

K-means and Hierarchical Clustering

Schwartz

August 13, 2017

Best. Music. EVAR

1. Elliott Smith
2. Sufjan Stevens
3. Iron and Wine
4. Damien Rice

1. Die Antwoord
2. Kendrick Lamar
3. Dan le Sac Vs Scroobius Pip

1. Rufus Wainwright
2. Lyle Lovett
3. Julie Doiron

1. D'Gary
2. Kishi Bashi
3. Christine and the Queens
4. Beirut

1. Rage Against the Machine
2. System of a Down
3. Smashing Pumpkins

1. Beck
2. Cake
3. Beastie Boys

Objectives

- ▶ Review *Supervised* versus *Unsupervised*

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ **K-means** (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*
 - ▶ Elbow, Silhouette, and Gap methods

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*
 - ▶ Elbow, Silhouette, and Gap methods
- ▶ Hierarchical clustering

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*
 - ▶ Elbow, Silhouette, and Gap methods
- ▶ Hierarchical clustering
- ▶ DBSCAN

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*
 - ▶ Elbow, Silhouette, and Gap methods
- ▶ Hierarchical clustering
- ▶ DBSCAN
- ▶ Bayesian mixture models

Objectives

- ▶ Review *Supervised* versus *Unsupervised*
- ▶ *K*-means (not to be confused with *KNN*)
 - ▶ and the curse of dimensionality
 - ▶ *norms* and the curse of dimensionality
 - ▶ more about the curse of dimensionality
- ▶ Choosing *K*
 - ▶ Elbow, Silhouette, and Gap methods
- ▶ Hierarchical clustering
- ▶ DBSCAN
- ▶ Bayesian mixture models
- ▶ Expectation-Maximization (EM) algorithm

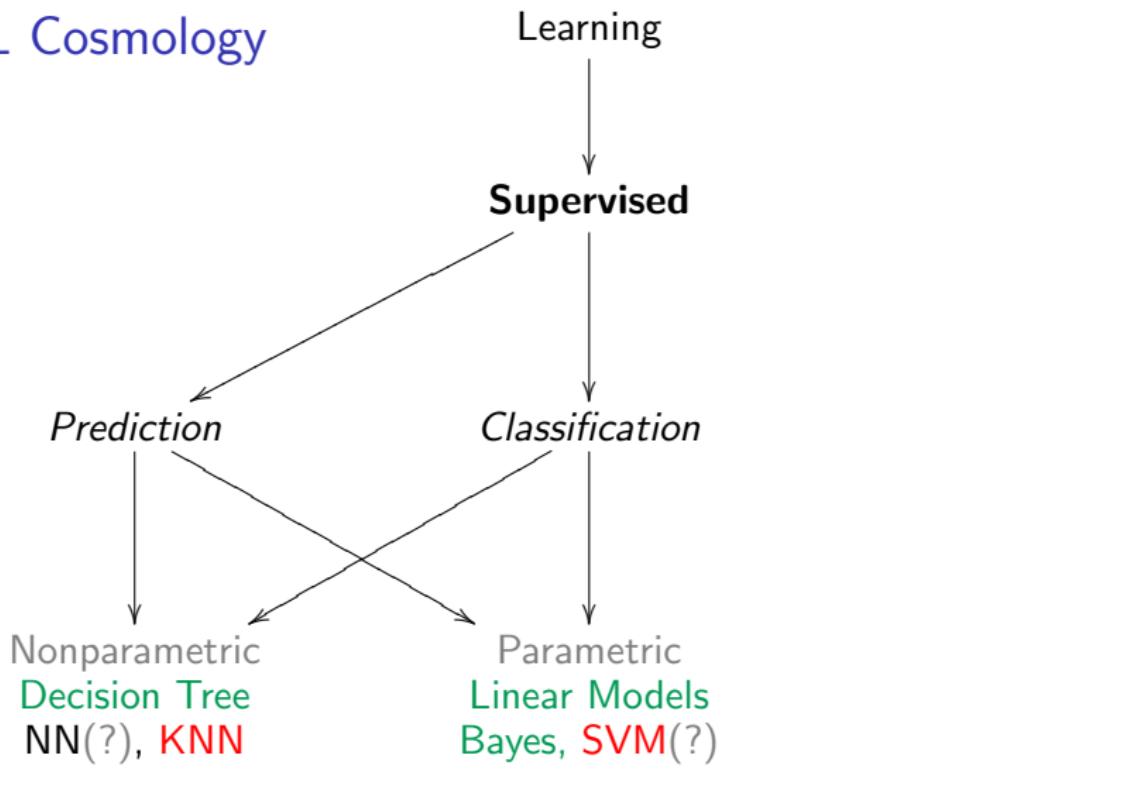
Learning



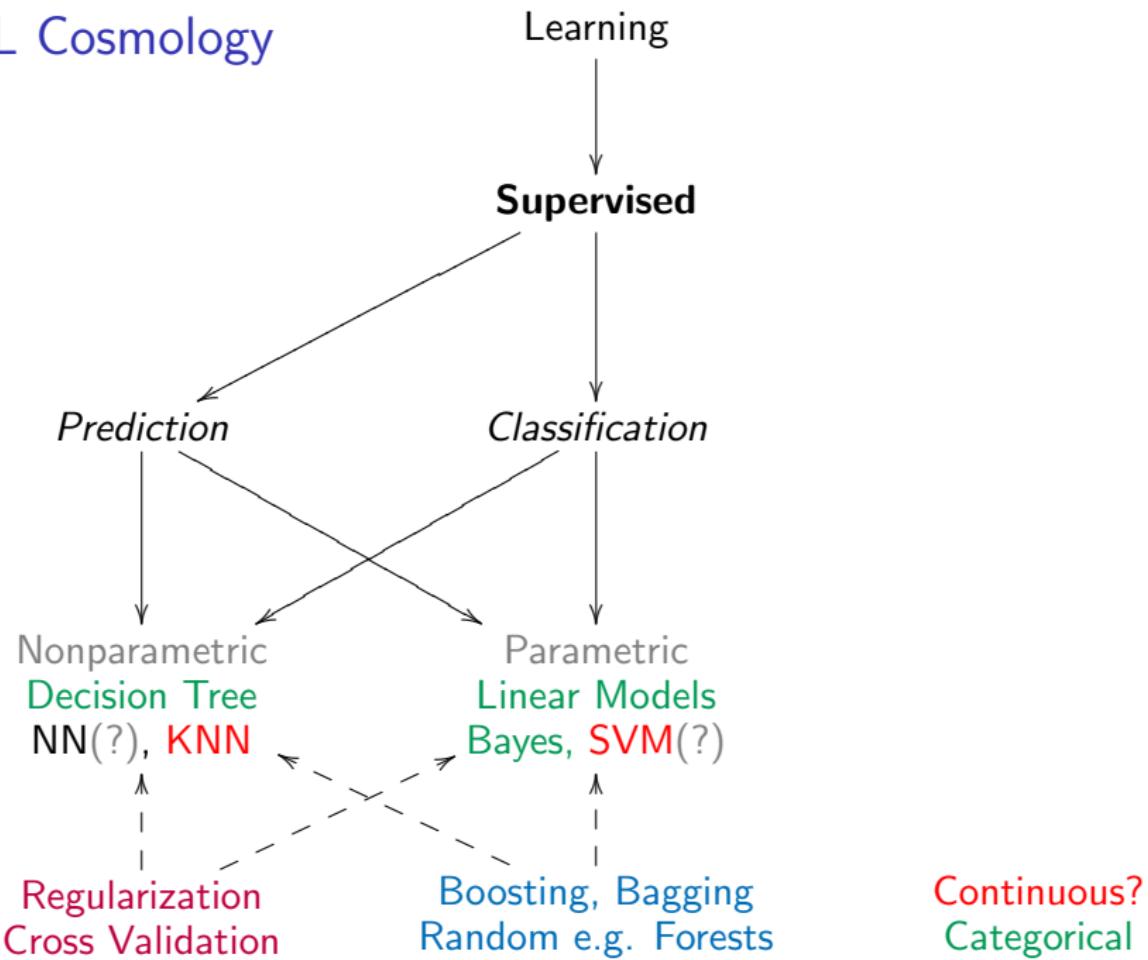
Supervised

Prediction

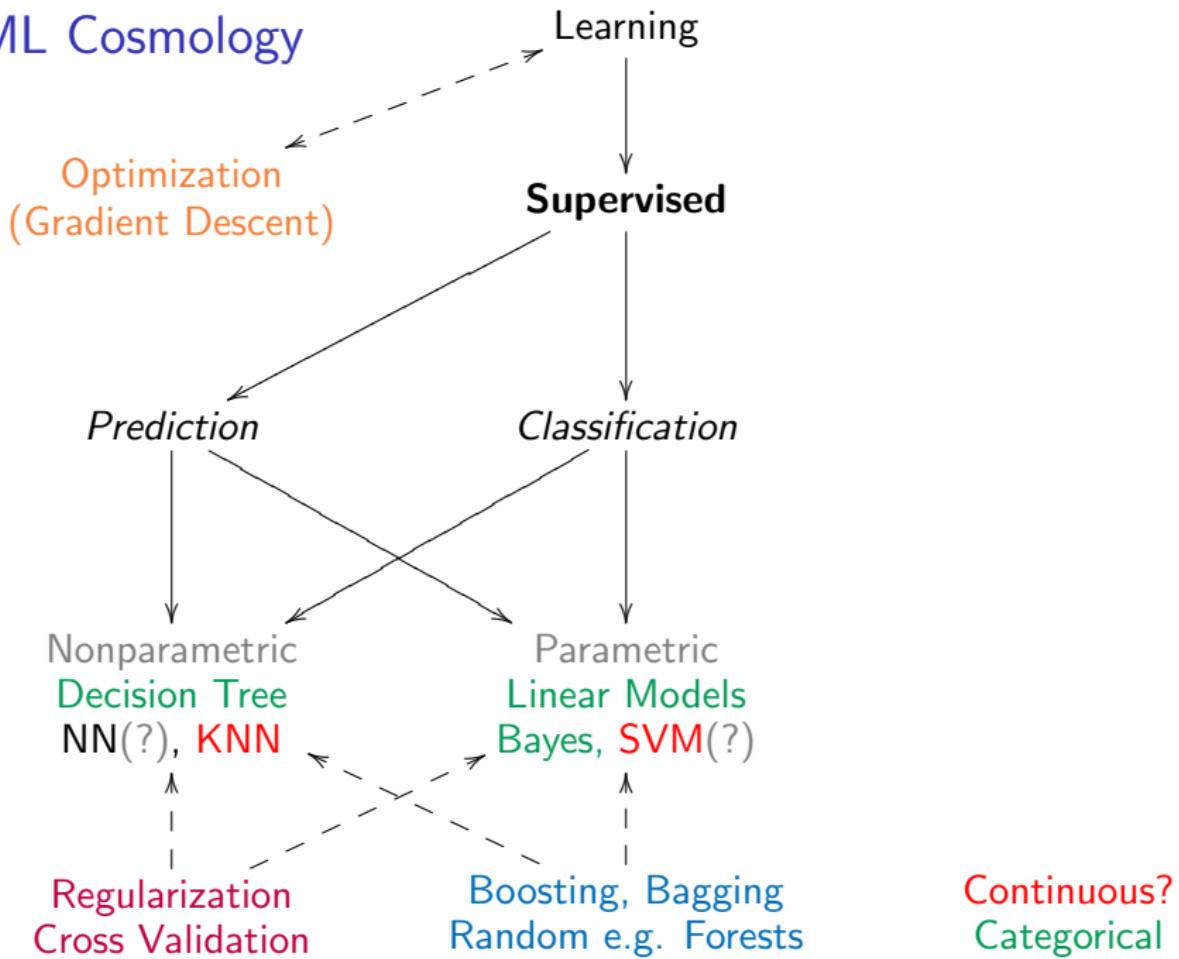
Classification



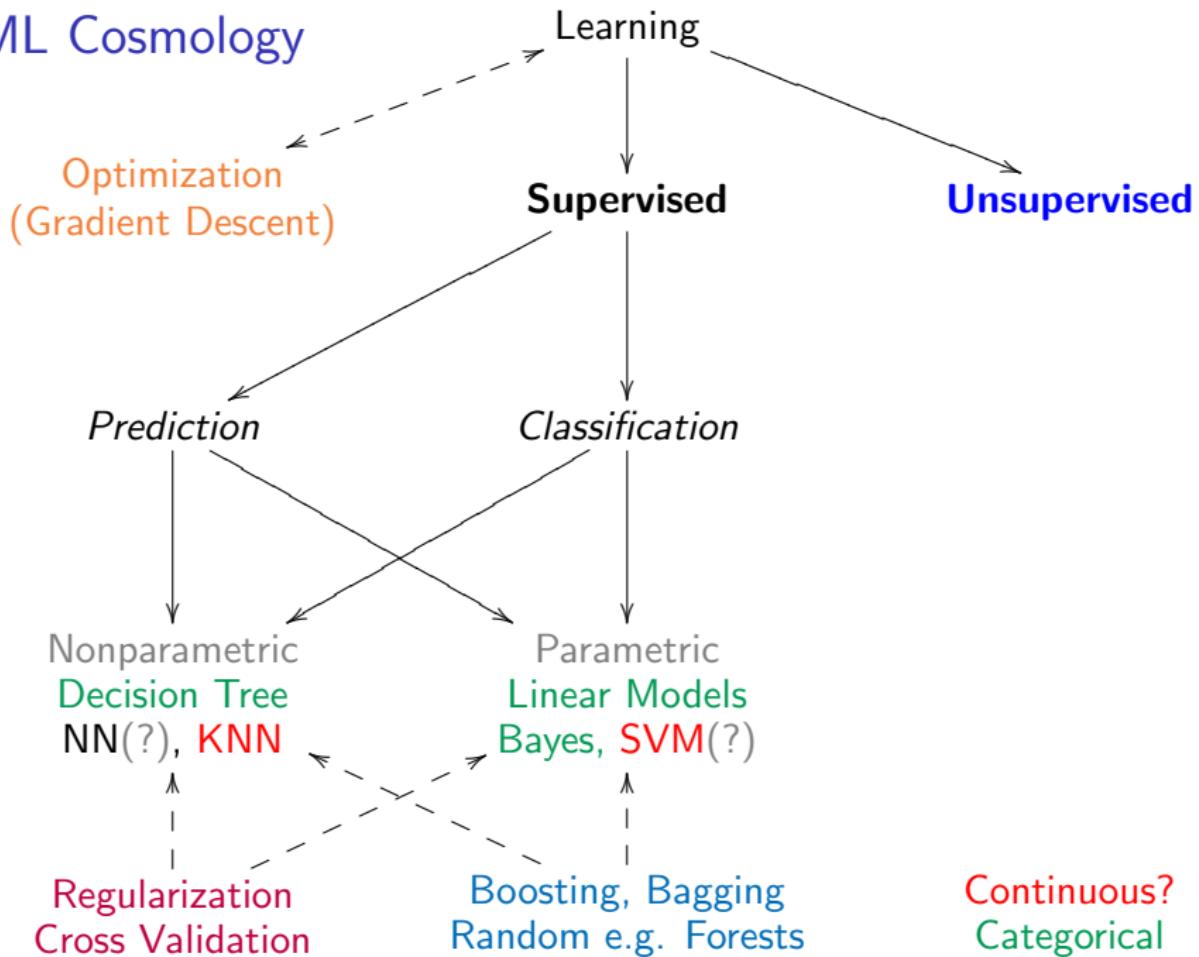
Continuous?
Categorical



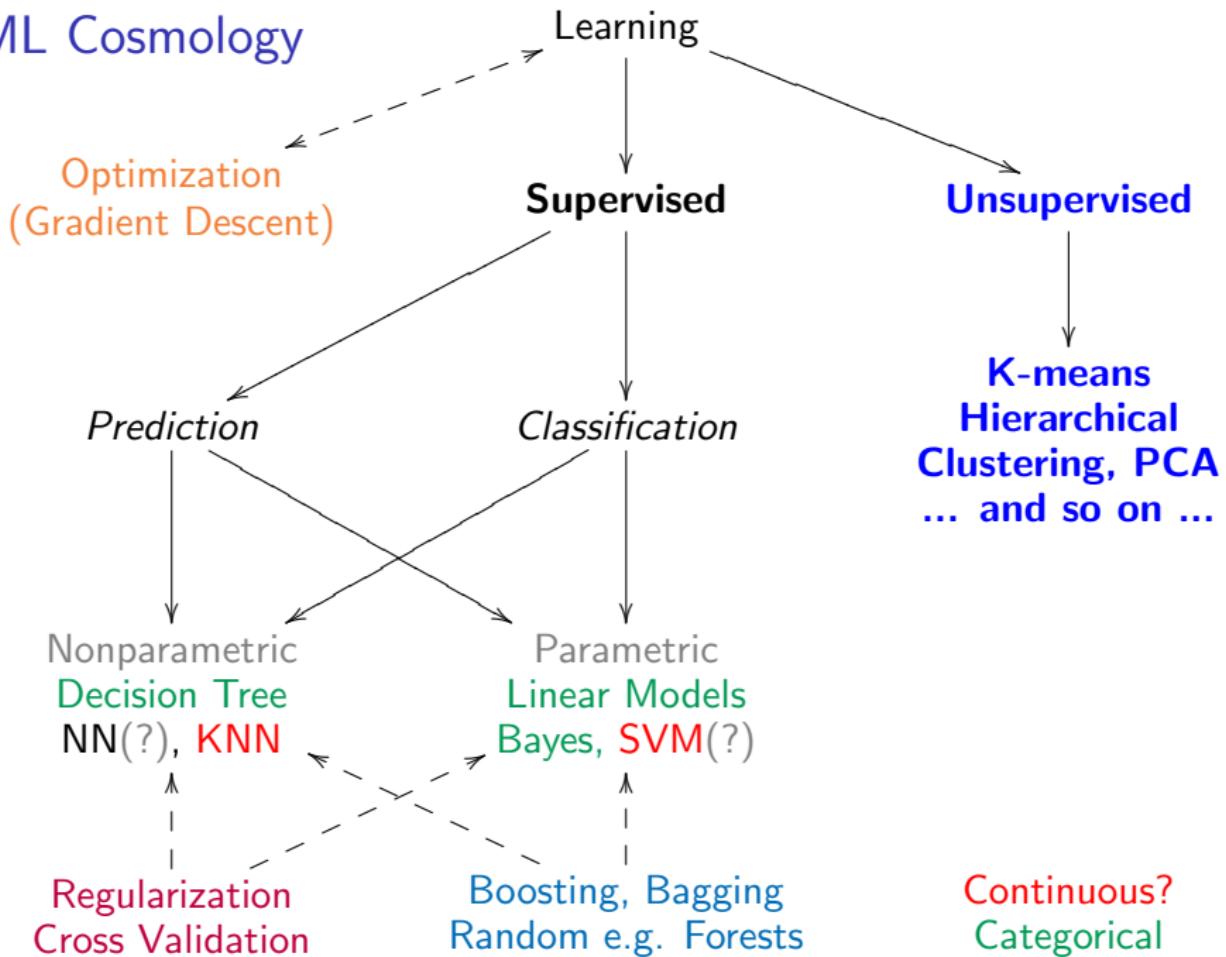
ML Cosmology



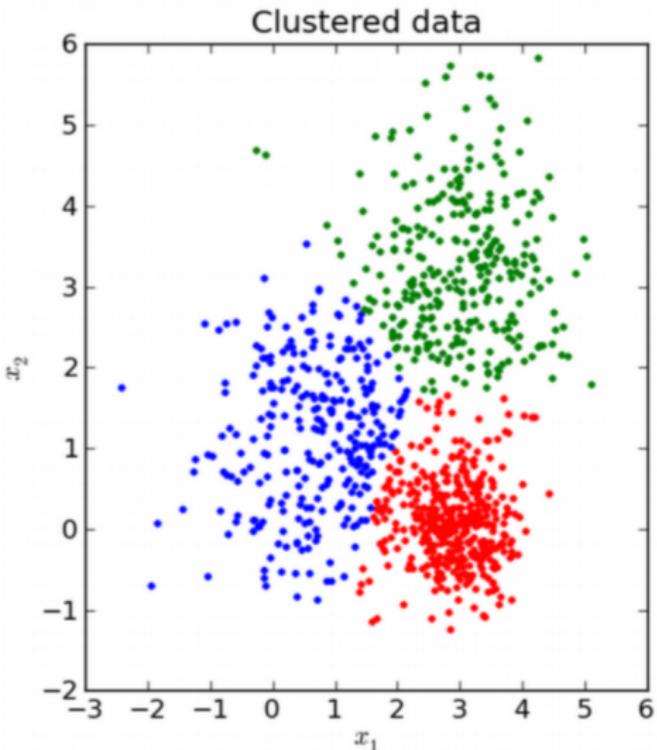
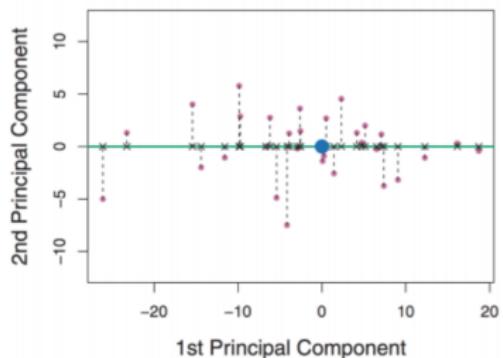
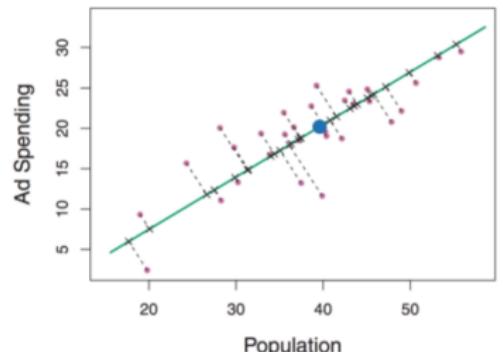
ML Cosmology



ML Cosmology



Unsupervised learning

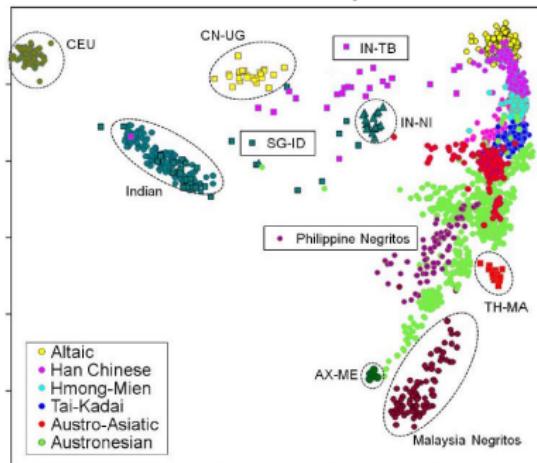


Low dimensional representations
of data capturing data variation

Homogeneous subgroups
capturing data substructure

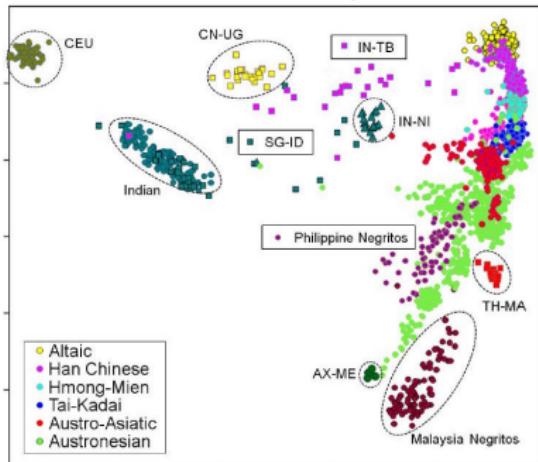
Supervised versus Unsupervised

Human Genetic Diversity in Asia Two-Dimensional Representation



Supervised versus Unsupervised

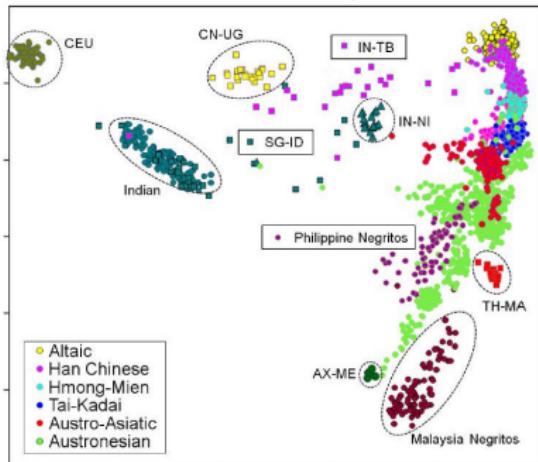
Human Genetic Diversity in Asia Two-Dimensional Representation



Unsupervised learning provides

Supervised versus Unsupervised

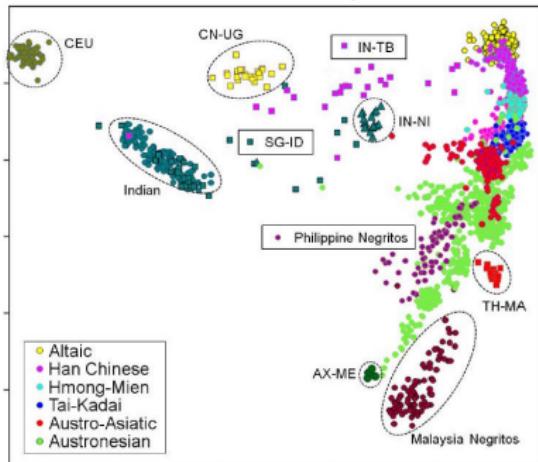
Human Genetic Diversity in Asia Two-Dimensional Representation



- Unsupervised learning provides
- ▶ exploratory data analysis (EDA) to look at/uncover feature structure

Supervised versus Unsupervised

Human Genetic Diversity in Asia Two-Dimensional Representation

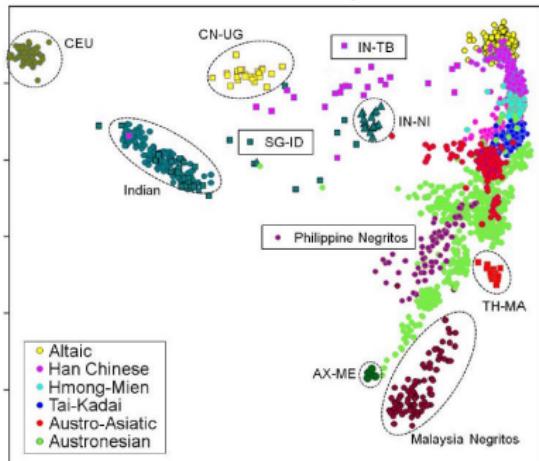


Unsupervised learning provides

- ▶ exploratory data analysis (EDA) to look at/uncover feature structure
- ▶ anomaly detection to provide data quality control (QC)

Supervised versus Unsupervised

Human Genetic Diversity in Asia Two-Dimensional Representation

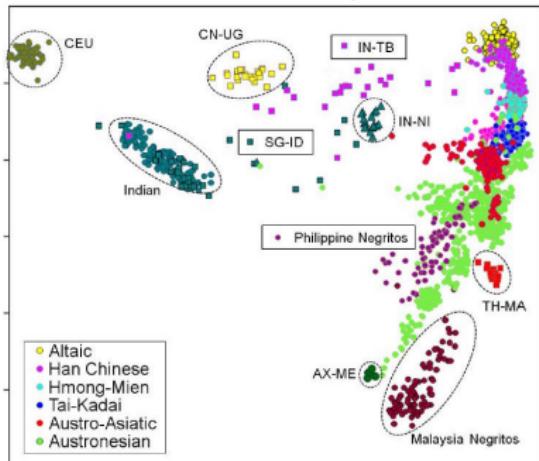


Unsupervised learning provides

- ▶ exploratory data analysis (EDA) to look at/uncover feature structure
- ▶ anomaly detection to provide data quality control (QC)
- ▶ dimensionality reduction to simplify large feature spaces

Supervised versus Unsupervised

Human Genetic Diversity in Asia Two-Dimensional Representation



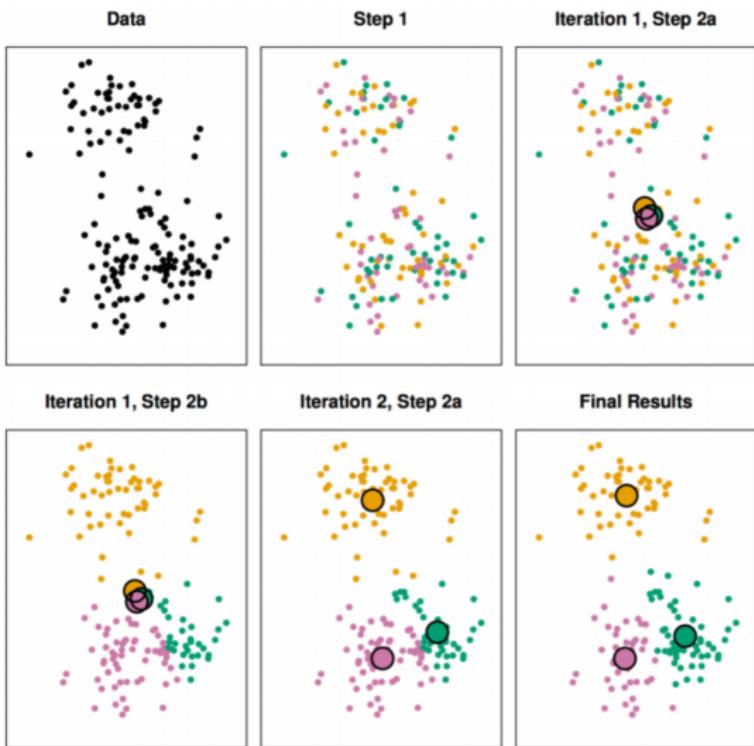
Unsupervised learning provides

- ▶ exploratory data analysis (EDA) to look at/uncover feature structure
- ▶ anomaly detection to provide data quality control (QC)
- ▶ dimensionality reduction to simplify large feature spaces

“Labels” are sometimes used to facilitate these objectives but unlike supervised learning, unsupervised learning is not necessarily immediately concerned with predicting outcomes (labels, targets, Y , dependent/endogenous variables)

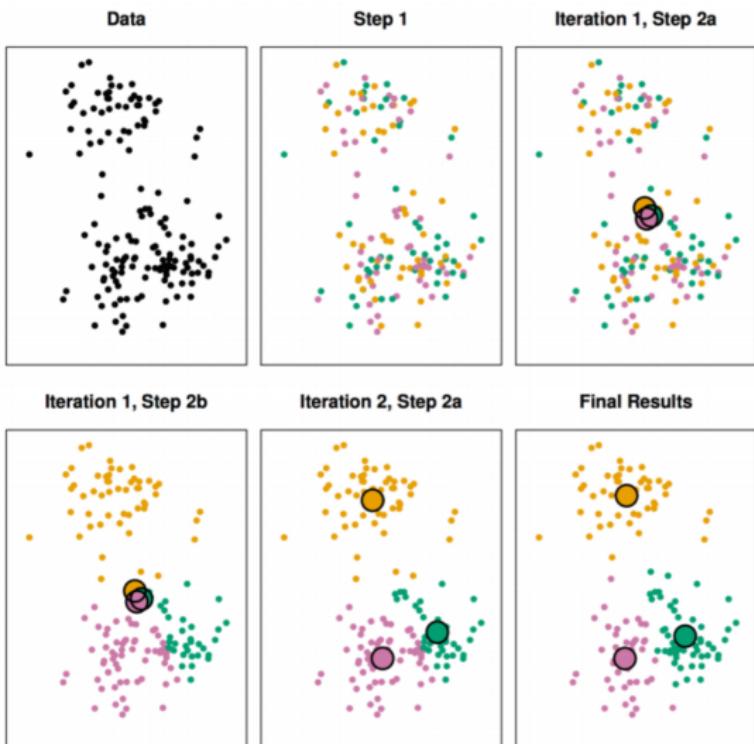
K-Means

1. Choose *number of clusters*, K
2. Randomly assign data to each cluster k
3. Compute the centroid for each cluster k
4. Assign data to the cluster with the closest centroid
5. Return to step 3, unless the centroids have stabilized



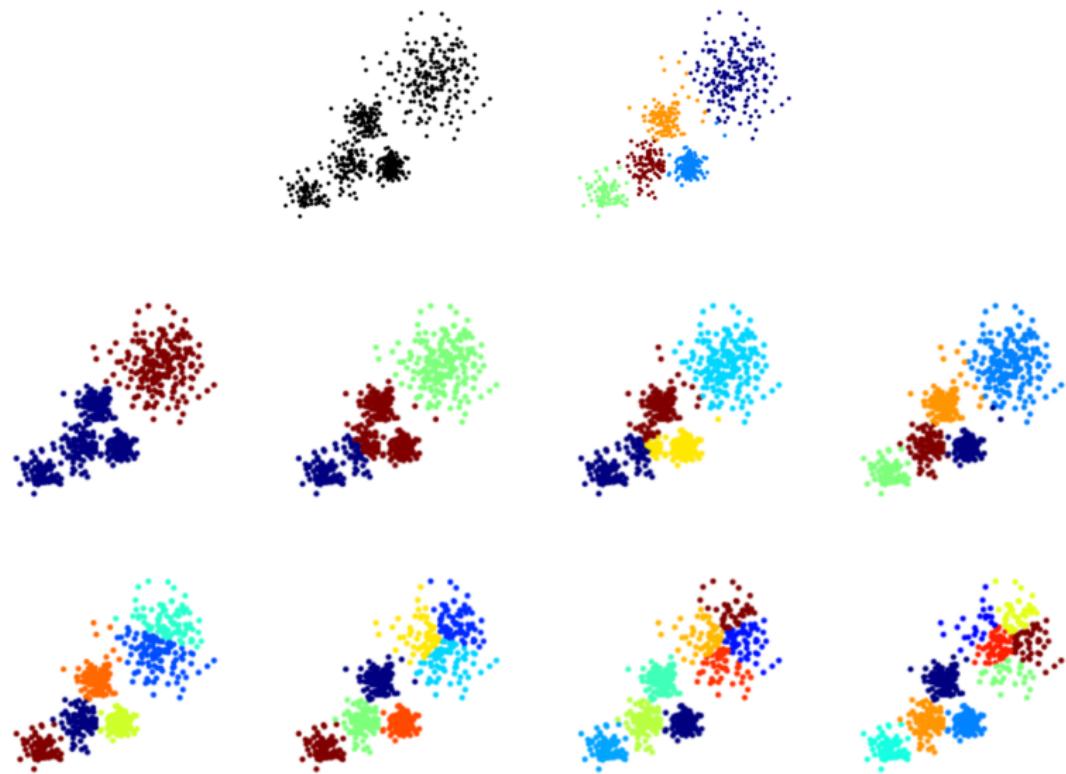
K-Means

1. Choose *number of clusters*, K
2. Randomly assign data to each cluster k
3. Compute the centroid for each cluster k
4. Assign data to the cluster with the closest centroid
5. Return to step 3, unless the centroids have stabilized



Other init. ideas?

K?



Elbow and Silhouette methods

For some clustering

$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

Elbow and Silhouette methods

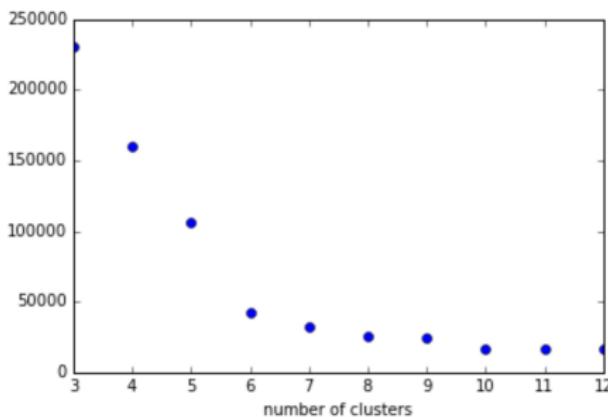
For some clustering

$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

Select based on diminishing returns



Elbow and Silhouette methods

For some clustering

$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

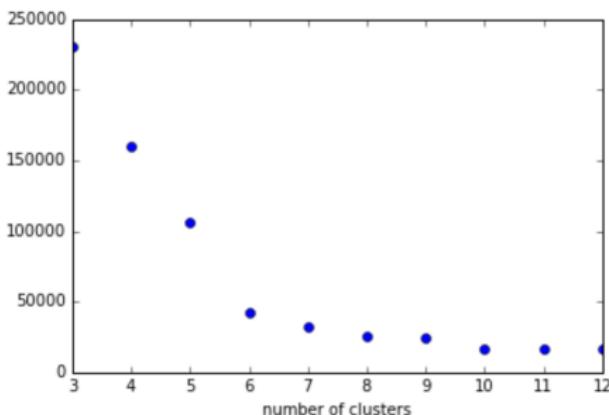
$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

$\bar{\Delta}_i^k$: x_i 's dissimilarity in cluster k

$\bar{\Delta}_i^{k'}$: x_i 's dissimilarity to cluster k'
(x_i in, & closest to clusters k, k')

$$\text{Silhouette}(i) = \frac{\bar{\Delta}_i^{k'} - \bar{\Delta}_i^k}{\max(\bar{\Delta}_i^{k'}, \bar{\Delta}_i^k)}$$

Select based on diminishing returns



Elbow and Silhouette methods

For some clustering

$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

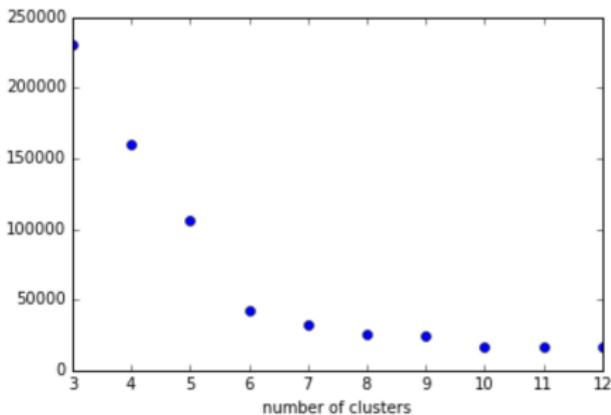
$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

$\bar{\Delta}_i^k$: x_i 's dissimilarity in cluster k

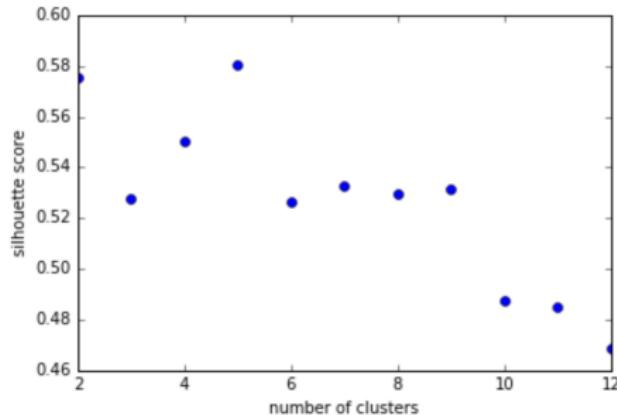
$\bar{\Delta}_i^{k'}$: x_i 's dissimilarity to cluster k'
(x_i in, & closest to clusters k, k')

$$\text{Silhouette}(i) = \frac{\bar{\Delta}_i^{k'} - \bar{\Delta}_i^k}{\max(\bar{\Delta}_i^{k'}, \bar{\Delta}_i^k)}$$

Select based on diminishing returns



Compare average silhouette scores



Elbow and Silhouette methods

For some clustering

$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

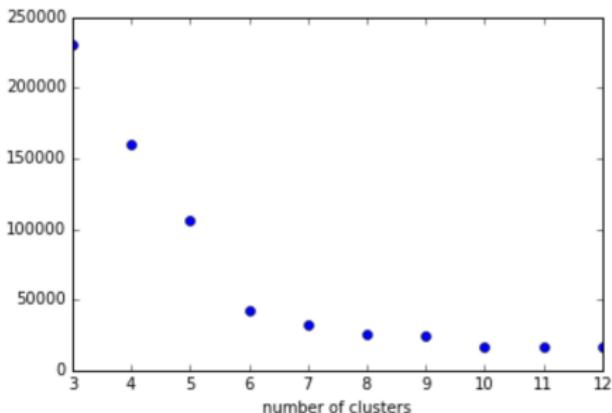
$\bar{\Delta}_i^k$: x_i 's dissimilarity in cluster k

$\bar{\Delta}_i^{k'}$: x_i 's dissimilarity to cluster k'
(x_i in, & closest to clusters k, k')

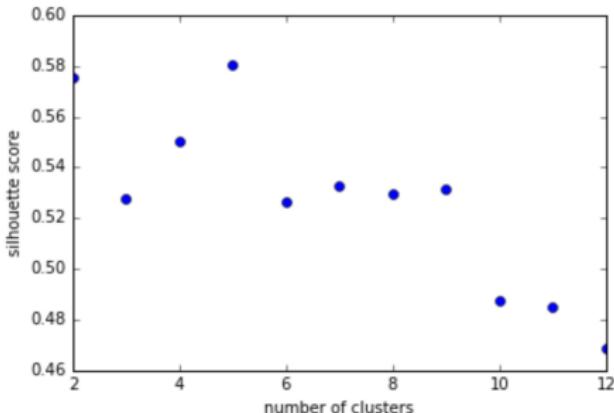
$$\text{Silhouette}(i) = \frac{\bar{\Delta}_i^{k'} - \bar{\Delta}_i^k}{\max(\bar{\Delta}_i^{k'}, \bar{\Delta}_i^k)}$$

Interpret silhouette scores?

Select based on diminishing returns



Compare average silhouette scores



Elbow and Silhouette methods

For some clustering

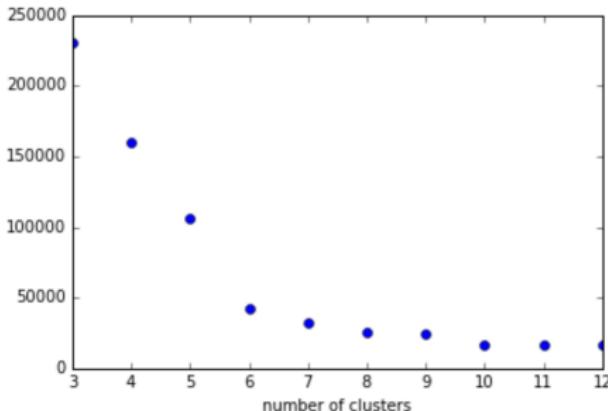
$$C_K(i) \mapsto \{1, 2, \dots, K\}$$

clustering fit can be measured as

$$W(C_K) = \frac{1}{K} \sum_{C_K(i)=C_K(j)=k} \|x_i - x_j\|^2$$

Does the scaling of the x matter?

Select based on diminishing returns



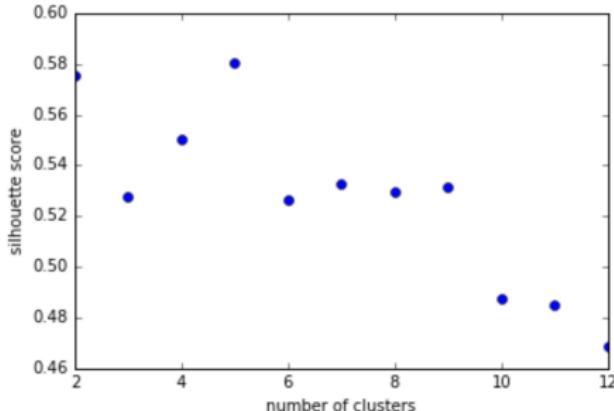
$\bar{\Delta}_i^k$: x_i 's dissimilarity in cluster k

$\bar{\Delta}_i^{k'}$: x_i 's dissimilarity to cluster k'
(x_i in, & closest to clusters k, k')

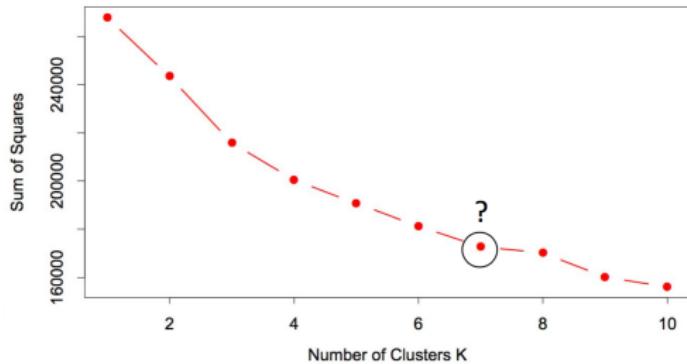
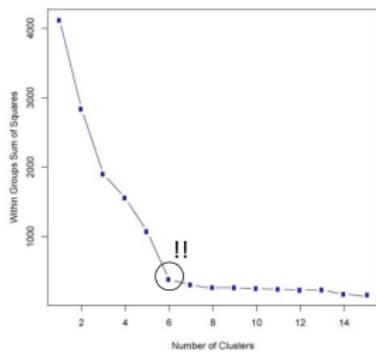
$$\text{Silhouette}(i) = \frac{\bar{\Delta}_i^{k'} - \bar{\Delta}_i^k}{\max(\bar{\Delta}_i^{k'}, \bar{\Delta}_i^k)}$$

Interpret silhouette scores?

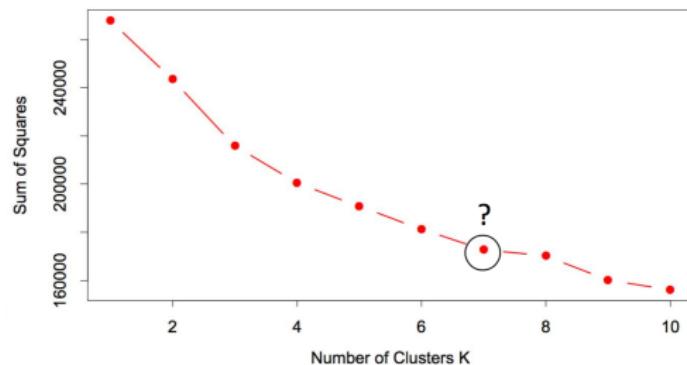
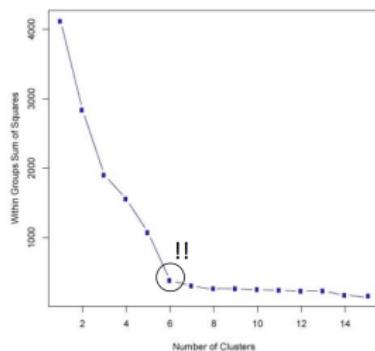
Compare average silhouette scores



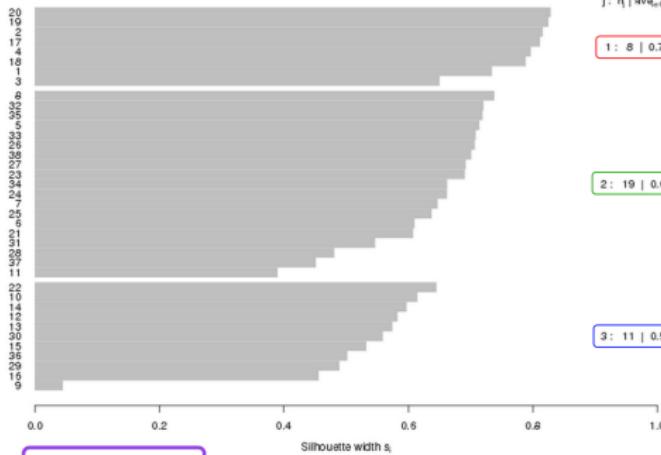
Elbow and Silhouette methods



Elbow and Silhouette methods



Silhouette plot of pam(x = cars.dist, k = 3)
n = 38



3 clusters C_1
j: 11 | avg $s_{C_1} s_j$

1: 19 | 0.64

3: 11 | 0.51

Guidelines for Overall Avg Silhouette

Range	Interpretation
0.71 – 1.0	Strong structure found
0.51 – 0.7	Reasonable structure
0.26 – 0.5	Structure weak/artificial
< 0.25	No substantial structure

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
                18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
                35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
                52, 53, 54, 55, 56, 57, 58, 59, 60])
```

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...

How could we test which class is doing better?

```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...
How could we test which class is doing better?

```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

- ▶ If the classes are doing equally well, the index doesn't matter

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...

How could we test which class is doing better?

```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

- ▶ If the classes are doing equally well, the index doesn't matter
↔ We could *permute* the index and see if it matters

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...

How could we test which class is doing better?

```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

- ▶ If the classes are doing equally well, the index doesn't matter

↔ We could *permute* the index and see if it matters

1. Repeatedly permute the index and recalculating the test statistic each time

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...

How could we test which class is doing better?

```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

- ▶ If the classes are doing equally well, the index doesn't matter

↔ We could *permute* the index and see if it matters

1. Repeatedly permute the index and recalculating the test statistic each time
2. These samples approximate the test statistic distribution under the null

The permutation test

- ▶ Students enrolling for a popular class get ids $i = 1, \dots, n$. As the class quickly fills, another section of the class is opened and students for that class are given ids $i = n + 1, \dots, 2n$.

```
id = np.array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
               18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34,
               35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
               52, 53, 54, 55, 56, 57, 58, 59, 60])
```

- ▶ The classes have the same tests...
How could we test which class is doing better?

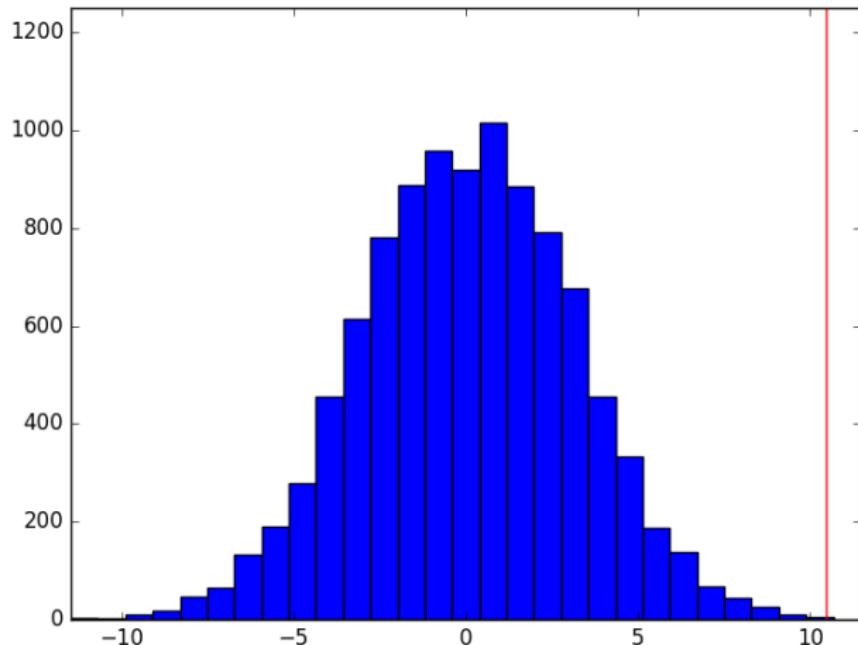
```
score = np.array([59, 92, 93, 83, 92, 61, 84, 92, 70, 93, 76, 70, 84, 61, 75, 91, 73,
                  67, 65, 64, 75, 80, 66, 56, 62, 53, 82, 69, 85, 94, 80, 86, 97, 99,
                  77, 84, 94, 70, 80, 87, 77, 85, 95, 80, 88, 90, 85, 84, 92, 75, 83,
                  89, 76, 95, 91, 97, 80, 70, 71, 94])
```

- ▶ If the classes are doing equally well, the index doesn't matter
↔ We could *permute* the index and see if it matters

1. Repeatedly permute the index and recalculating the test statistic each time
2. These samples approximate the test statistic distribution under the null
3. Compare the test statistic to this null distribution
to suggest how "strange" the actual observed test statistic is if the null is true

The *permutation test* (*This is my favorite test, btw*)

1. Permute the ids (i.e., believe the null is true: ids don't matter)
2. Recalculate the test statistic each time (under null)
3. See how strange your observed statistic is compared to nulls



The permutation test and the null distribution

- ▶ What was needed in the permutation test was really the distribution of the test statistic under the null hypothesis

The permutation test and the null distribution

- ▶ What was needed in the permutation test was really the distribution of the test statistic under the null hypothesis
- ▶ We used permutation to get it

The permutation test and the null distribution

- ▶ What was needed in the permutation test was really the distribution of the test statistic under the null hypothesis
- ▶ We used permutation to get it
but that was just a means to an end

The *permutation test* and the *null distribution*

- ▶ What was needed in the permutation test was really the distribution of the test statistic under the null hypothesis
- ▶ We used permutation to get it
but that was just a means to an end
- ▶ We can test against ANY null distribution we wish to propose

The Gap statistic

412

R. Tibshirani, G. Walther and T. Hastie

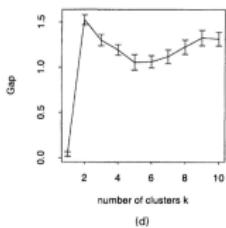
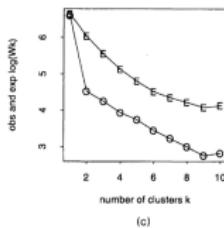
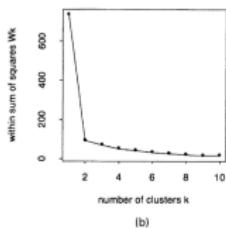
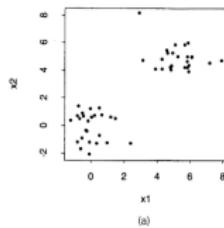


Fig. 1. Results for the two-cluster example: (a) data; (b) within sum of squares function W_k ; (c) functions $\log(W_k)$ (O) and $\hat{E}_n(\log(W_k))$ (E); (d) gap curve

416

R. Tibshirani, G. Walther and T. Hastie

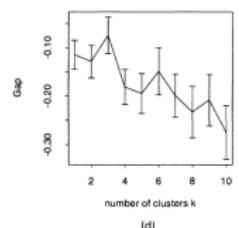
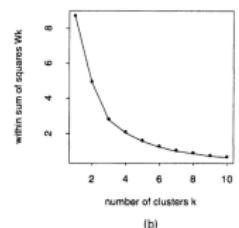
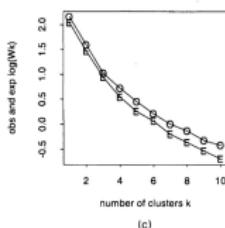
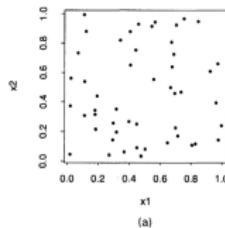
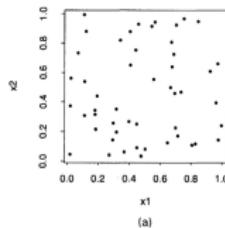


Fig. 2. Results for the uniform data example: (a) data; (b) within sum of squares function W_k ; (c) functions $\log(W_k)$ (O) and $\hat{E}_n(\log(W_k))$ (E); (d) gap curve

The Gap statistic

For M null distribution samples, calculate

$$\text{Gap}(K) = \bar{I} - \log W^{(*)}(C_K) \quad \text{and} \quad s_K = \sqrt{\frac{1}{R} \sum_{j=1}^M (\log W^{(r)}(C_K) - \bar{I})^2}$$

$$\text{where } W^{(r)}(C_K) = \sum_{C_K(i)=C_K(j)=k} \|x_i^{(r)} - x_j^{(r)}\|^2 \quad \text{and} \quad \bar{I} = \frac{1}{R} \sum_{r=1}^R \log W^{(r)}(C_K)$$

Then choose the smallest K such that $\text{Gap}(K) \geq \text{Gap}(K+1) - s_{K+1}$

412 R. Tibshirani, G. Walther and T. Hastie

416 R. Tibshirani, G. Walther and T. Hastie

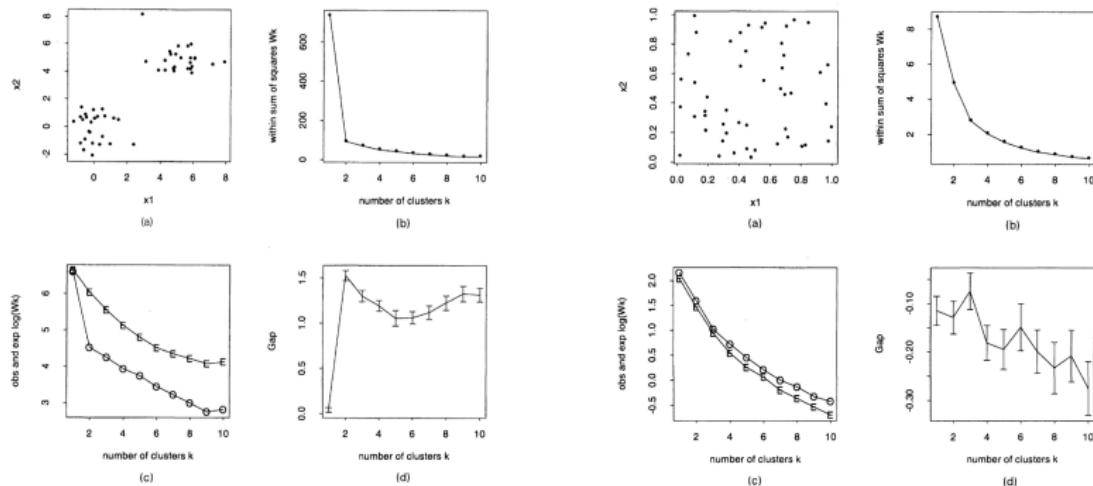


Fig. 1. Results for the two-cluster example: (a) data; (b) within sum of squares function W_k ; (c) functions $\log(W_k)$ (O) and $E_n(\log(W_k))$ (E); (d) gap curve

Fig. 2. Results for the uniform data example: (a) data; (b) within sum of squares function W_k ; (c) functions $\log(W_k)$ (O) and $E_n(\log(W_k))$ (E); (d) gap curve



The Gap statistic

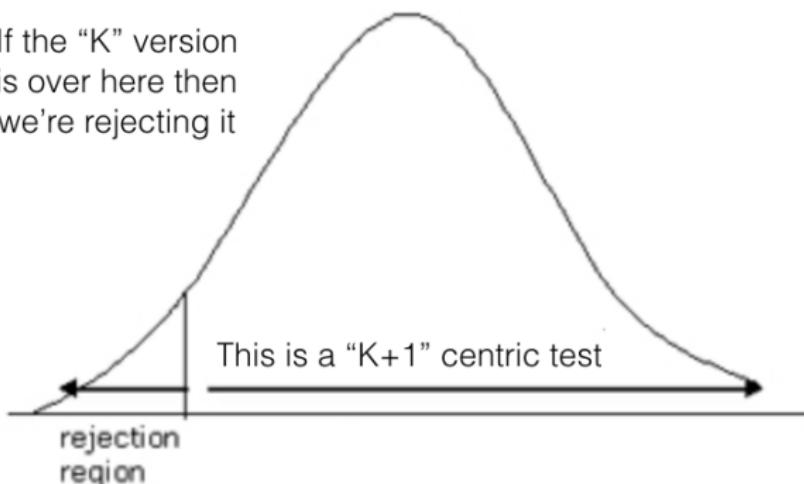
For M null distribution samples, calculate

$$\text{Gap}(K) = \bar{I} - \log W^{(*)}(C_K) \quad \text{and} \quad s_K = \sqrt{\frac{1}{R} \sum_{j=1}^M (\log W^{(r)}(C_K) - \bar{I})^2}$$

$$\text{where } W^{(r)}(C_K) = \sum_{C_K(i)=C_K(j)=k} \|x_i^{(r)} - x_j^{(r)}\|^2 \quad \text{and} \quad \bar{I} = \frac{1}{R} \sum_{r=1}^R \log W^{(r)}(C_K)$$

Then choose the smallest K such that $\text{Gap}(K) \geq \text{Gap}(K+1) - s_{K+1}$

If the “K” version
is over here then
we’re rejecting it



The Curse of Dimensionality

Just as nearest neighbors breaks down in high dimensional space...
Distance based clustering breaks down in high dimensional space...

The Curse of Dimensionality

Just as nearest neighbors breaks down in high dimensional space...

Distance based clustering breaks down in high dimensional space...

Go see Ryan's great slides motivating *The Curse of Dimensionality*

They're in the KNN slides

(THIS is the KMEANS lecture, NOT the KNN lecture)

The Curse of Dimensionality

Just as nearest neighbors breaks down in high dimensional space...

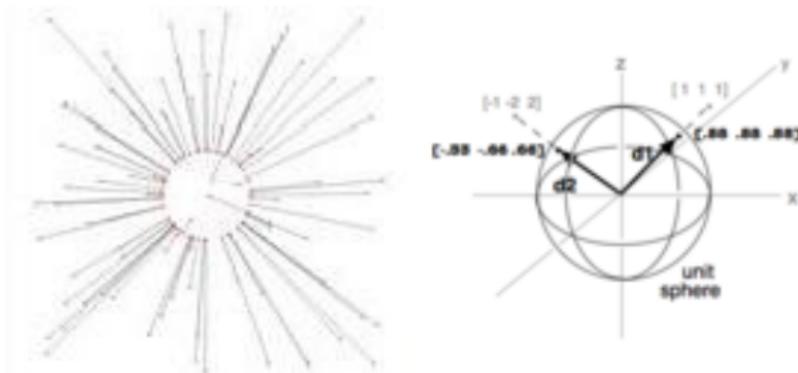
Distance based clustering breaks down in high dimensional space...

Go see Ryan's great slides motivating *The Curse of Dimensionality*

They're in the KNN slides

(THIS is the KMEANS lecture, NOT the KNN lecture)

However, when we normalize a vector (e.g., turn bag-of-words into probabilities, or make it a unit vector) we project all vectors onto a *manifold* where “locality” and “neighborhoods” make sense again.



The Curse of Dimensionality

So who is effected by *The Curse of Dimensionality* and *why (not)?*

Hint: *Overfitting* and
The Curse of Dimensionality are
slightly different things...

Challenge

A

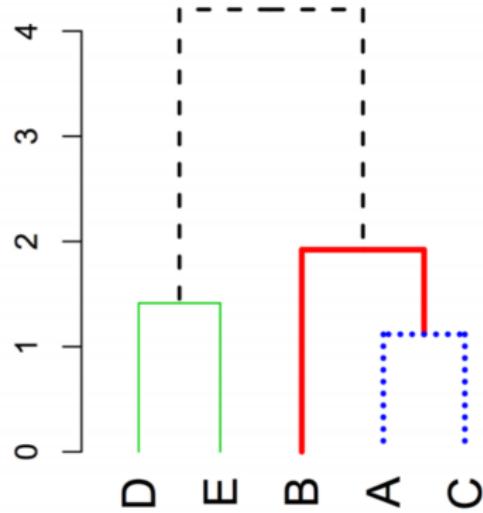
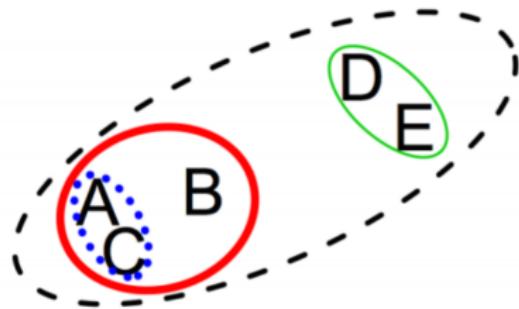
B

C

D

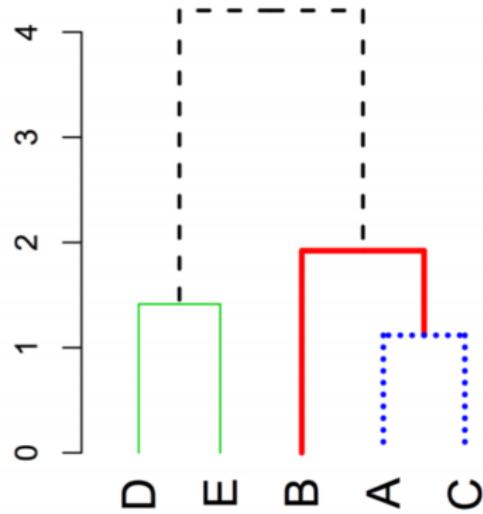
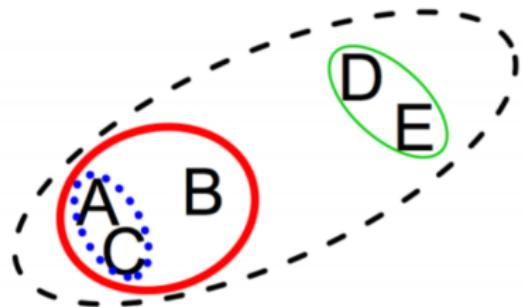
E

Hierarchical clustering



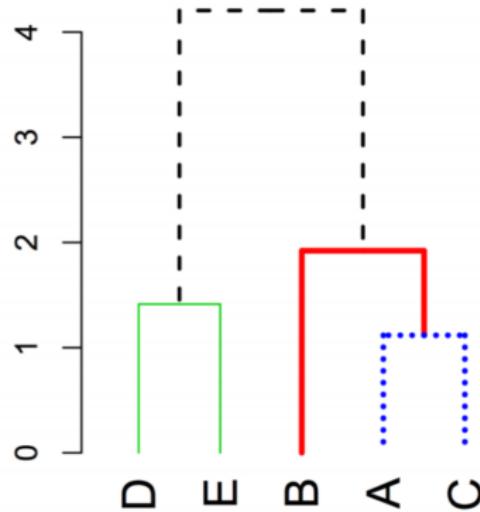
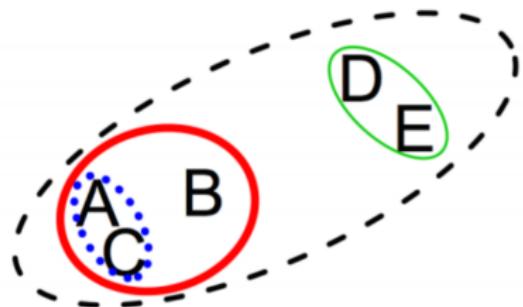
1. Assign each point to its own cluster

Hierarchical clustering



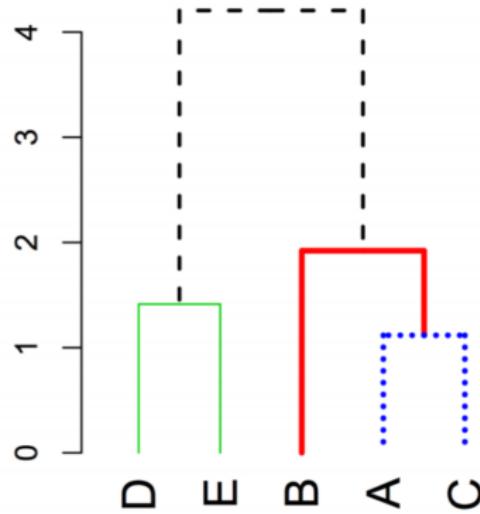
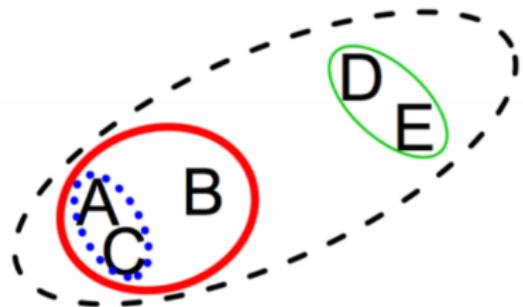
1. Assign each point to its own cluster
2. Computer pairwise cluster distances

Hierarchical clustering



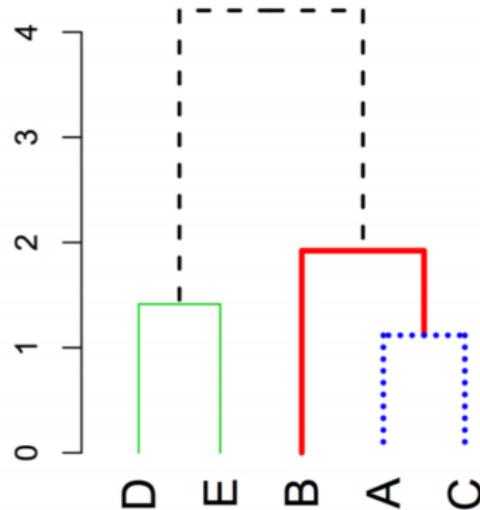
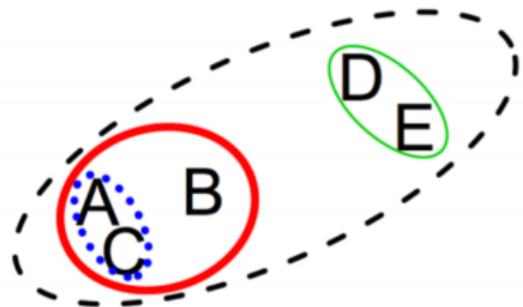
1. Assign each point to its own cluster
2. Computer pairwise cluster distances
3. Merge *closest two clusters*

Hierarchical clustering



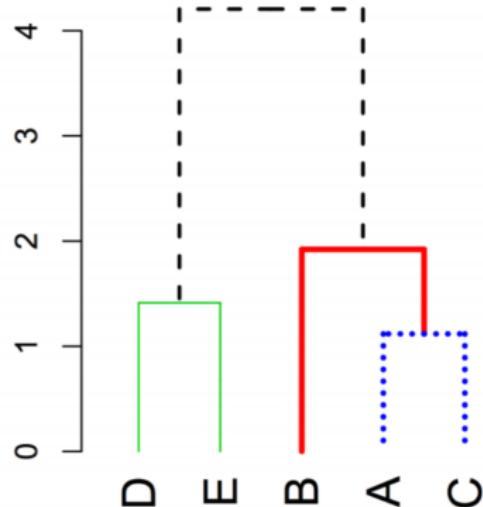
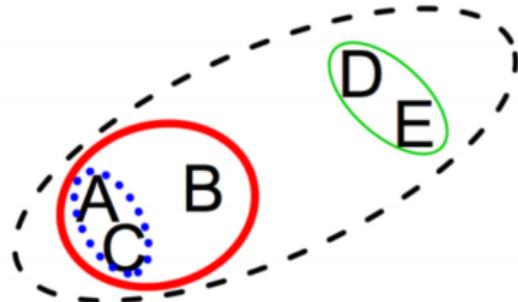
1. Assign each point to its own cluster
2. Computer pairwise cluster distances
3. Merge *closest two clusters*
4. Return to 2 until all clusters merged

Hierarchical clustering

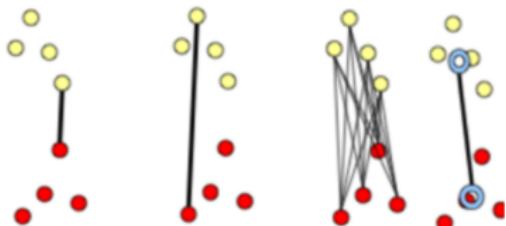


1. Assign each point to its own cluster
 2. Compute pairwise cluster distances
 3. Merge *closest two* clusters
 4. Return to 2 until all clusters merged
- ▶ Single: minimum pairwise point dissimilarity
 - ▶ Complete: maximum pairwise point dissimilarity
 - ▶ Average: average pairwise point dissimilarity
 - ▶ Centroid: centroid dissimilarity

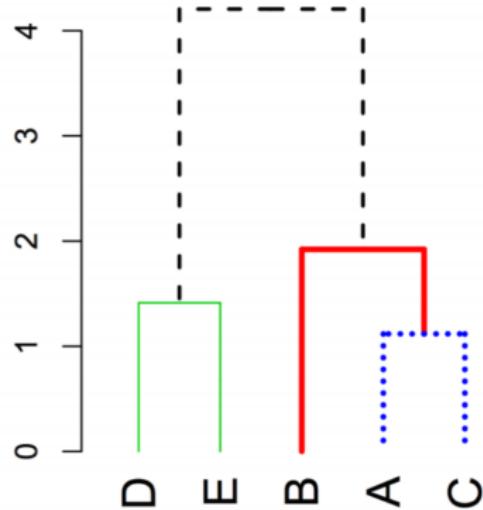
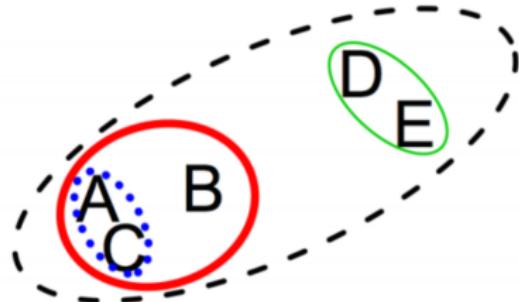
Hierarchical clustering



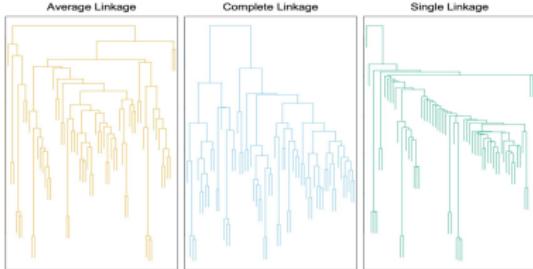
1. Assign each point to its own cluster
2. Compute pairwise cluster distances
3. Merge *closest two* clusters
4. Return to 2 until all clusters merged



Hierarchical clustering



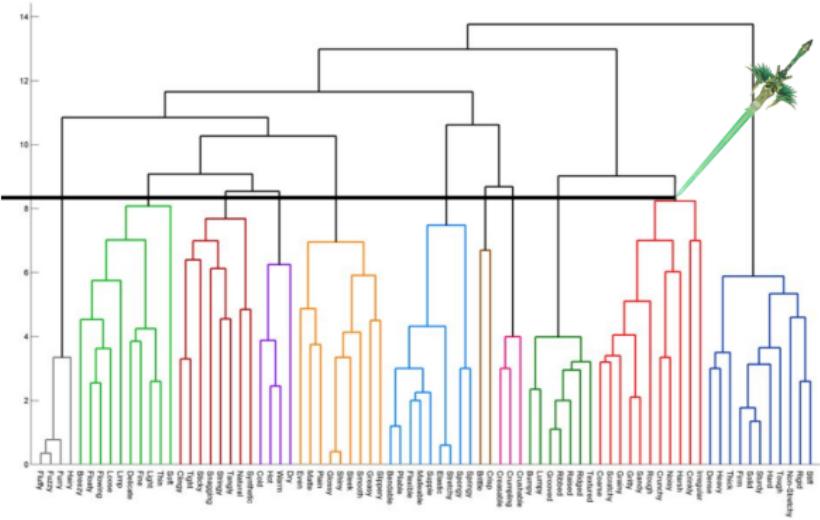
1. Assign each point to its own cluster
2. Compute pairwise cluster distances
3. Merge *closest two* clusters
4. Return to 2 until all clusters merged



Hierarchical clustering

Clusterings are created based on distance between proposed merges

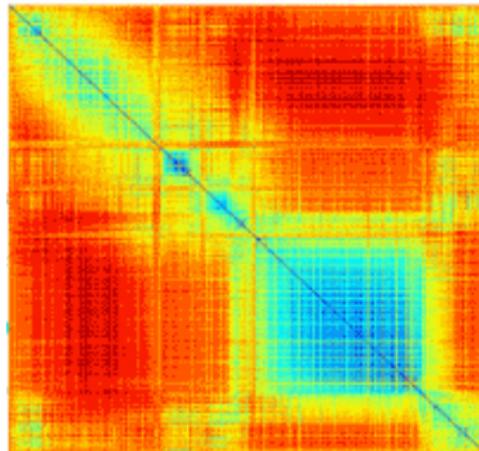
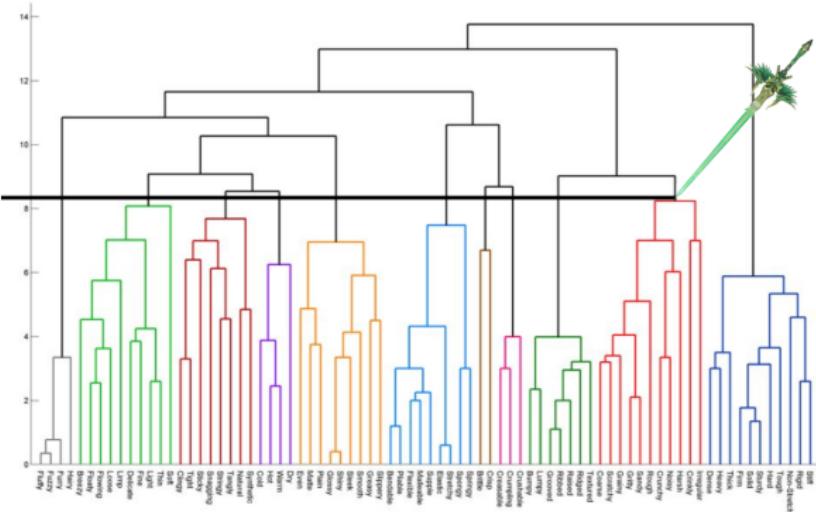
And then the number of clusters can be chosen *just like K-means...*



Hierarchical clustering

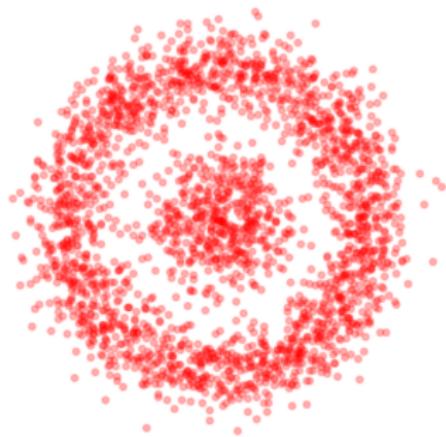
Clusterings are created based on distance between proposed merges

And then the number of clusters can be chosen *just like K-means...*

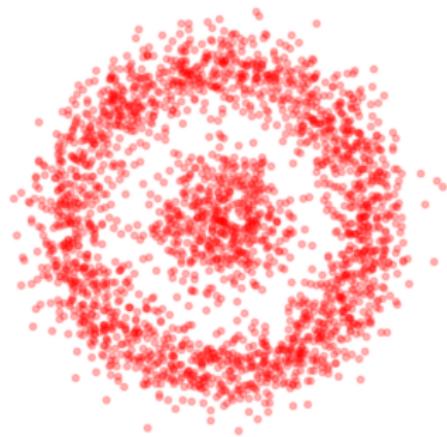


This is super cool, but
unlike K-means, hierarchical clustering requires all pairwise
comparisons so it doesn't scale gracefully with increasing data...

DBSCAN

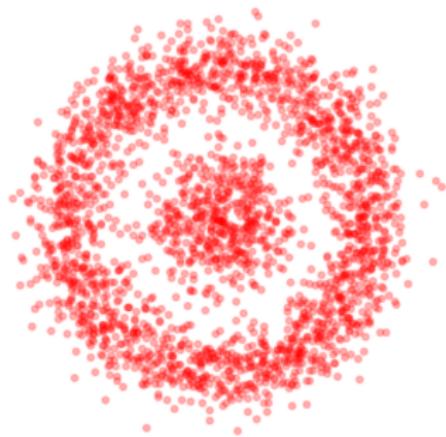


DBSCAN



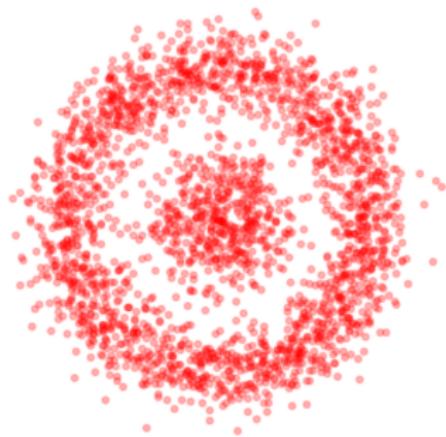
- ▶ Specify connection distance ϵ and minimum number of points for core m

DBSCAN



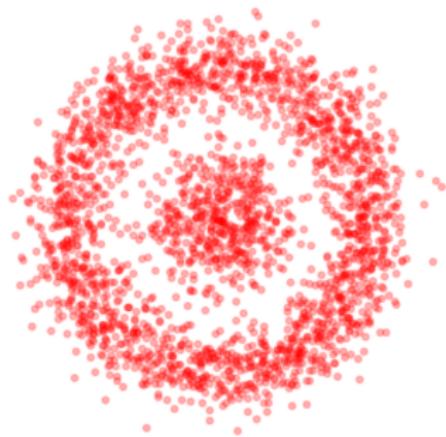
- ▶ Specify connection distance ϵ and minimum number of points for core m
- ▶ A cluster is all connected *core points*

DBSCAN

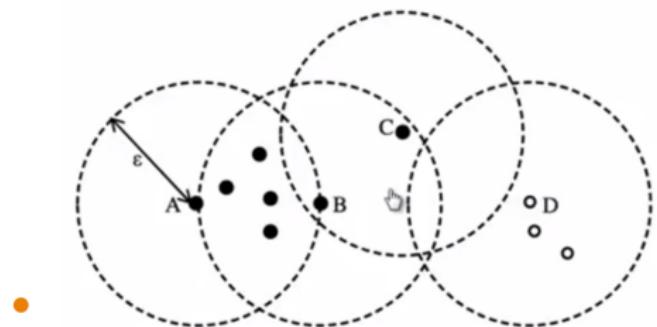


- ▶ Specify connection distance ϵ and minimum number of points for core m
- ▶ A cluster is all connected *core points*
- ▶ All other points are *noise*

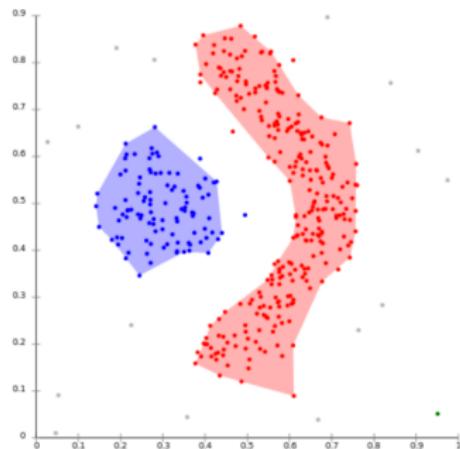
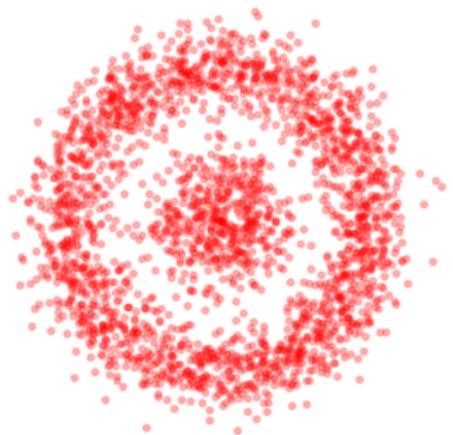
DBSCAN



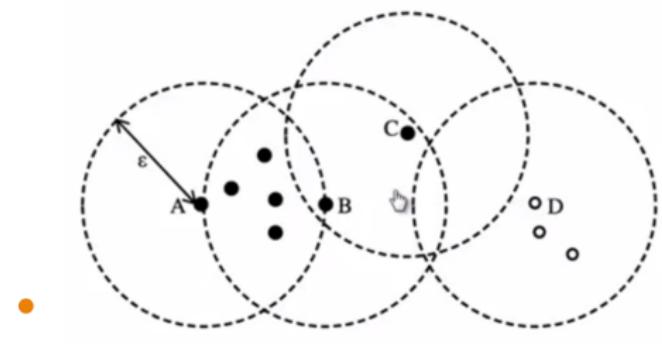
- ▶ Specify connection distance ϵ and minimum number of points for core m
- ▶ A cluster is all connected *core points*
- ▶ All other points are *noise*



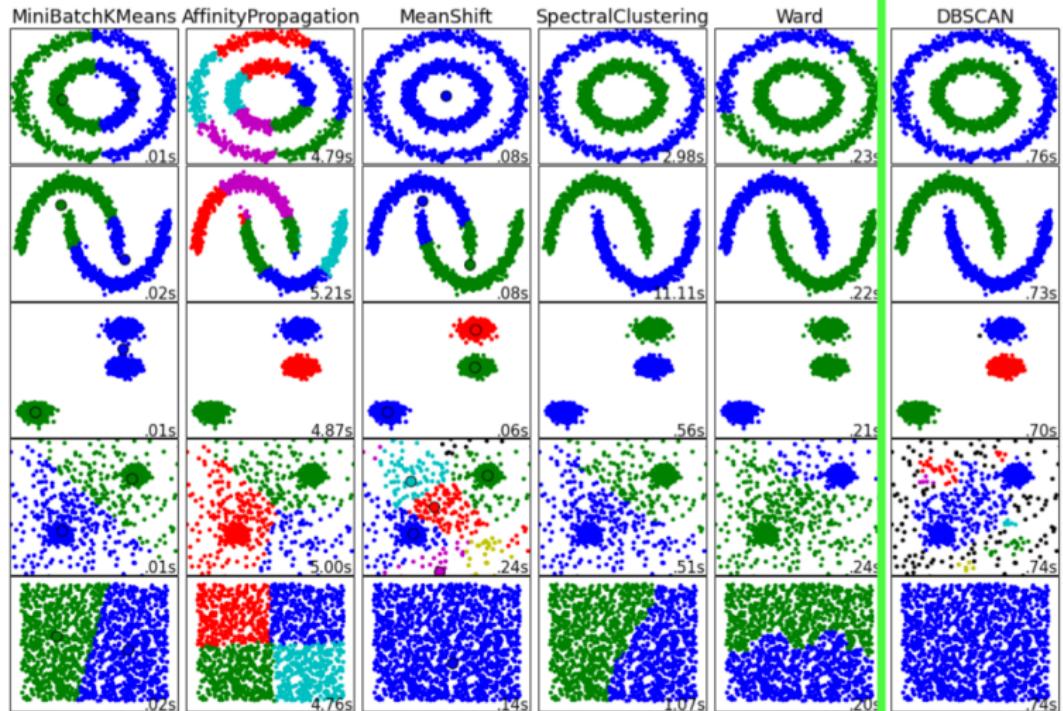
DBSCAN



- ▶ Specify connection distance ϵ and minimum number of points for core m
- ▶ A cluster is all connected *core points*
- ▶ All other points are *noise*

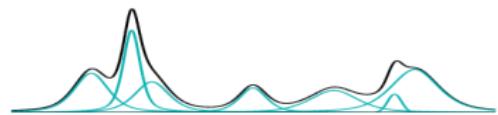


DBSCAN



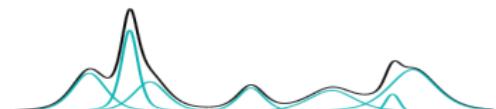
Bayesian Mixture Models (*xtra: my grad school jam*)

$$f(X_i|\mu, \sigma^2, \pi, \pi) = \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2)$$



Bayesian Mixture Models (xtra: my grad school jam)

$$f(X_i|\mu, \sigma^2, \pi, \omega) = \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2)$$



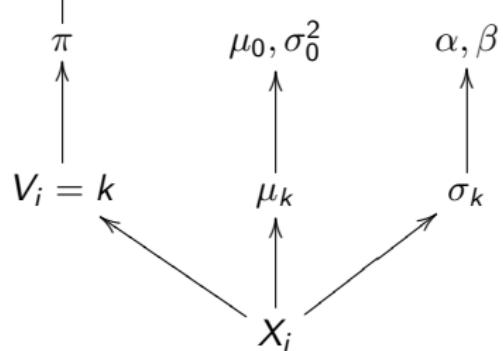
$$f(X_i|V_i, \mu, \sigma^2, \pi) = N(\mu_{V_i}, \sigma_k^2)$$

$$\Pr(V_i) = \text{Multinomial}(\pi, n=1)$$

$$f(\pi) = \text{Dirichlet}(\omega)$$

$$f(\mu_k) = N(\mu_0, \sigma_0^2)$$

$$f(\sigma_k^{-2}) = \text{Gamma}(\alpha, \beta)$$



Bayesian Mixture Models (xtra: my grad school jam)

$$f(X_i|\mu, \sigma^2, \pi, \omega) = \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2)$$



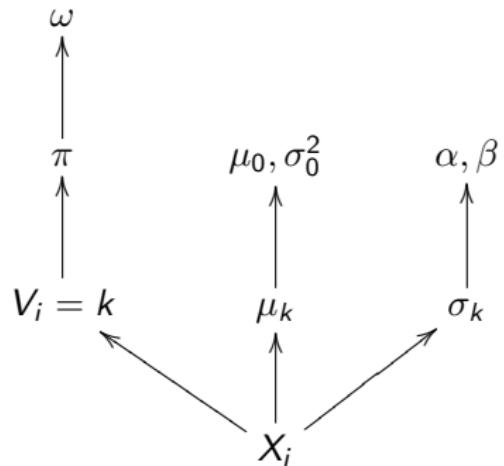
$$f(X_i|V_i, \mu, \sigma^2, \pi) = N(\mu_k, \sigma_k^2)$$

$$\Pr(V_i) = \text{Multinomial}(\pi, n=1)$$

$$f(\pi) = \text{Dirichlet}(\omega)$$

$$f(\mu_k) = N(\mu_0, \sigma_0^2)$$

$$f(\sigma_k^{-2}) = \text{Gamma}(\alpha, \beta)$$



$$f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \omega) =$$

$$\prod_{i=1}^n \left[\left(\sum_{k=1}^K \mathbb{1}_{[V_{ik}=1]} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2}\left(\frac{X_i - \mu_k}{\sigma_k}\right)^2} \right) \left(\prod_{k=1}^K \pi_k^{V_{ik}} \right) \right]$$

$$\times \left(\prod_{k=1}^K \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2}\left(\frac{\mu_k - \mu_0}{\sigma_0}\right)^2} \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{1}{\sigma^2}\right)^{\alpha-1} e^{-\beta \frac{1}{\sigma^2}} \right) \left(\frac{1}{\mathbf{B}(\omega)} \prod_{k=1}^K \pi_k^{\omega_k-1} \right)$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \mu, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \mu, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \mu, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \omega)$$

$$\propto f(\mathbf{X}, \mathbf{V}, \mu, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \omega) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi} | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega})$$

$$\propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) posterior sampling

A Gibbs sampler for the posterior

$$f(\mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Is made by cycling through the *full conditional distributions*

$$f(\mu_k | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$f(\sigma_k^2 | \mathbf{V}, \sigma^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \sigma^2, \pi, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

$$\Pr(\pi | \mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \\ \propto f(\mathbf{X}, \mathbf{V}, \boldsymbol{\mu}, \sigma^2, \pi | \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \quad (\text{which is proportional to the joint distribution})$$

Markov Chain Monte Carlo (MCMC) *full conditionals*

$$\Pr(V_{ik} = 1 | \mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\sigma}^2, \boldsymbol{\pi}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) \propto \pi_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2}\left(\frac{X_i - \mu_k}{\sigma_k}\right)^2}$$

$$f(\pi | \mathbf{V}, \mathbf{m}\boldsymbol{\mu}, \boldsymbol{\sigma}^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) = Dirichlet(\{\omega_k + n_k : k = 1, \dots, K\})$$

$$n_k = \sum_{V_{ik}=1} 1$$

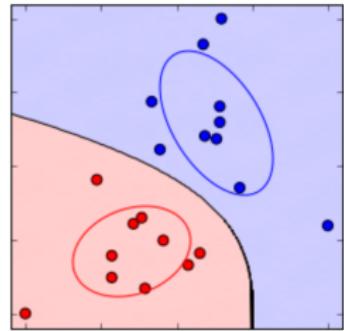
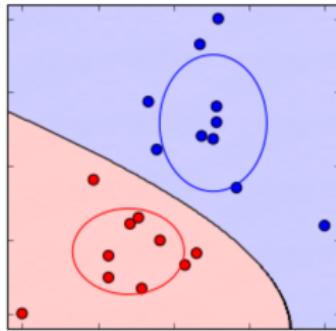
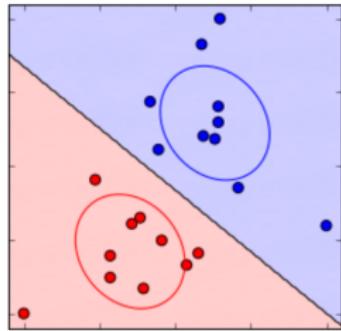
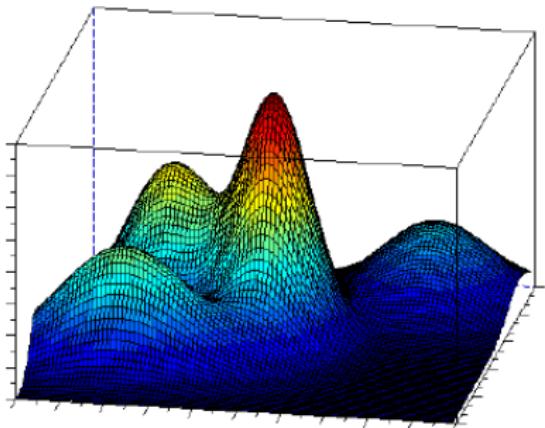
$$f(\sigma_k^2 | \mathbf{V}, \boldsymbol{\mu}, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) = Gamma\left(\frac{n_k}{2} + \alpha, \frac{1}{2} \sum_{V_{ik}=1} (X_i - \mu_k)^2 + \beta\right)$$

$$f(\mu_k | \mathbf{V}, \boldsymbol{\sigma}^2, \mathbf{X}, \mu_0, \sigma_0^2, \alpha, \beta, \boldsymbol{\omega}) =$$

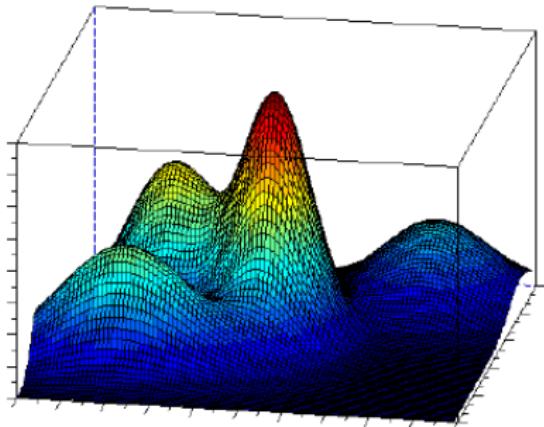
$$N\left(\hat{\sigma}_k^2 \left(\frac{\sum_{V_{ik}=1} X_i}{\sigma_k^2} + \frac{\mu_0}{\sigma_0^2}\right), \hat{\sigma}_k^2 = \left(\frac{n_k}{\sigma_k^2} + \frac{1}{\sigma_0^2}\right)^{-1}\right)$$

Mixture Models!

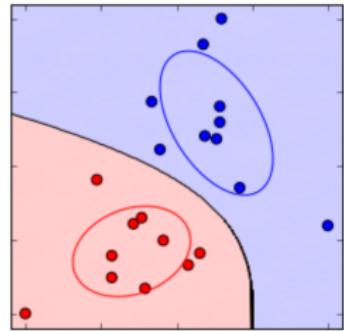
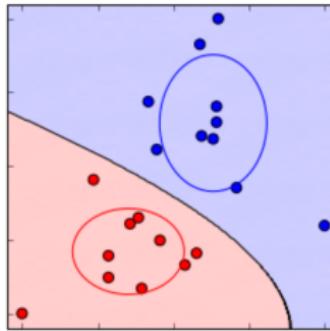
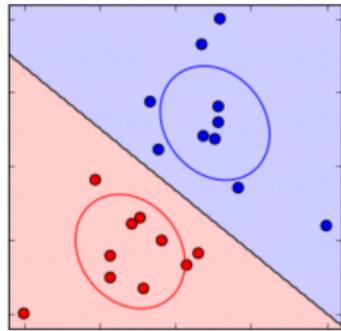
- ▶ Used to model “subpopulations”



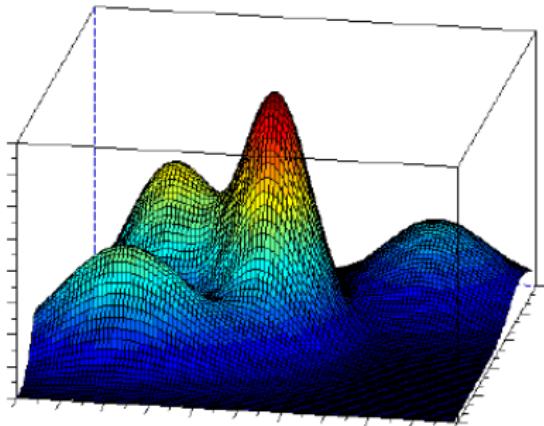
Mixture Models!



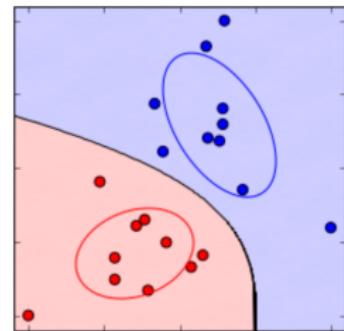
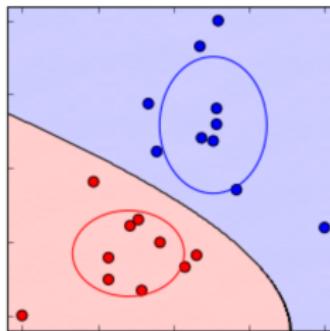
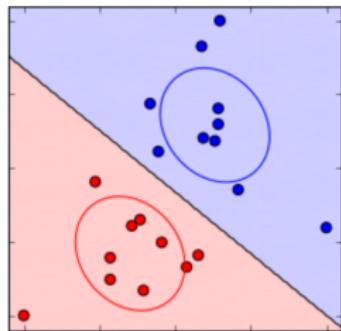
- ▶ Used to model “subpopulations”
- ▶ Or simply complex distributional shapes



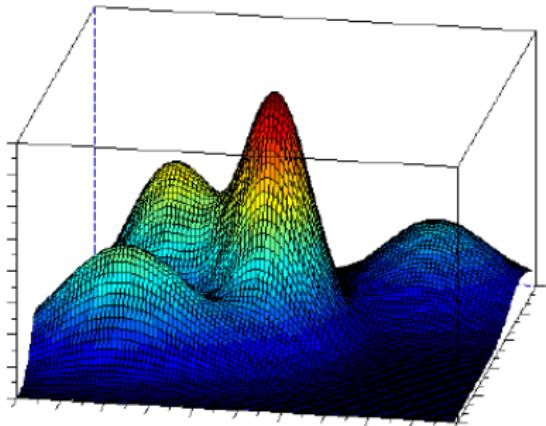
Mixture Models!



- ▶ Used to model “subpopulations”
- ▶ Or simply complex distributional shapes
- ▶ It's *almost nonparametric* like kernel density estimation



Mixture Models!



- ▶ Used to model “subpopulations”
- ▶ Or simply complex distributional shapes
- ▶ It's *almost nonparametric* like kernel density estimation
- ▶ Expectation-Maximization (EM) algorithm is another way to fit mixture models

