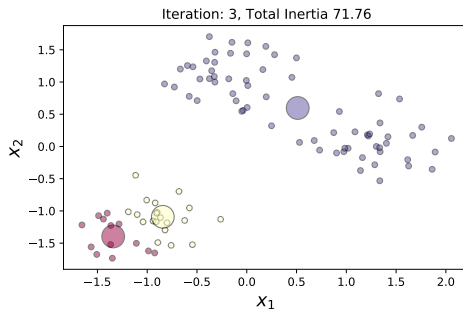
K-Means

Bad Initialization



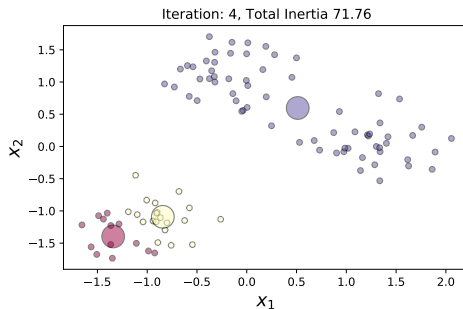
K-Means

Bad Initialization



K-Means

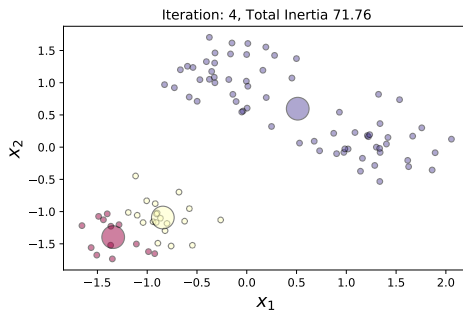
Bad Initialization



K-Means

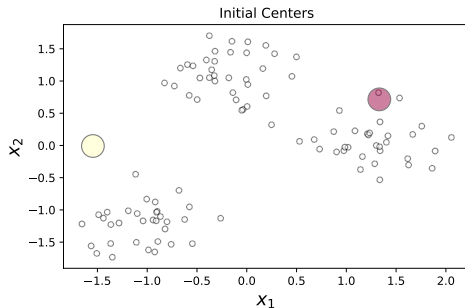
Bad Initialization

Done!



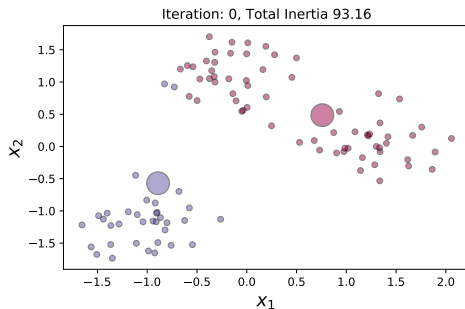
K-Means

Not Enough Clusters



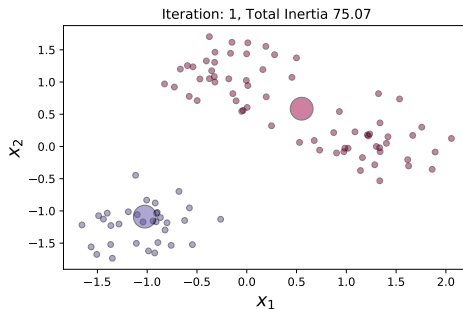
K-Means

Not Enough Clusters



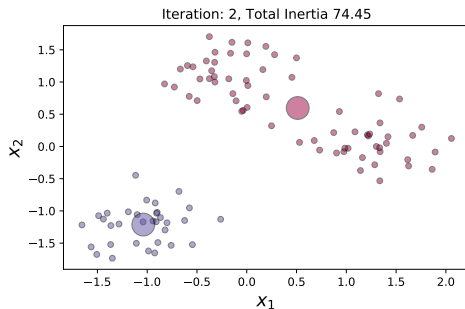
K-Means

Not Enough Clusters



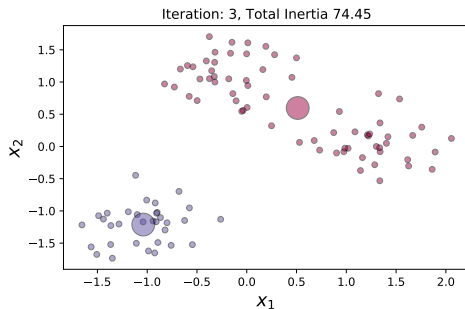
K-Means

Not Enough Clusters



K-Means

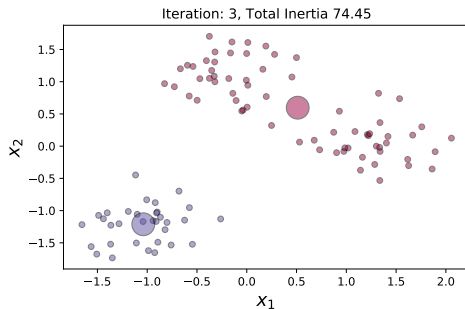
Not Enough Clusters



K-Means

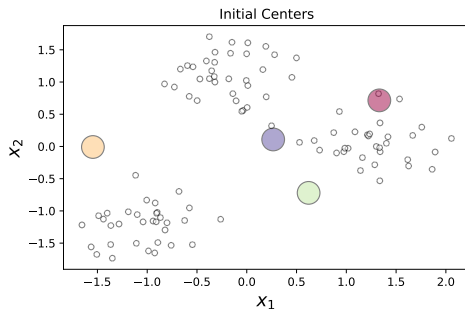
Not Enough Clusters

Done!



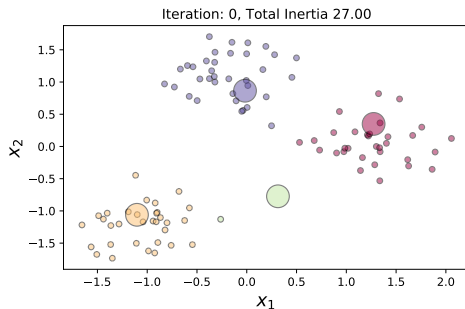
K-Means

Too Many Clusters



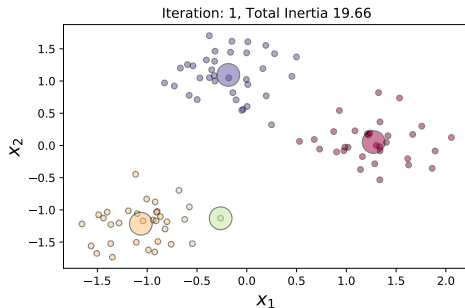
K-Means

Too Many Clusters



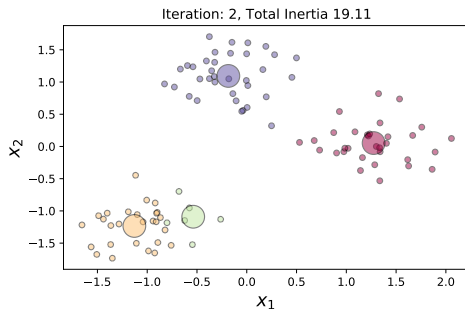
K-Means

Too Many Clusters



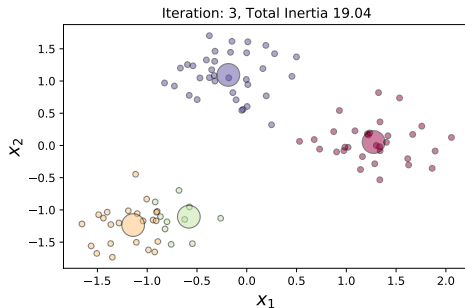
K-Means

Too Many Clusters



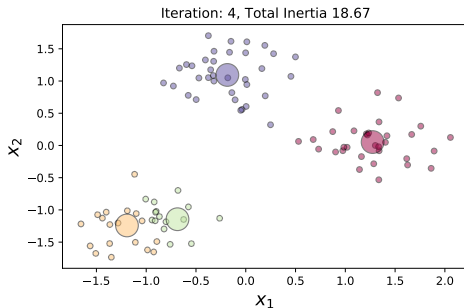
K-Means

Too Many Clusters



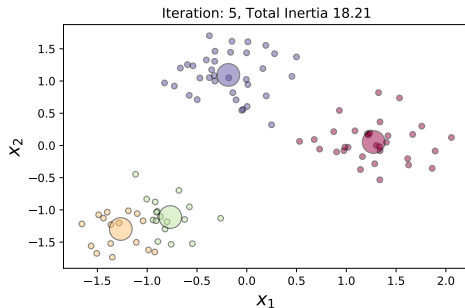
K-Means

Too Many Clusters



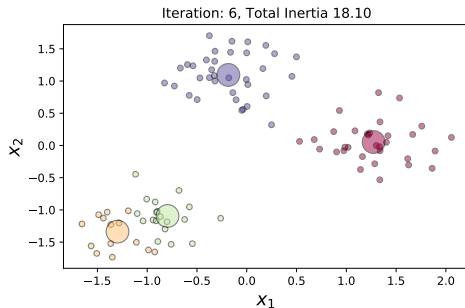
K-Means

Too Many Clusters



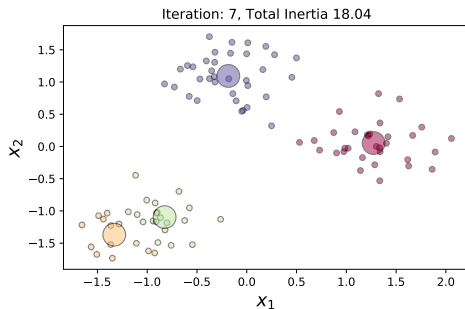
K-Means

Too Many Clusters



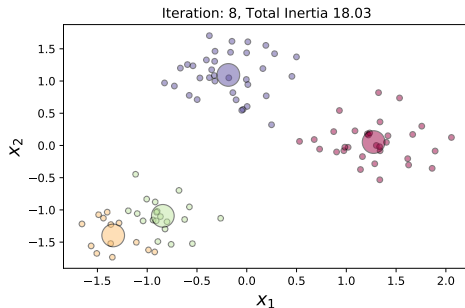
K-Means

Too Many Clusters



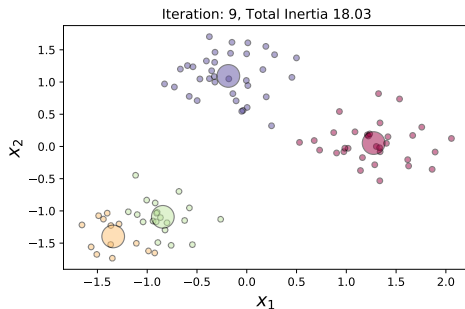
K-Means

Too Many Clusters



K-Means

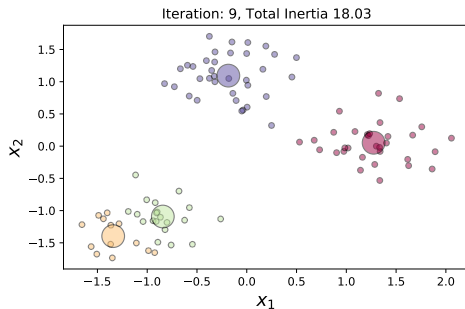
Too Many Clusters



K-Means

Too Many Clusters

Done!

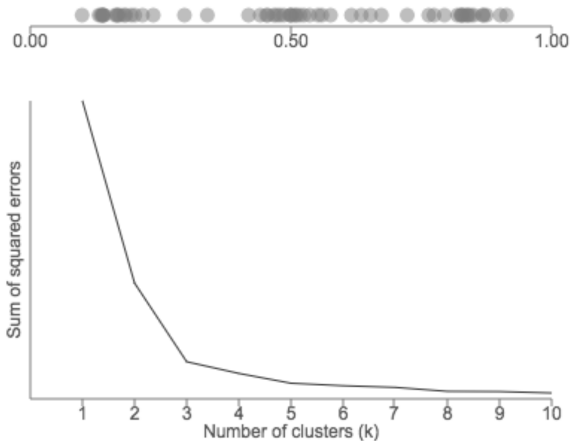


K-Means

- Sensitive to initial centroid positions.
- Sensitive to scaling of data dimensions.
- How to choose K ?

K-Means

Elbow Method



Hierarchical Agglomerative Clustering

- *Hierarchical Clustering*: greedy tree-based clustering methods
- *Hierarchical Agglomerative Clustering* (HAC): the most popular type
 - 1 Start with all data cases assigned to own clusters.
 - 2 Greedily and recursively merge pairs of clusters.

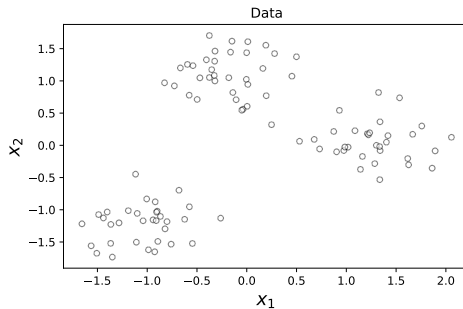
Hierarchical Agglomerative Clustering

Algorithm

HAC

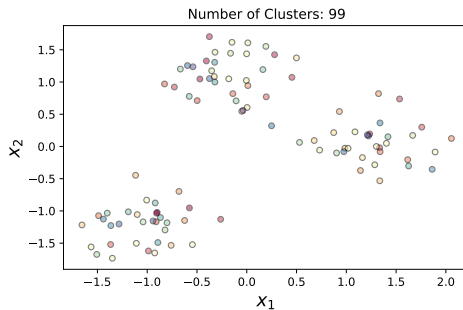
- 1 Start with all data cases assigned to own clusters.
- 2 Calculate all pairwise distances.
- 3 for $i = N, N - 1, \dots, 2$: $\leftarrow i = \text{number of clusters}$
 - 1 Merge the two closest clusters among i clusters.
 - 2 Calculate the pairwise distances between all $i - 1$ clusters.

HAC



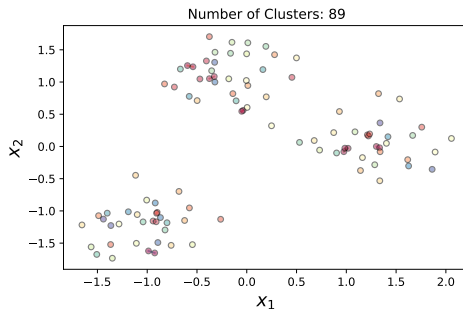
HAC

Example



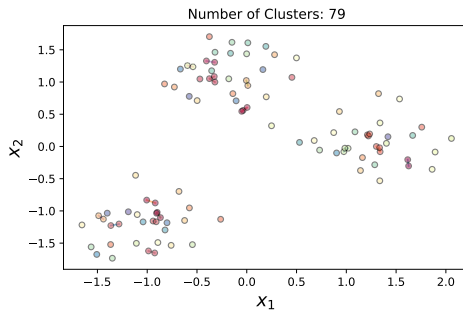
HAC

Example



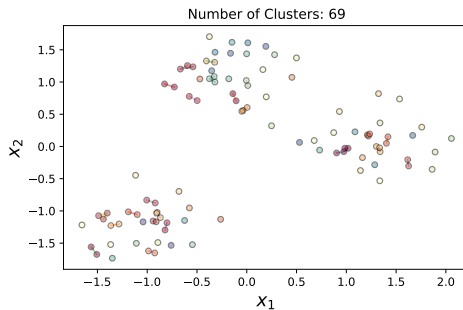
HAC

Example



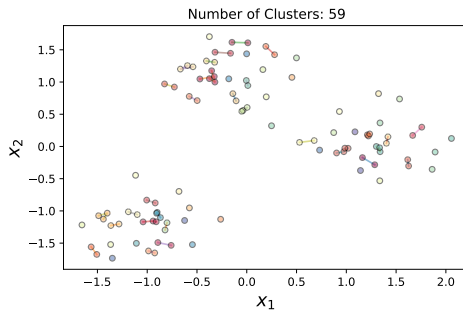
HAC

Example



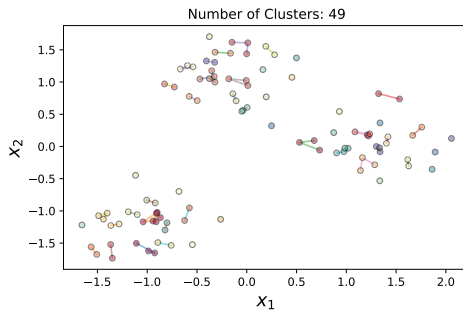
HAC

Example



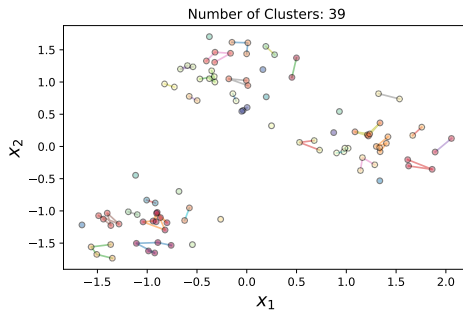
HAC

Example



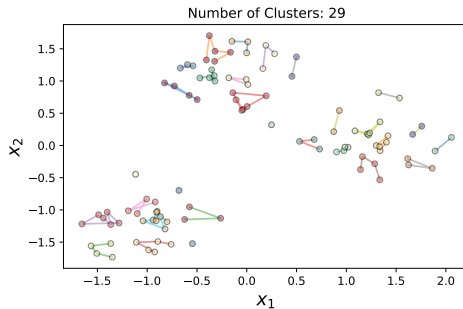
HAC

Example



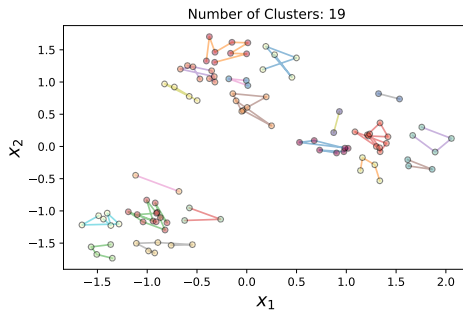
HAC

Example



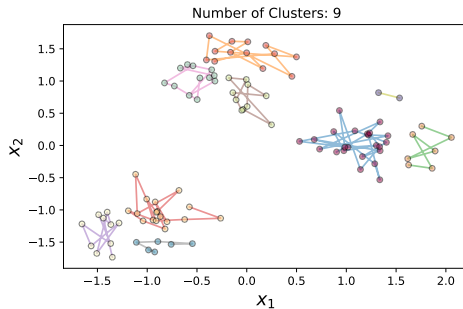
HAC

Example



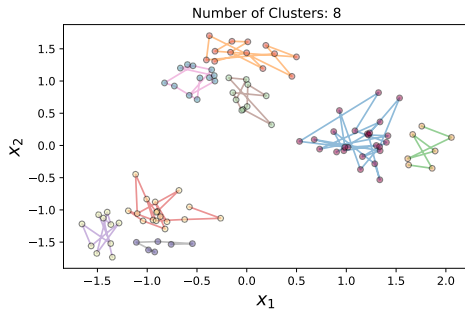
HAC

Example



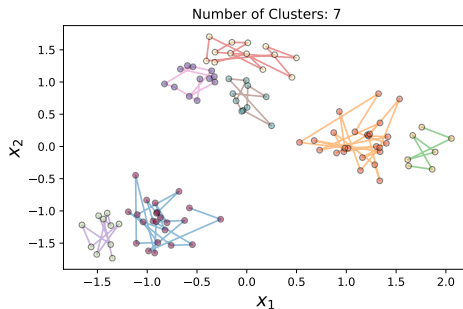
HAC

Example



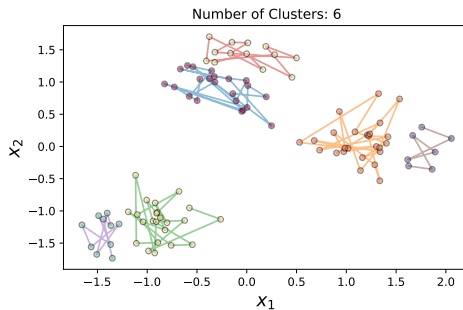
HAC

Example



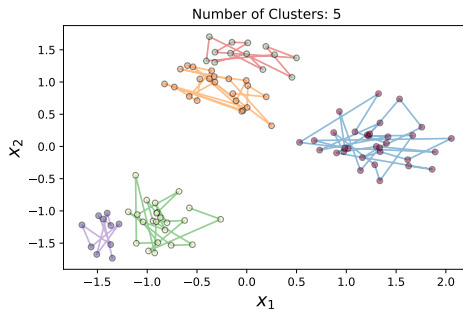
HAC

Example



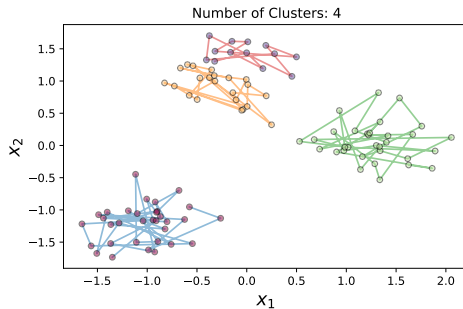
HAC

Example



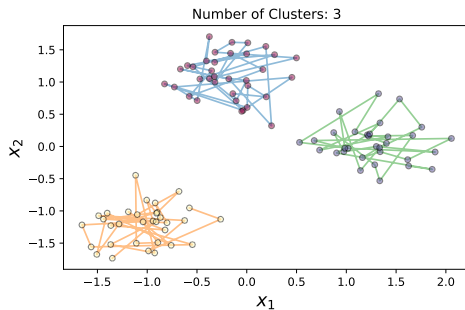
HAC

Example



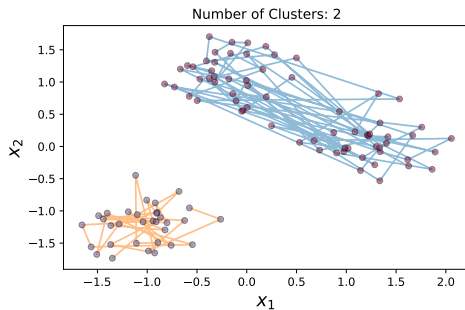
HAC

Example



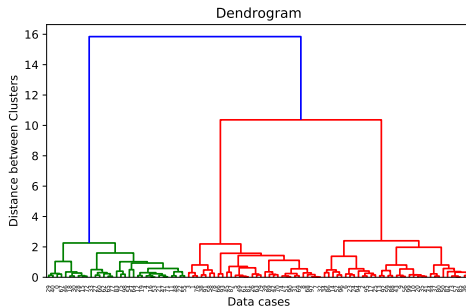
HAC

Example



HAC

Dendrogram



Hierarchical Agglomerative Clustering

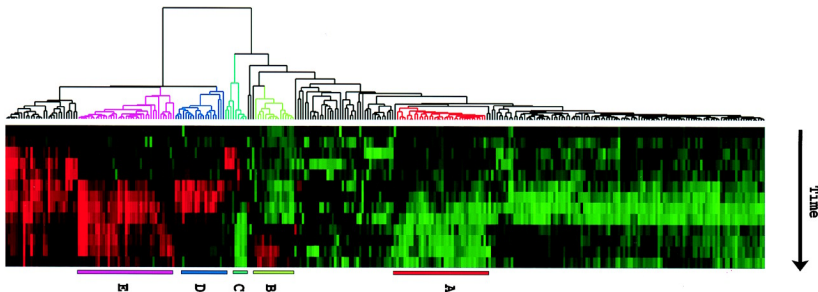
Issues

- 1 (like K-Means) Need good notion of similarity between clusters
- 2 Choose good 'linkage' function
- 3 (like K-Means) Sensitive to data scaling
- 4 Caution when interpreting results!

Clustering Examples

Gene Expression

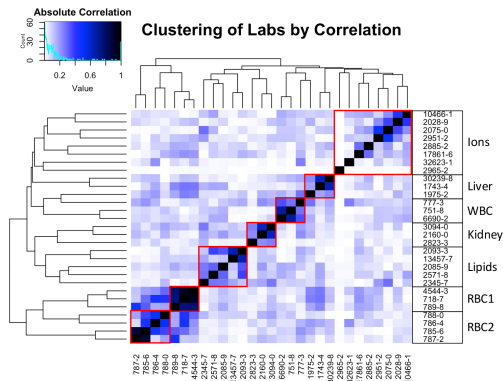
HAC used in gene expression analysis to visualize how genes group



Clustering Examples

Clinical Lab Measures

HAC and K-Means visualize lab measures in Geisinger EHR data
(Bauer et. al. 2016)

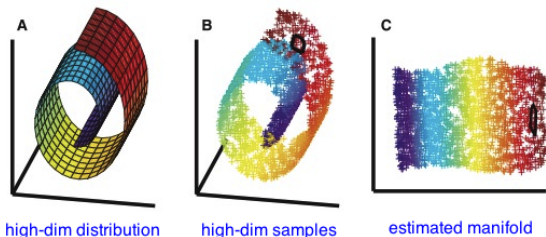


Unsupervised Learning

- Set of data: $\{\mathbf{x}_i, i = 1 \dots N\}$ with d features

Definition (Dimensionality Reduction)

Given a set of data $\mathbf{x} \in \mathbb{R}^d$, map the feature vectors into a lower dimensional space \mathbb{R}^k where $k < d$ while preserving certain properties of the data.



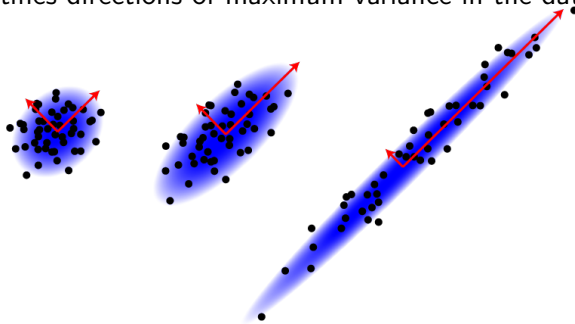
Principal Component Analysis

- 1 PCA assumes $\mathbf{x}_i \in \mathbb{R}^d$ lies on a k -dimensional linear manifold within \mathbb{R}^d .
- 2 In math,

$$\mathbf{X} = \mathbf{Z} \times \mathbf{B}$$

Principal Component Analysis

PCA identifies directions of maximum variance in the data.



Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}

Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}
- 2 Compute covariance matrix $\Sigma = \mathbf{X}^T \mathbf{X}$

Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}
- 2 Compute covariance matrix $\Sigma = \mathbf{X}^T \mathbf{X}$
- 3 Compute the k leading eigenvectors of Σ , $w_1 \dots w_k$

Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}
- 2 Compute covariance matrix $\Sigma = \mathbf{X}^T \mathbf{X}$
- 3 Compute the k leading eigenvectors of Σ , $w_1 \dots w_k$
- 4 Stack the eigenvectors into a $d \times k$ matrix \mathbf{W} where each column is an eigenvector

Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}
- 2 Compute covariance matrix $\Sigma = \mathbf{X}^T \mathbf{X}$
- 3 Compute the k leading eigenvectors of Σ , $w_1 \dots w_k$
- 4 Stack the eigenvectors into a $d \times k$ matrix \mathbf{W} where each column is an eigenvector
- 5 Compute the k -dimensional projection $\mathbf{Z} = \mathbf{XW}$

Principal Component Analysis

Steps:

- 1 Given centered $N \times d$ data matrix \mathbf{X}
- 2 Compute covariance matrix $\Sigma = \mathbf{X}^T \mathbf{X}$
- 3 Compute the k leading eigenvectors of Σ , $w_1 \dots w_k$
- 4 Stack the eigenvectors into a $d \times k$ matrix \mathbf{W} where each column is an eigenvector
- 5 Compute the k -dimensional projection $\mathbf{Z} = \mathbf{XW}$
- 6 To reconstruct \mathbf{X} from \mathbf{Z} and \mathbf{W} , compute $\hat{\mathbf{X}} = \mathbf{ZW}^T$

Principal Component Analysis

Why does this work?

- 1 Insight: Any real, symmetric matrix (like $\Sigma = \mathbf{X}^T \mathbf{X}$) can be decomposed into *eigenvectors* with corresponding *eigenvalues*

$$\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

Principal Component Analysis

Why does this work?

- 1 Insight: Any real, symmetric matrix (like $\Sigma = \mathbf{X}^T \mathbf{X}$) can be decomposed into *eigenvectors* with corresponding *eigenvalues*

$$\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

- 2 The maximum direction of variance in \mathbf{X} is the eigenvector of $\mathbf{X}^T \mathbf{X}$ with the largest eigenvalue.

Principal Component Analysis

Why does this work?

- 1 Insight: Any real, symmetric matrix (like $\Sigma = \mathbf{X}^T \mathbf{X}$) can be decomposed into *eigenvectors* with corresponding *eigenvalues*

$$\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

- 2 The maximum direction of variance in \mathbf{X} is the eigenvector of $\mathbf{X}^T \mathbf{X}$ with the largest eigenvalue.
- 3 The k biggest directions of variance in \mathbf{X} are the eigenvectors of $\mathbf{X}^T \mathbf{X}$ with the k largest eigenvalues.

Principal Component Analysis

Why does this work?

- 1 Insight: Any real, symmetric matrix (like $\Sigma = \mathbf{X}^T \mathbf{X}$) can be decomposed into *eigenvectors* with corresponding *eigenvalues*

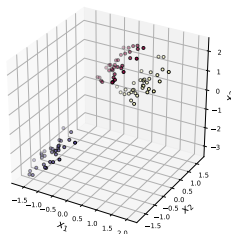
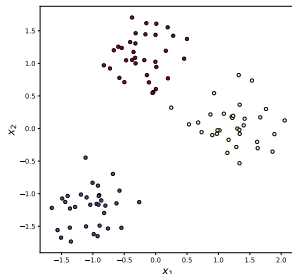
$$\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

- 2 The maximum direction of variance in \mathbf{X} is the eigenvector of $\mathbf{X}^T \mathbf{X}$ with the largest eigenvalue.
- 3 The k biggest directions of variance in \mathbf{X} are the eigenvectors of $\mathbf{X}^T \mathbf{X}$ with the k largest eigenvalues.
- 4 eigenvectors are orthogonal to each other, so the data projected into \mathbf{Z} will be linearly independent.

Principal Component Analysis

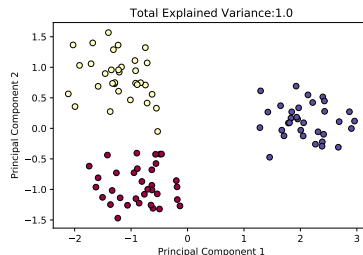
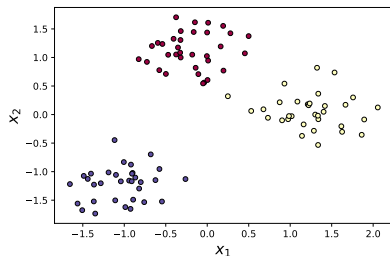
Example

Our cluster data has an extra dimension ($x_3 = x_1 + x_2$)



Principal Component Analysis

Example



Principal Component Analysis

Uses

- 1 Each principal component is a linear combination of columns in \mathbf{X} .

Principal Component Analysis

Uses

- 1 Each principal component is a linear combination of columns in \mathbf{X} .
- 2 PCA is useful for visualization, because it can plot high-dimensional data along its first two directions of maximum variance.

Principal Component Analysis

Uses

- 1 Each principal component is a linear combination of columns in \mathbf{X} .
- 2 PCA is useful for visualization, because it can plot high-dimensional data along its first two directions of maximum variance.
- 3 Data should be centered and scaled to unit-variance.

Principal Component Analysis

Uses

- 1 Each principal component is a linear combination of columns in \mathbf{X} .
- 2 PCA is useful for visualization, because it can plot high-dimensional data along its first two directions of maximum variance.
- 3 Data should be centered and scaled to unit-variance.
- 4 “Variance explained” gives a measure of how well the data dimensionality can be reduced.

Principal Component Analysis

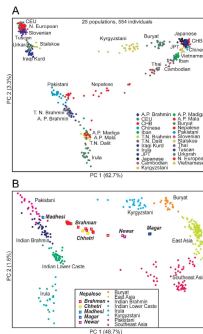
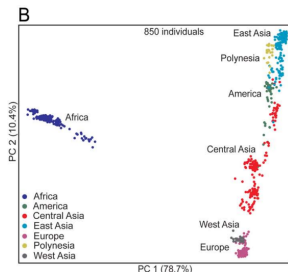
Uses

- 1 Each principal component is a linear combination of columns in \mathbf{X} .
- 2 PCA is useful for visualization, because it can plot high-dimensional data along its first two directions of maximum variance.
- 3 Data should be centered and scaled to unit-variance.
- 4 “Variance explained” gives a measure of how well the data dimensionality can be reduced.
- 5 *loading*: how much each variable in \mathbf{x} corresponds to the principal components

Examples

Genomes

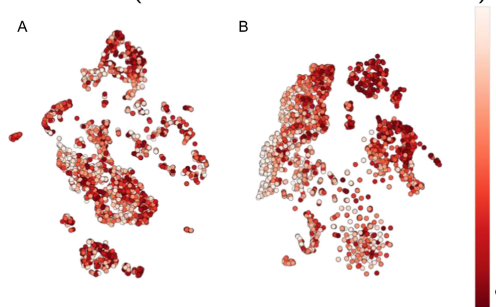
Large-scale variation in human genomes corresponds to the history of our species (Xing et. al. 2010)



Examples

ALS

Learning low-order representations can improve supervised learning methods (Beaulieu-Jones et. al. 2016)



A: PCA followed by t-SNE. B: t-SNE of a 250 node auto-encoder. Color: days survived.

Conclusions

- Other unsupervised learning algorithms
 - t-SNE
 - Mixture Models
 - Multidimensional Scaling (MDS)
 - Non-negative Matrix Factorization (NMF)
 - Auto-encoders
- mail me for lecture examples or questions
lacava@upenn.edu