CSCI E-82

Advanced Machine Learning,
Data Mining & Artificial Intelligence
Lecture 6

Dynamic Time Warping Clustering

Peter V. Henstock Fall 2018

© 2018 Peter V. Henstock

Dynamic Time Warping

Problem with Distances

- What if we want to compare things that are different lengths?
- What if they are the same length but have intermediate states like phonemes within speech?
- What can we do?

© 2018 Peter V. Henstock

Dynamic Time Warping

- Stocks are easy to compare since start at same time and end at same time
- Speech is almost never aligned

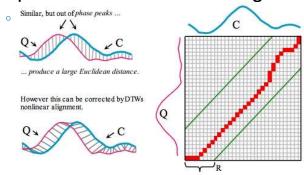
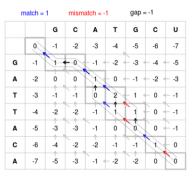


Figure 3: left) Two time series which are similar but out of phase, right) To align the sequences we construct a warping matrix, and search for the optimal warping path (red/solid squares). Note that Sakoe-Chiba Band with width R is used to constrain the warping path

http://4.bp.blogspot.com/-7QKyVmAJv21/UIxFDTlcjMI/AAAAAAAACew/AUbn2u-rS8A/s1600/dtw2.jpg

Dynamic Time Warping

- Stocks are easy to compare since start at same time and end at same time
- Bioinformatics pairwise sequence analysis is the same



https://en.wikipedia.org/wiki/Needleman%E2%80%93Wunsch_algorithm

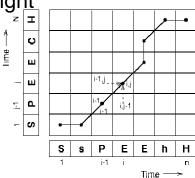
2018 Peter V. Hensto

Dynamic Time Warping

- Sequences x₁...x_n and y₁...y_m
- Construct warping matrix where i,j element is dist (x_i,y_j)
- Define cost: min {sqrt(sum w_k)}
 - w_k = matrix element of the warping path through the warping matrix
- Warping path
 - $z(i,j) = d(x_i,y_i) + min\{z(i-1,j-1), z(i-1,j), z(i,j-1)\}$

Dynamic Time Warping

- Construct matrix of distances
- Choose lowest cost path from bottom left to top right



$$z(i,j) = d(x_i,y_i) + min\{ z(i-1,j-1), z(i-1,j), z(i,j-1) \}$$

http://www.cnel.ufl.edu/~kkale/dtw.html

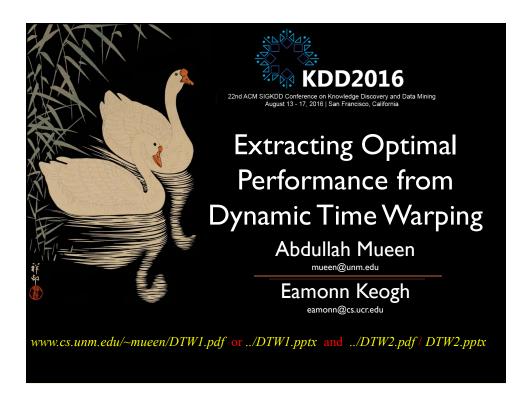
© 2018 Peter V. Henstock

What are the problem with this?

- 1) How big will the distance matrix be for sequences of 5000 and 4000?
- 2) Assumed that they start and end at the same place. What if they don't?

Dynamic Time Warping + KNN

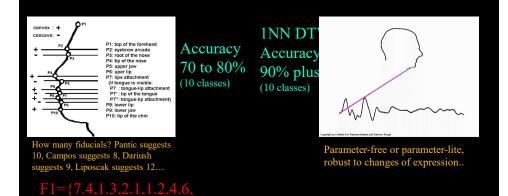
- Nearest Neighbor DTW is very hard to beat across fields
- When it loses, it barely loses but requires extreme complexity/coding
- DTW arises in many different problems

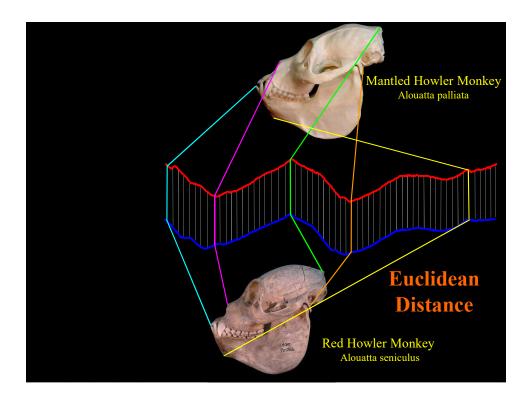


An Unstated Assumption (alterative)

For all time series, heartbeats, gait, fish and faces, we could design and extract features, and just represent the objects as feature vectors. However, in the dozen of cases where we compared this type of approach to just using DTW on the raw data, DTW almost always wins!

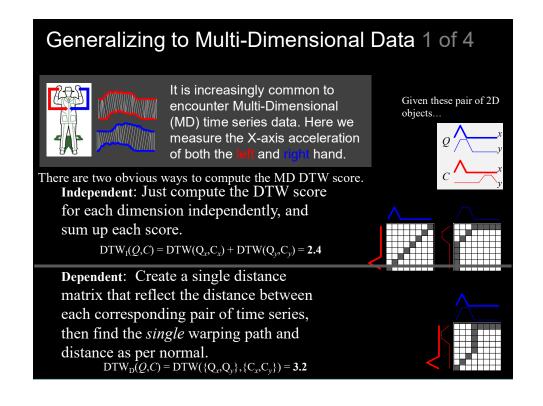
So, at the very least, try the simplest ideas first, and DTW is very simple.





An Unstated Assumption (alterative) For all time series, heartbeats, gait, faces, and even fish!, we could design and extract features, and just represent the objects as feature vectors. However, in the dozen of cases where we compared this type of approach to just using DTW on the raw data, DTW almost always wins! So, at the very least, try the simplest ideas first, and DTW is very simple. Accuracy 75.7% Accuracy 86.0%

rsal fin and adipose fin



Generalizing to Multi-Dimensional Data 4 of 4

Critical Point: You *can* generalize DTW to 2,3,4,...1,000 dimensions.

However, it is *very* unlikely that more than 2 to 4 is useful, after that, you are almost certainly condemned to the curse of dimensionality.

PAMAP Physical Activity

Consider a physical activity dataset containing 36 axis synchronous measurements from three Inertial Measurement Units (IMUs) located on the wrist, chest and ankle. This dataset has eight subjects performing activities such as: rope-jumping, running, folding laundry, ascending-stairs.

Using *all* dimensions is a disaster! Using the best three is a lot better than using the best one, so there is evidence that Multi-Dimensional DTW really does help.

PAMAP, Physical Activity Monitoring for Aging People www.pamap.org/ (we used DTW with w = 0 for simplicity here, we could do better by tuning w)



Which Dimensions?	Accuracy
All dimensions	0.19
A single random dimension (on average)	0.51
The best single dimension (as predicted by CV)	0.72
The best 3 dimensions (as predicted by CV)	0.89

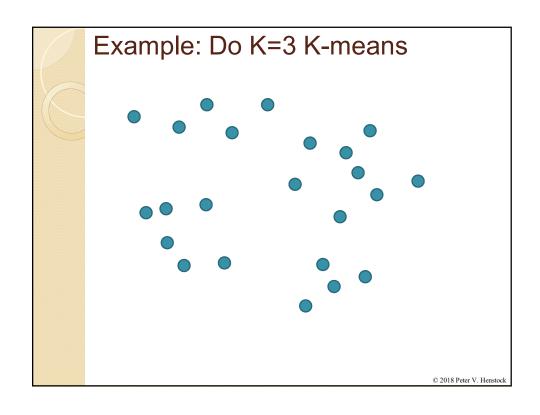
Clustering

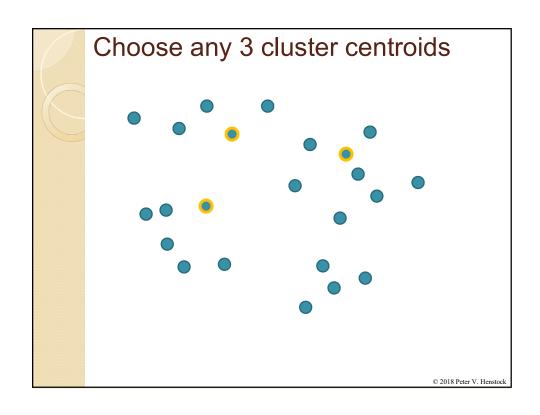
Clustering Approaches

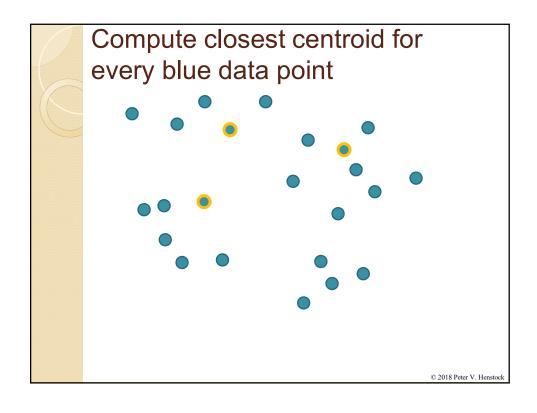
- Partitioning: non-overlapping subsets
- Hierarchical
- Density-based: DBSCAN
- Grid-based: CLIQUE
- (Correlation)
- Spectral
- (Gravitational)
- (Hard)

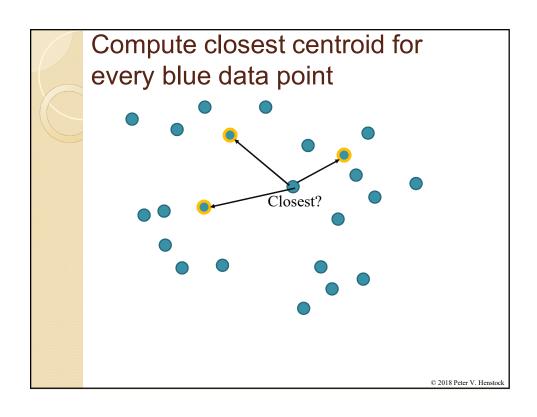
© 2018 Peter V. Henstock

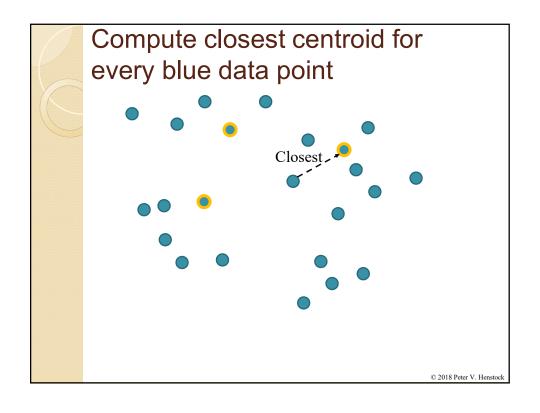
K-Means

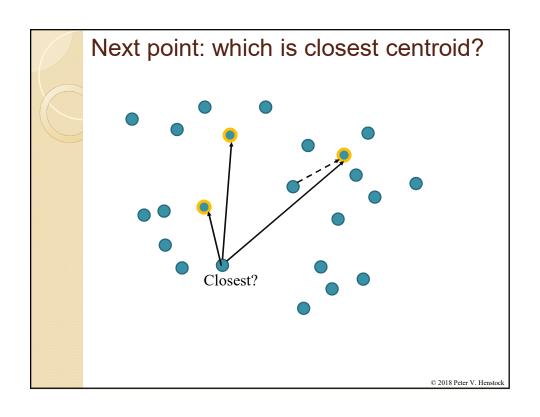


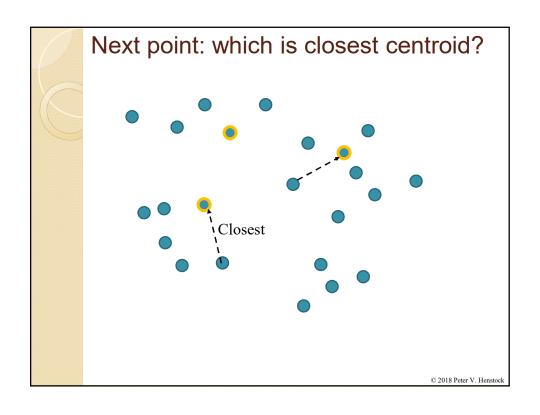


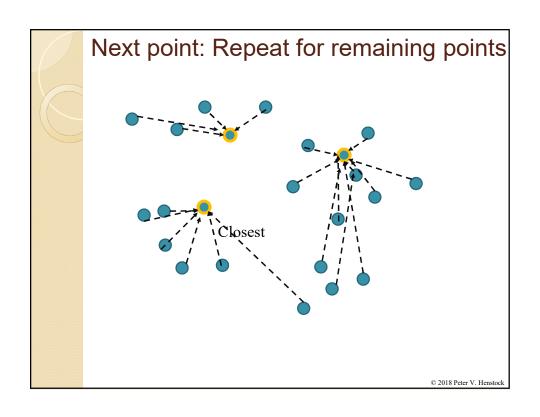


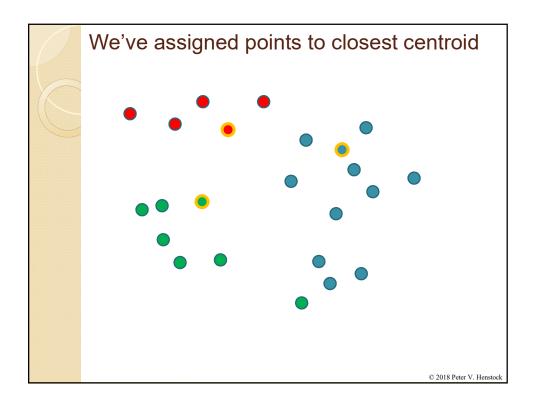


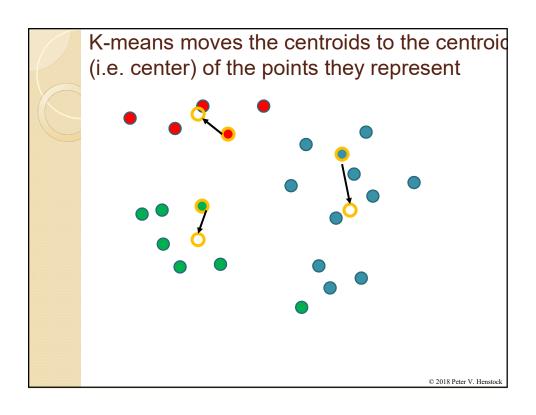


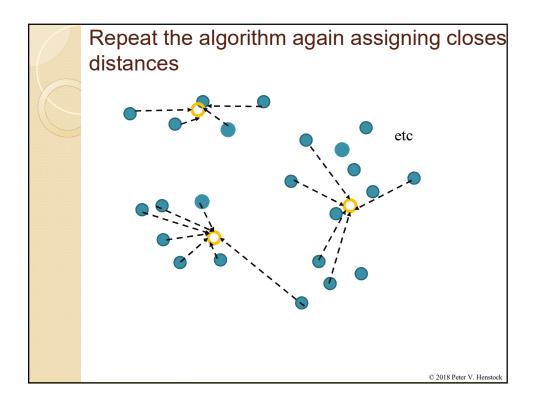


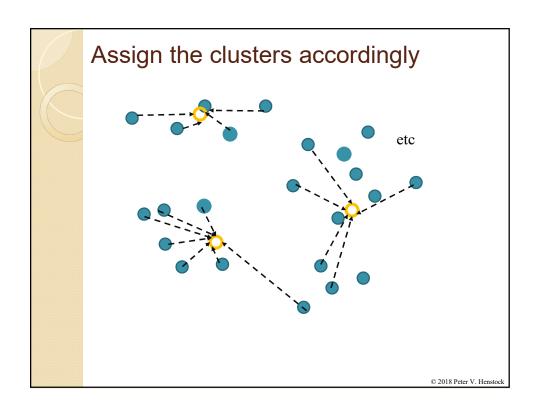


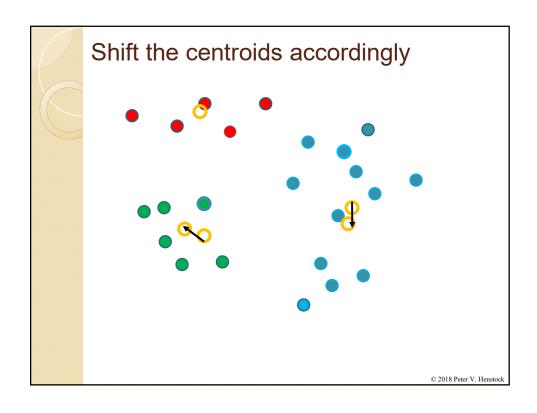


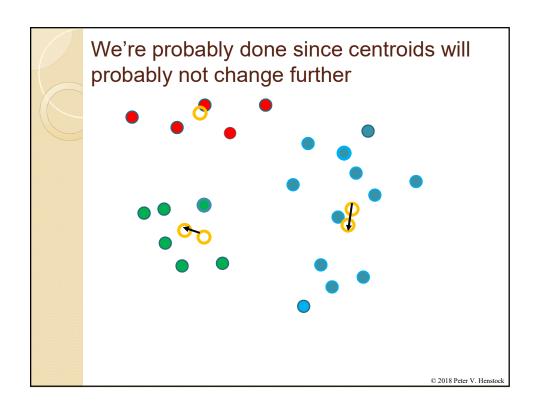


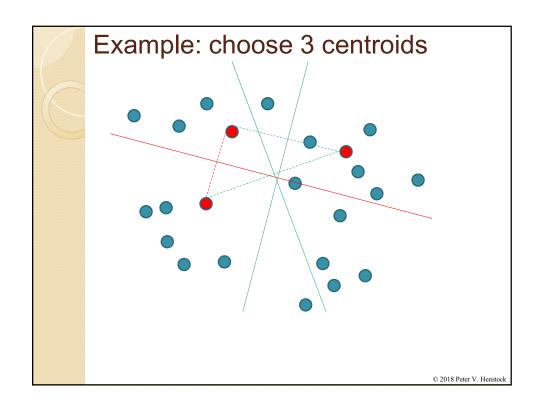


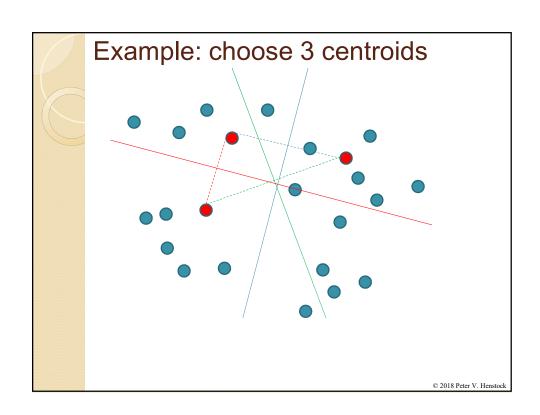


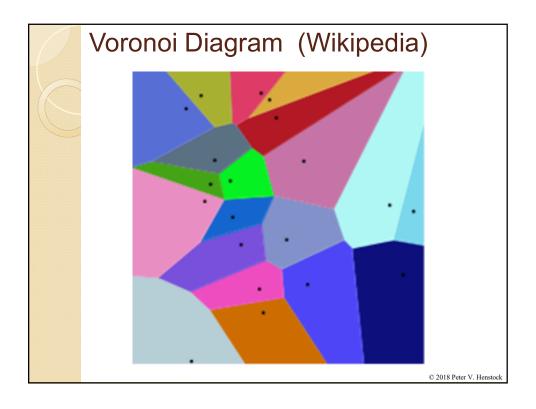












Optimization

- Machine learning is all about optimization so...
- What is the optimization criteria for K-Means?

K-Means Clustering

- Data vector x₁..x_N (with M features)
- Clusters $\{\mu_1 \dots \mu_K\}$
- Which x map to which cluster?
 - ∘ {c₁...c_N} is the cluster label for each x₁..x_N
 - · Encode vector where 3 indicates cluster 3
 - Alternative encoding for cluster labels is "one hot coding" where
 - K vectors with 1 only when it applies
 - ∘ r_i = responsibility vector
 - r_{ij} = 1 if x_i maps to cluster j else 0

© 2018 Peter V. Henstock

Coding Representation







• C coding: [??????]

Coding Representation







- C coding: [1 3 2 1 2 1]
- One hot coding with r
 - r1: [?]
 - r2: [?]
 - r3: [?]

© 2018 Peter V. Henstock

Coding Representation







- C coding: [1 3 2 1 2 1]'
- One hot coding with r
 - r1: [1 0 0 1 0 1]'
 - o r2: [0 0 1 0 1 0]'
 - ° r3: [0 1 0 0 0 0]'

Cost (or Loss) Function

- $J_{xi}(r, \mu) = \Sigma_k r_{ik} [x_i \mu_k]^2$
- This describes how far away point x_i is from the particular cluster to which is assigned
 - How many non-zero terms are in the sum?

© 2018 Peter V. Henstoc

Cost (or Loss) Function

- $J_{xi}(r, \mu) = \Sigma_k r_{ik} [x_i \mu_k]^2$
- This describes how far away point x_i is from the particular cluster to which is assigned
 - How many non-zero terms are in the sum?



Cost (or Loss) Function

- $J_{xi}(r, \mu) = \Sigma_k r_{ik} [x_i \mu_k]^2$
- This describes how far away point x_i is from the particular cluster to which is assigned
 - How many non-zero terms are in the sum?



- J (all r, μ 1.. $\mu\kappa$) = $\Sigma_i \Sigma_k r_{ik} [x_i \mu_k]^2$
 - Perform this for all points xi

© 2018 Peter V. Henstock

Optimization

- Similar case to regression where we set up a cost function
 - Solved it using a gradient search
 - Solved it analytically
- J (all r, μ 1.. μ κ) = $\Sigma_i \Sigma_k r_{ik} [x_i \mu_k]^2$
- Formally called the "K-Means Problem"
- Only problem is that it is NP-hard

Lloyd's Algorithm

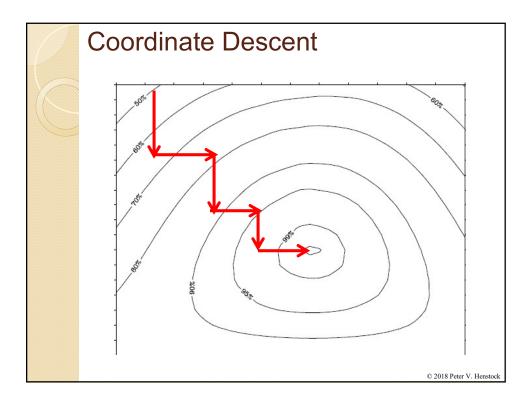
- Computationally feasible solution to solve the K-Means problem
- Coordinate descent

© 2018 Peter V. Henstock

Coordinate Descent

- Gradient descent took the steepest descent across multiple variables at one time using the slope
- Coordinate descent optimizes one variable at a time





Lloyd's Algorithm

- Computationally feasible solution to solve the K-Means problem
- Coordinate descent
- Optimize the r's and the μ 's

Two Steps of Lloyd's Algorithm

- 1) For each x_i , set r_{ik} = argmin_k $[x_i \mu_k]^2$
 - Cluster assignment step
- 2) Optimize the centroids μ
 - $_{\circ}$ Compute Gradient(J_{xi}(r,\,\mu))_{\mu k}\, and set to 0
 - = Gradient $\Sigma_i r_{ik} [x_i \mu_k]^2$
 - \circ -2 $\Sigma_i r_{ik} [x_i \mu_k] = 0$
 - $\circ \ \Sigma_{i} \ r_{ik} \ x_{i} \ = \Sigma_{k} \ r_{ik} \ \mu_{k}$
 - $\circ \Sigma_i r_{ik} x_i = N_k \mu_k$ where $N_k = \#pts in k$
 - $\mu_k = \sum_i r_{ik} x_i / N_k$
 - $_{\circ}$ μ_{k} = average of x's assigned to cluster k

© 2018 Peter V. Henstock

K-Means Clustering

- Requires the number of clusters k
- Takes a set of input vectors x₁..x_N
- Take k unique arbitrary input vectors as initial centroids μ₁..μ_k
- Loop {
 - Assign each input to the closest centroid
 - Shift centroid to center of points it represents

Critical aspects of K-Means?

- Choice of K
- Initial conditions
- Distance function (usually Euclidian)

© 2018 Peter V. Henstoc

Computation & Initialization

• What is the big O() for K-Means?

Computation & Initialization

- What is the big O() for K-Means?
 - O(#loops [kN dist calc + N re-centering])
 - O(#loops k N #features)
- Is it a fast method or a slow method?

© 2018 Peter V. Henstock

Computation & Initialization

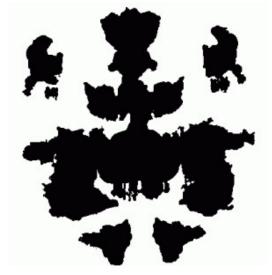
- What is the big O() for K-Means?
 - O(#loops [kN dist calc + N recentering])
 - Assume k << N
 - O(loop kN)
- Is it a fast method or a slow method?

Computation & Initialization

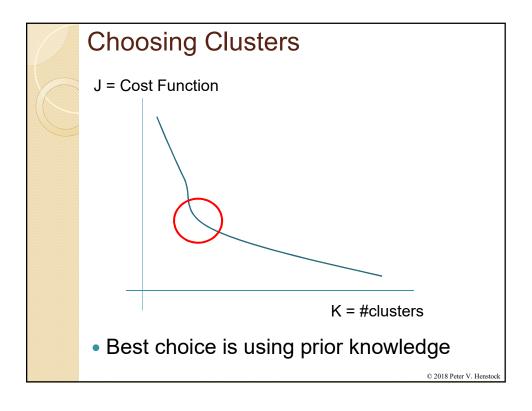
- What is the big O() for K-Means?
 - O(#loops [kN dist calc + N recentering])
 - Assume k << N</p>
 - O(loop kN)
- Is it a fast method or a slow method?
 - Pretty fast
 - Can consider multiple initializations

© 2018 Peter V. Henstock

How many clusters?



- Easy case when only 2D
- http://www.obeyclothing.com/blog/?p=3901



Issues With K-Means **Entry x1 x2 x3 x4 x5 x6** 6.6 data1 2.3 1.5 3.5 4.9 9.4 data2 2 1 3 4 7 9 1 7 3 5 6 data3 10 2 4 NA data4 © 2018 Peter V. Henstock

Issues With K-Means

Name	Height (feet)	Weight (lbs)	GPA
joe	5.7	160	3.5
fred	6.1	180	3.9
david	5.9	192	3.0
mike	5.95	210	4.0

© 2018 Peter V. Henstock

Issues With K-Means

Name	Race	Major	grade
joe	white	IT	Α
fred	asian	IT	Α
david	hispanic	Econ	В
mike	hispanic	ΙΤ	Α

K-Medoid Variation

- Medoid = represents object of a data set
- K-Means
 - Mean value of each parameter within cluster
- K-Mode
 - Mode value of each parameter
- K-Medoid
 - Best representative of a cluster that is also a data point

© 2018 Peter V. Henstock

K-Means

Good

Bad

Ugly

K-Means

- Good
 - Fast
 - Easy to interpret
- Bad
 - Need to predict the structure with K
 - Not optimal
- Ugly
 - Strongly depends on initial K
 - Varies with distance function approach
 - Creates spherical/circular clusters

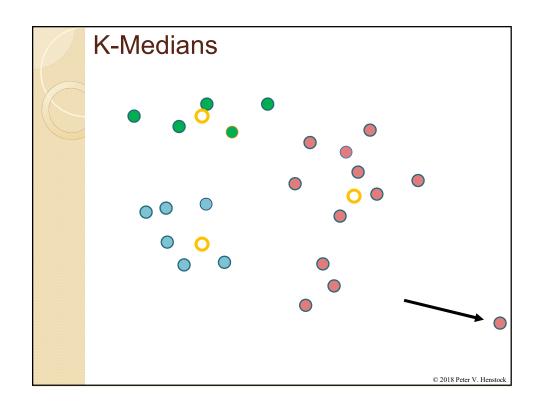
© 2018 Peter V. Henstock

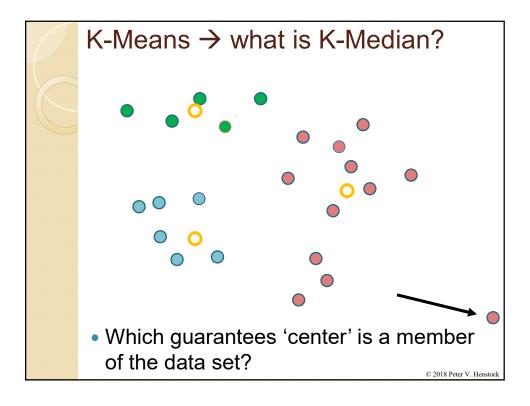
K-Means++ is which of the following?

- a) K-means algorithm implemented in the C++ language for speed
- b) K-means++ > K-means+ > K-means-but you only use the latest: K-means+++
- K-means for large data that increments values at each iteration
- d) K-means with a probabilistic spacing of the initial centroids
- e) K-means with improved cost function so it will reach the global minimum

K-Means++

- Initialization is different
 - Randomly choose one point as first centroid
 - Repeat until K centroids chosen
 - · Compute D(x) for all points
 - $D(x_i)$ = distance of x_i to closest centroid so far
 - Choose new centroid with probability $\frac{D(x)^2}{\sum D(x)^2}$
- Use K-Means algorithm for the rest
- What points does "D² weighting" favor?





K-Medoid Variation

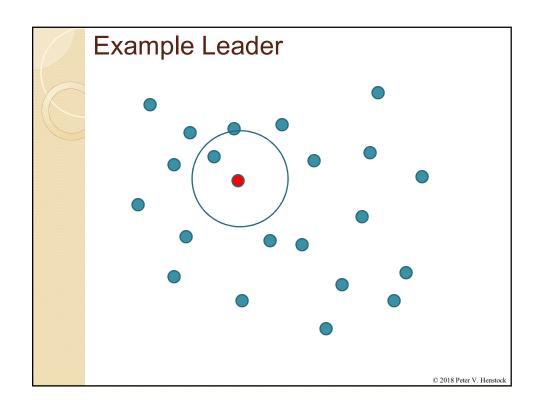
- Medoid
 - = representative object of a data set
 - = member of data set closest to middle
- K-Means
 - Mean value of each parameter within a cluster
- K-Mode
 - Mode value of each parameter
- K-Medoid
 - Best representative of a cluster that is also a data point

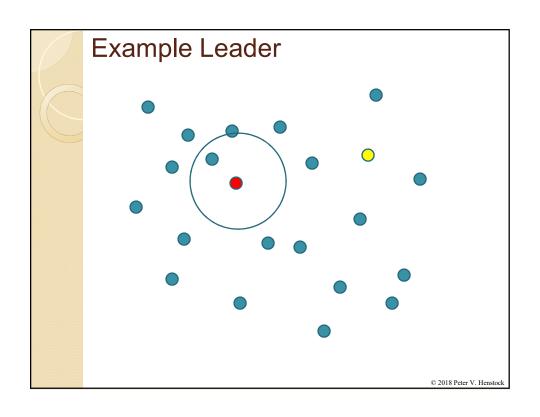
Leader Clustering

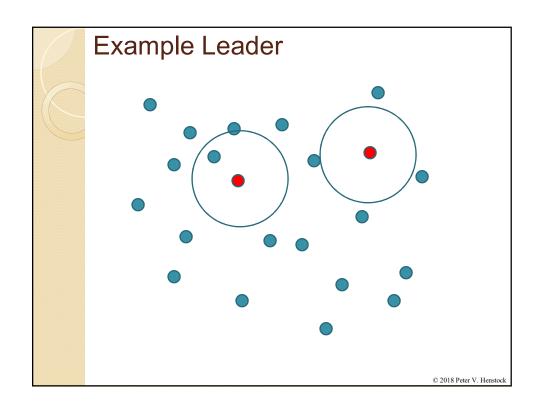
© 2018 Peter V. Henstock

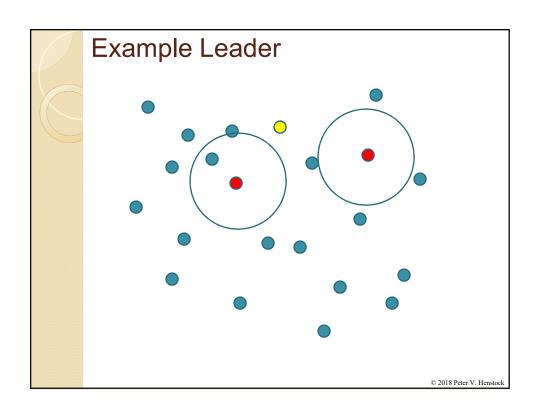
Leader Clustering

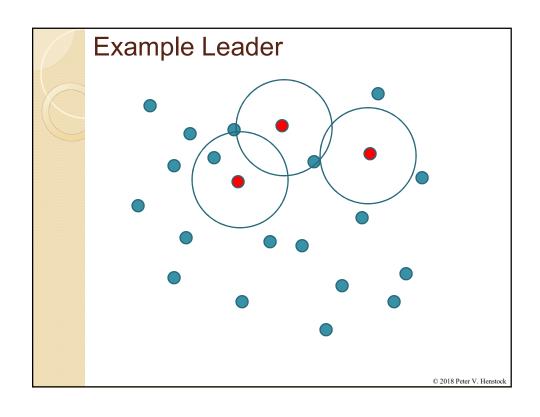
- ClusterCenters = {}
- For each point x_i in randomly ordered data set
 - closestDist = dist(closestCluster, x_i)
 - If closestDist < threshold
 - Assign x_i as a member of closestCluster
 - Else {
 - Add x_i to ClusterCenters

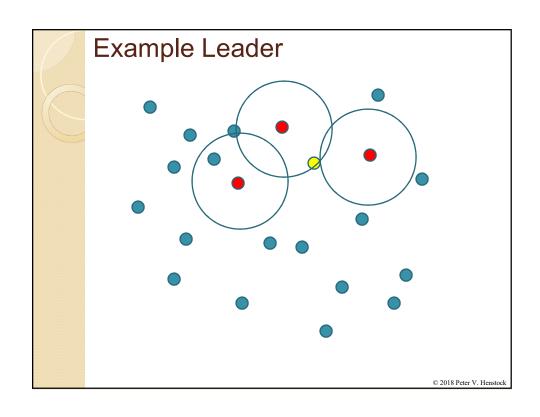


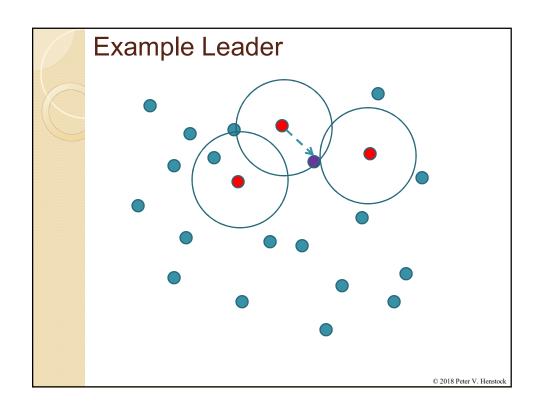


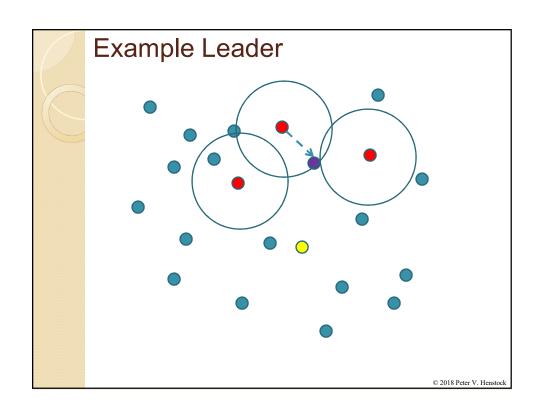


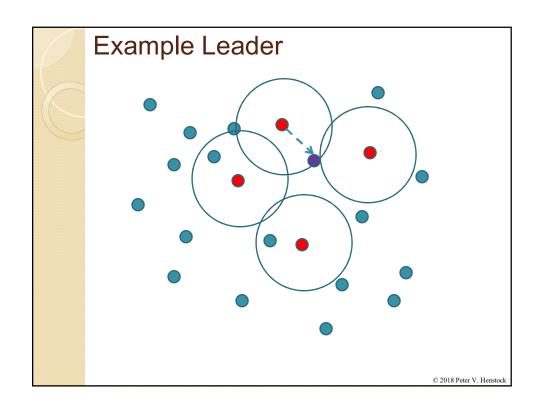


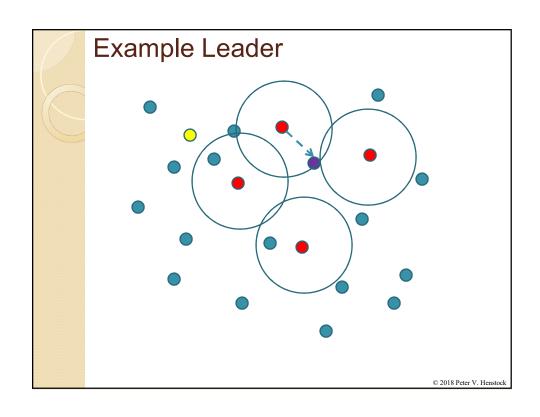


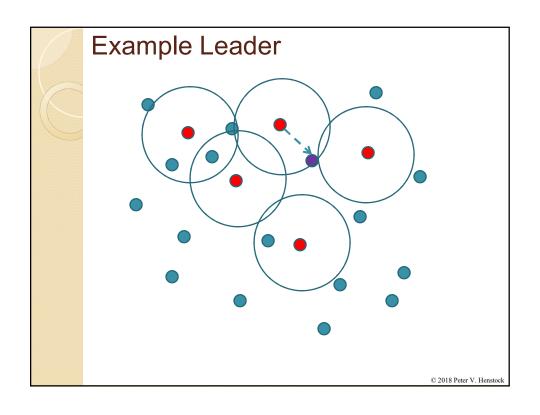


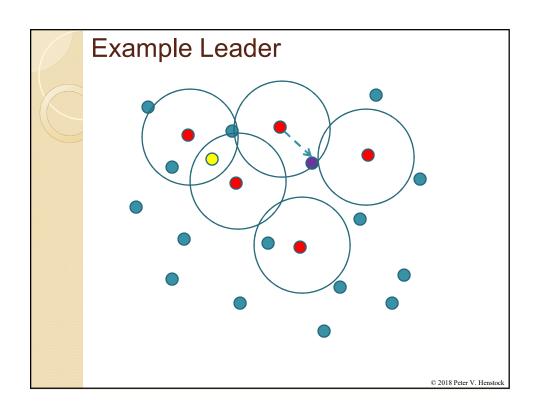












K-Means & Leader

- What do they have in common?
- Which is better?
- Which is faster?

© 2018 Peter V. Henstock

Space & Time

• How much time?

• How much memory?

OptiSim

© 2018 Peter V. Henstock

How to get a diverse set

- K-means works poorly if the initial values are too close together
- How would you go about approaching the problem?

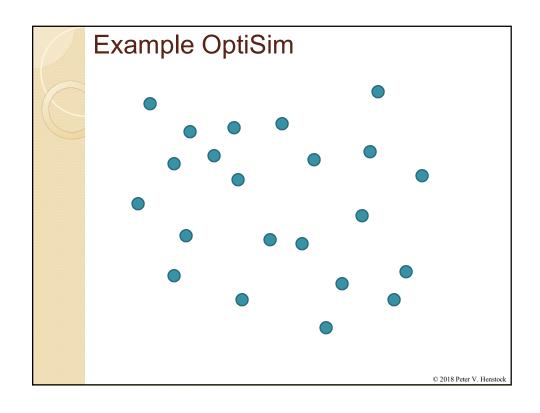
OptiSim

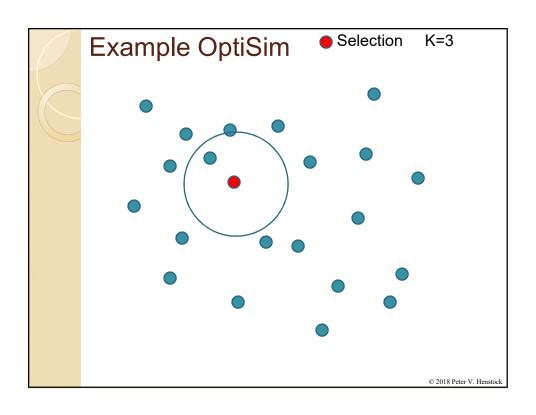
- Clark 1997, Clark & Langton 1998
- Randomly order data points
- Take 1st data point → {selection}
- SampleSet ← {}
- Repeat until obtained "enough" samples
 - dataPoint ← next data point
 - If dist(dataPoint, sample in selection) > Thr
 - Add dataPoint to SampleSet
 - · if |SampleSet| == K
 - · Add most dissimilar point to {selection}
 - · SampleSet ← {}

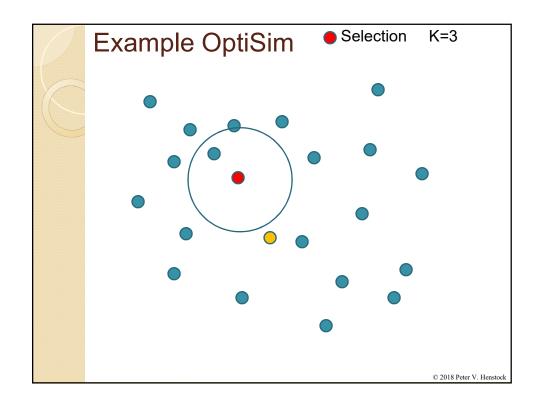
© 2018 Peter V. Henstock

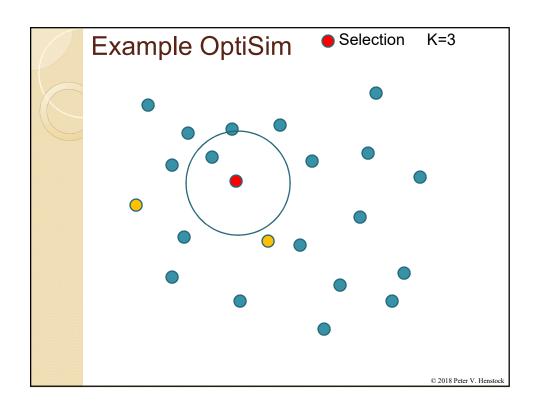
OptiSim Explained

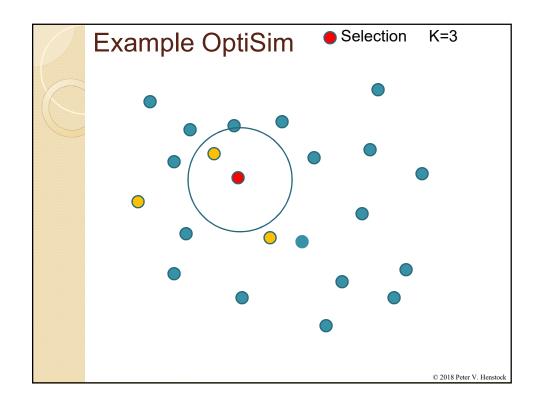
- All samples are ≥ threshold apart
 - Throw out all the close samples
- Sample set is essentially a random sample from the data set
 - Avoids extreme values that are rare
 - Chooses data representative of the data set
- Optimizing from sample set by choosing most distant point

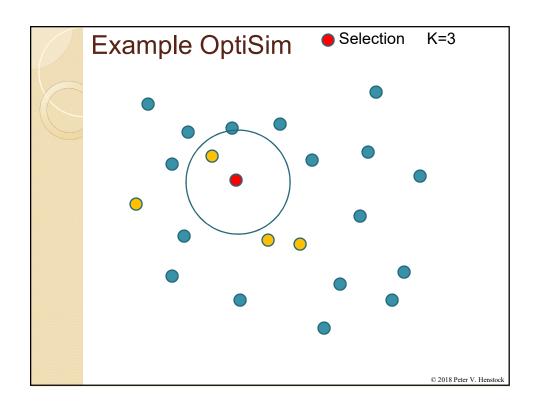


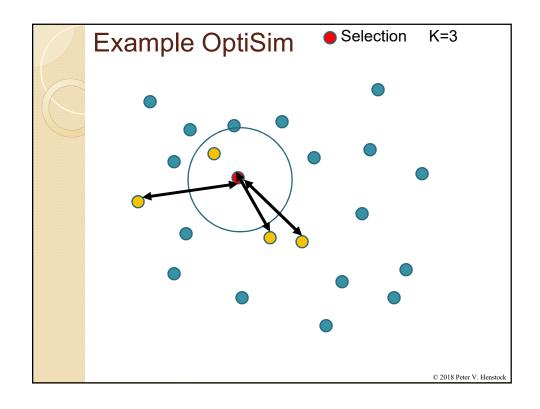


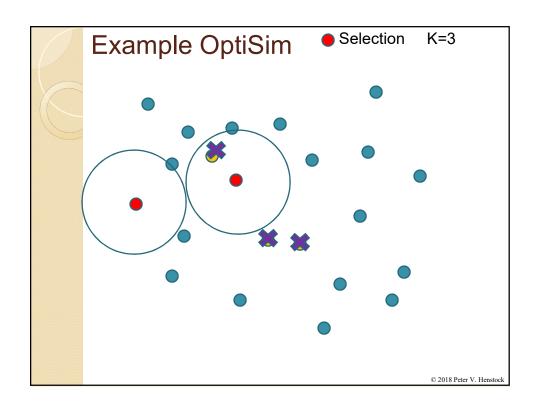


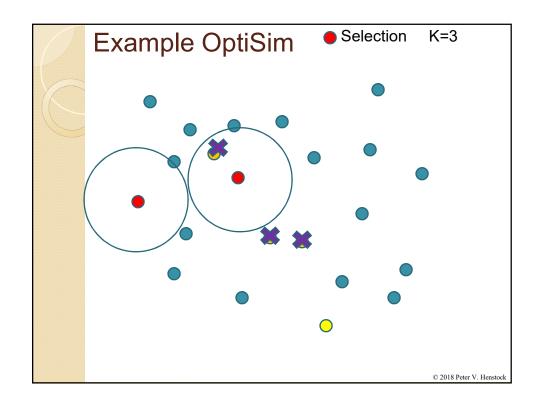


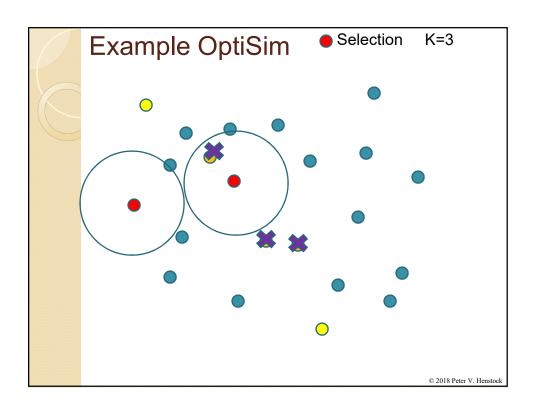


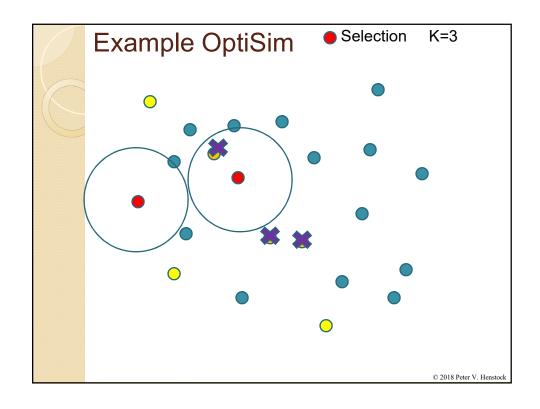


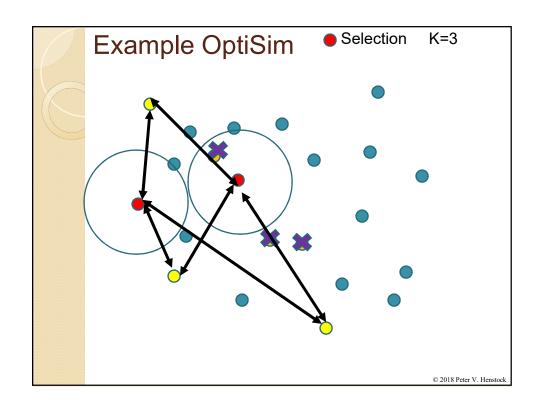


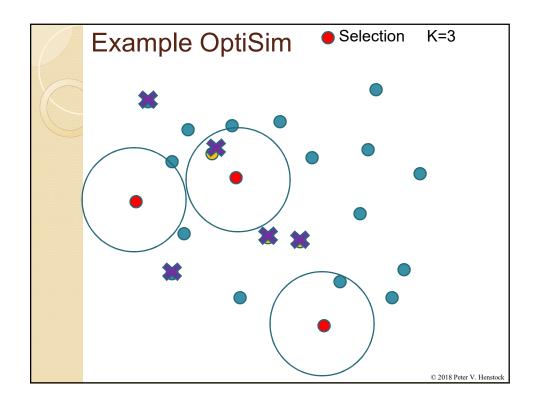










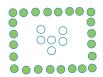


Jarvis Patrick

Irregularly Shaped Clustering

- Jarvis-Patrick and many more in literature
- Tend to be more application-dependent and less useful for statistical analysis
- Good to know they exist



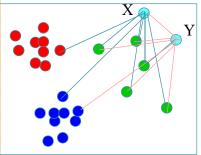


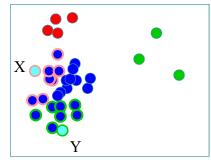
- Jarvis-Patrick example definition:
 - X and Y are in the same cluster if they share N% of their nearest M neighbors

© 2018 Peter V. Henstock

Jarvis-Patrick

- X and Y are in the same cluster if they share
 N% of their nearest M neighbors (Left)
- But, there is a gradual merging of clusters through transitive property (Right)





Clustering Using a Similarity Measure Based on Shared Near Neighbors. Jarvis & Patrick, IEEE Trans. Computers Vol. C-22 No. 11, November 1973

Good, Bad & Ugly of Jarvis-Patrick

- Good?
- Bad?
- Ugly?

© 2018 Peter V. Henstock

Good, Bad & Ugly of Jarvis-Patrick

- Good
 - Allows for interesting shapes that may be useful
 - Exploits peer-pressure relationships
 - Found to be useful for many applications
 - Provides a definition for a cluster
- Bad
 - Similar in complexity to hierarchical since have to find the closest compounds that is similar to all pairwise calculation
 - Nearest neighbor in a sparse region is may be far away
 - · Some implementations include distance threshold as well
 - Simple cluster definition evolves into complex relationship as clusters grow
- Ugly
 - Quite complex to understand how distant points are related, particularly in > 2 dimensions

Hierarchical Clustering

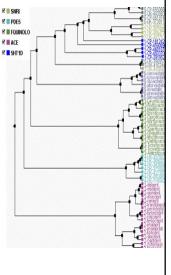
© 2018 Peter V. Henstock

Approaches to Hierarchical

- Agglomerative
 - Start: each data point is a separate cluster
 - Merge clusters until 'done'
- Divisive
 - Start: all data points are in single cluster
 - Divide and conquer until 'done'
- ~80% clustering algorithms are bottomup since there are efficient methods
 - Almost all hierarchical are agglomerative

Hierarchical Clustering

- Produces "dendrogram"
 - Tree-representation
- 2 children per parent
 - (almost always)
- Represented in .nhx files
 - New Hampshire string format:
 - (((a,b),(c,(d,e)),(f,g))



© 2018 Peter V. Henstoo

Two Approaches to Hierarchical

- Distance matrix approach
 - Pre-compute a distance matrix
 - Store distance matrix in memory
 - Efficient for smallish data < 30K data rows
- Reciprocal Nearest Neighbor (RNN)
 - Computes distances on the fly
 - Efficient for large data sets

Distance Matrix Method

- Compute all pairwise distances
- Using an appropriate "distance metric"
 - Triangle inequality, dist >= 0, dist(x,y)=dist(y,x)
 - Only need to compute upper/lower triangle

	Α	В	С	D	Е
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
Е					0

© 2018 Peter V. Henstock

Distance Matrix Method

- Repeat until there is only one tree:
 - Find the two "closest" sub-trees
 - Merge these "closest" sub-trees
 - Update the distances to all other nodes

	Α	В	С	D	Ε
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
Ε					0

- Most similar entries are C & D
- Merge C and D; C' = (C,D)

Distance Matrix Method

- Repeat until there is only one tree:
 - Find the two "closest" sub-trees
 - Merge these "closest" sub-trees
 - Update the distances to all other nodes

	Α	В	С	D	Е
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
Ε					0

- Most similar entries are C & D
- Merge C and D; C' = (C,D)
- Now what do we do?

© 2018 Peter V. Henstock

Updating distances

 How to calculate the new distance between a leaf and a subtree? Or two subtrees?





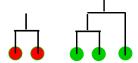


- How would you approach the problem?
- What would you do?

Updating distances

 How to calculate the new distance between a leaf and a subtree? Or two subtrees?



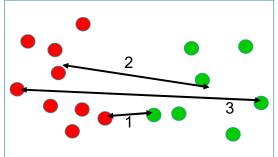


- Answer is defined as the "linkage" and there are many choices:
 - Min
 - Max
 - Median
 - Centroid
 - Average pairwise
 - Minimum variance

© 2018 Peter V. Henstock

Linkage Criteria

- Linkage is the criteria that is used as the basis of calculating distance between subtrees and/or data
- What is the distance between red and green cluster?



Linkage Criteria

- Linkage is the criteria that is used as the basis of calculati distance between subtrees and/or data
- Distance between two trees:
 - Closest pair = single linkage
 - · Great for identifying fully separable clusters
 - · Tends to suffer from "chaining problem"
 - Furthest pair = complete linkage
 - Often creates very small clusters
 - Ward's pair = minimize variance of merge
 - · Tends to work well for compounds
 - www.chemaxon.com/conf/Eurocombi_poster_Ltr.pdf
 - Average linkage: average all pairwise distances
 - · Centroid linkage: centroids of clusters

. . .

© 2018 Peter V. Henstock

Update Formula (Euclidian)

rigure 2 Aggion	nerative clustering schemes.	
Name	Distance update formula	Cluster dissimilarity
1	Formula for $d(I \cup J, K)$	between clusters A and B
single	$\min(d(I,K),d(J,K))$	$\min_{a \in A, b \in B} d[a, b]$
complete	$\max(d(I,K),d(J,K))$	$\max_{a \in A, b \in B} d[a, b]$
average	$n_I d(I, K) + n_J d(J, K)$	$\frac{1}{\sqrt{\sum d[a,b]}}$
average	$n_I + n_J$	$\frac{1}{ A B } \sum_{a \in A} \sum_{b \in B} d[a, b]$
weighted	d(I,K) + d(J,K)	
	2	VII.
Ward $\sqrt{\frac{(n_I - n_I)^2}{n_I}}$	$+ n_K)d(I, K) + (n_J + n_K)d(J, K) - n_Kd(I, J)$	$\sqrt{\frac{2 A B }{ A + B }} \cdot \vec{c}_A - \vec{c}_B _2$
V	$n_I + n_J + n_K$	$\sqrt{ A + B }$
centroid	$n_I d(I, K) + n_J d(J, K) = n_I n_J d(I, J)$	117 7 11
сепатон	$\sqrt{\frac{n_I d(I, K) + n_J d(J, K)}{n_I + n_J} - \frac{n_I n_J d(I, J)}{(n_I + n_J)^2}}$	$\ \vec{c}_A - \vec{c}_B\ _2$
median	d(I,K) + d(J,K) - d(I,J)	$\ \vec{w}_A - \vec{w}_B\ _2$
mount	V 2 2 4	$\ w_A - w_B\ _2$

Legend: Let I, J be two clusters joined into a new cluster, and let K be any other cluster. Denote by n_I , n_J and n_K the sizes of (i.e. number of elements in) clusters I, J, K, respectively.

 Modern hierarchical, agglomerative clustering methods by Daniel Műllner

Lance-Williams Update Formula

dist(i \cup j, k) = α_i dist(i, k) + α_j dist(j, k) + β dist(i, j) + γ |dist(i, k) - dist(j, k)|

- Different linkage functions provide different values of $\alpha_i \alpha_i \beta \gamma$
- Idea: merge two clusters i and j
- Update distance to k based on the distances between (i,j), (i,k) and (j,k)
- Assume Euclidian distance

© 2018 Peter V. Henstock

Lance-Williams Update Formula

dist(i \cup j, k) = α (i) dist(i, k) + α (j) dist(j, k) + β dist(i, j) + γ |dist(i, k) – dist(j, k)|

Hierarchical clustering methods (and	Lance and Williams dissimilarity update formula	McQuitty's method (WPGMA)	$\alpha_i = 0.5$ $\beta = 0$ $\gamma = 0$
Single link (nearest	$\alpha_i = 0.5$ $\beta = 0$	Median method (Gower's, WPGMC)	$\alpha_i = 0.5$ $\beta = -0.25$ $\gamma = 0$
neighbor)	$\gamma = -0.5$ (More simply: $min\{d_{ik}, d_{jk}\}$)	Centroid (UPGMC)	$\alpha_i = \frac{ i }{ i + j }$ $\beta = -\frac{ i j }{(i + j)^2}$ $\gamma = 0$
Complete link (diameter)	$\alpha_i = 0.5$ $\beta = 0$ $\gamma = 0.5$ (More simply: $max\{d_{ik}, d_{jk}\}$)	Ward's method (minimum var- iance, error	$\alpha_i = \frac{ i + k }{ i + j + k }$ $\beta = -\frac{ k }{ i + j + k }$ $\gamma = 0$
Group average (average link,	$\alpha_i = \frac{ i }{ i + j }$ $\beta = 0$ $\alpha_i = 0$	tn://arxiv.org/pd	F/1105 0121 mds

Hierarchical Clustering Algorithm

- Find all pairwise distances
- Repeat until there is only one tree:
 - Find the two "closest" sub-trees
 - Merge these "closest" sub-trees
 - Update the distances to all other nodes

	Α	В	С	D	Ε
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
					^

- Most similar entries are C & D
- Merge C and D; C' = (C,D)
- Recalculate the distance between C' and {A,B,E}

© 2018 Peter V. Henstock

Example update for Centroid

	Α	В	С	D	E
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
Е					0

Centroid (UPGMC)

$$\alpha_i = \frac{|i|}{|i|+|j|}$$

$$\beta = -\frac{|i||j|}{(|i|+|j|)^2}$$

$$\gamma = 0$$

- dist(i \cup j, k) = α_i dist(i, k) + α_j dist(j, k) + β dist(i, j) + γ |dist(i, k) dist(j, k)|
- dist(C \cup D, A) = α dist(C, A) + α dist(D, A) + β dist(C, D) + γ |dist(C, A) dist(D, A)|
- $\frac{1}{2}(0.8) + \frac{1}{2}(0.4) \frac{1}{4}(0.2) + 0 = 0.55$
- Repeat for dist(C∪D, B) & dist(C∪D, E)

Updating distances

- Potentially extremely computationally complex
 - Average pairwise would require recomputing distance between all nodes in the subtree for each candidate
- There are fast-update equations that update distance calculations every merge quickly
 - But...Only work where triangle equality holds!
- Example: Centroid:
 - d(p+q,l) = np/(np+nq)*d(p,l) + (nq/(np+nq)*d(q,l) ni*d(p,q)
 - www.soziologie.wiso.uni-erlangen.de/ koeln/script/chap4.pdf
- Notion is you can quickly update your distance to all other objects after every merge from pre-merge distances

© 2018 Peter V. Henstock

Practical Data Structures

Distance matrix upper triangle: ~O(N²/2)

	Α	В	С	D	E
Α	0	0.9	0.8	0.4	0.5
В		0	0.7	0.3	0.4
С			0	0.2	0.3
D				0	0.8
Ε					0

- Sorted minimum value of each column
 - Replace N² search for minimum with N
 - ∘ [0.9, 0.7 0.2 0.3] needs updating each merge
- Table: ordered merges x, y (and distance)
- Indexed array of cluster members and #members in each

Practical Data Structures

- Distance matrix upper triangle
- Sorted min of columns
 - · [0.9, 0.7 0.2 0.3]
 - needs updating each merge
- Table of merges

 Table of clusters 	(D active?)
---------------------------------------	-------------

$A \rightarrow A$	1	$A \rightarrow A$	1
$B \rightarrow B$	1	$B \rightarrow B$	1
$C \rightarrow C$	1	$C \rightarrow C,D$	
$D \rightarrow D$	1	D → null	0
$E \rightarrow E$	1	$E \rightarrow E$	1

X	Y	Dist
С	D	0.20
С	E	0.24
В	С	0.56
Α	В	0.62

© 2018 Peter V. Henstoc

0.3 0.4 0.2 0.3

0.5

0.8

0

0 0.9 0.8 0.4

0 0.7

Big O applied to Hierarchical

- Find all pairwise distances
 - O(n²) in space; O(n²) in time Repeat until there is only one tree:
 - Find the two "closest" sub-trees
 - $-O(n^2)$ in time but can do in O(n) = $O(n^2)$
 - ▶ Merge these "closest" sub-trees
 - -O(1) O(n)
 - Update the distances to all other nodes
 - $-O(n) \times n = O(n^2)$
- Conclusion O(n²) in space and time
- 10,000 data points →~1E8 time/space units
- 1,000,000 data points → ~1E12 time/space units

Time and Space



- How much time?
- How much memory?

© 2018 Peter V. Henstock

Time and Space

- How much time?
 - Compute distance matrix O(N²)
 - For each data: N times
 - Find min distance → O(N)
 - Update N distances → O(N)
- How much memory?
 - ∘ Distance matrix $O(N-1)^2/2)$ → $O(N^2)$

Good, Bad and Ugly of Hierarchical Good? Bad? Ugly?

Good, Bad and Ugly of Hierarchical Good Very intuitive for end users **₩ M** ACE Very few parameters required ₹¶5HT1D distance function and linking function · Cut the tree into "clusters" afterwards Define clusters after clustering Bad Takes a lot of memory and speed Do about 25,000 on a laptop in an hour Challenging to do significantly more Ugly Can't understand how clusters relate No definition of a cluster within the hierarchy Spatial arrangement is meaningless between clusters © 2018 Peter V. Henstock

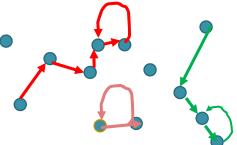
Reciprocal Nearest Neighbor

- Only works for a subset of hierarchical clustering methods
- On-the-fly clustering
- Avoids pre-computing the distance matrix
- Repeat until done
 - Take a data point or cluster u1
 - ∘ Find closest point of u1 → u2
 - ∘ Find closest point of u2 → u3
 - 0
 - If closest(ui) == u(i-1) merge them

© 2018 Peter V. Henstoc

Reciprocal Nearest Neighbor

- Form a RNN-chain that identifies the leaf of the dendrogram
- Lots of mathematical/geometric proofs that it works

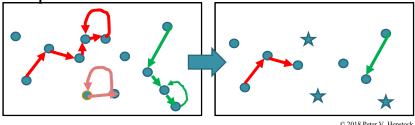


What's next step after finding each RNN?

2018 Peter V. Henstock

Update formula

- Previously computed the distance to all the other nodes d(aUb,*)
- However, we haven't computed the distances of d(a,*) or d(b,*)
- Replace the merged cluster with a replacement node



Update Formulae

			17
Median method (Gower's, WPGMC)	$\alpha_i = 0.5$ $\beta = -0.25$ $\gamma = 0$	$g = \frac{g_i + g_j}{2}$	$\ \mathbf{g}_i - \mathbf{g}_j\ ^2$
Centroid (UPGMC)	$\alpha_i = \frac{ i }{ i + j }$ $\beta = -\frac{ i }{(i + j)^2}$ $\gamma = 0$	$g = \frac{ i g_i + j g_j}{ i + j }$	$\ \mathbf{g}_i - \mathbf{g}_j\ ^2$
Ward's method (minimum var- iance, error sum of squares)	$\alpha_i = \frac{ i + k }{ i + j + k }$ $\beta = -\frac{ k }{ i + j + k }$ $\gamma = 0$	$g = \frac{ i g_i + j g_j}{ i + j }$	$\frac{ i j }{ i + j } \ \mathbf{g}_i - \mathbf{g}_j\ ^2$

 Use the g's to compute the new centroid from those merged

Recommended paper on RNN Methods of Hierarchical Clustering by Murtagh & Contreras 2011 http://arxiv.org/pdf/1105.0121.pdf

Recent Work

- If you assume you will cut into clusters anyway
 - Pre-cluster using a fast technique (I.e. leader)
 - Perform hierarchical clustering on clusters with RNN
 - Takes kO(m²) where m<n so avoid full n x n matrix
- Reciprocal Nearest Neighbor (RNN)
 - Only works for some distance/linkage combos
 - Take entry (j) and find nearest neighbor of j: k
 - Repeat to find NN(k) until NN(k)=NN(j) which is RNN
 - NN(x1)=x2; NN(x2)=x3; NN(x3)=x4; ...NN(x14)=x13
 - RNN chain will form an independent subtree

© 2018 Peter V. Henstock

RNN Movie Link & References

- http://www.juergenwiki.de/work/wiki/dok u.php?id=public:reciprocal_nearest_neig hbour clustering rnn
- Recommended paper on RNN
- Methods of Hierarchical Clustering by Murtagh & Contreras 2011 http://arxiv.org/pdf/1105.0121.pdf

Assessing Clustering

© 2018 Peter V. Henstock

Cluster Validation

- "The validation of clustering structures is the most difficult and frustrating part of cluster analysis."
- "Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage."

Algorithms
for Clustering Data

And C. Jain
Plate Date University

Page 2144

General Assessment

- Do we know the correct clustering?
 - "Extrinsic" = Supervised
 - Classification rates
 - Information theoretic or Entropy
 - Purity
 - Rand Index
- No ground truth: "Intrinsic"
 - Assess goodness of clustering properties
 - Separation of data
 - Compact/density of clusters

© 2018 Peter V. Henstock

Extrinsic Measures

- Properties of good clusters
 - ∘ Homogeneity / Purity → clusters same truth
 - ∘ Completeness: same truth → same clusters
 - Rag bag
 - Tradeoff of spoiling pure cluster vs. "other" cluster
 - Small cluster preservation
- Set based
- Entropy
- Contingency Table
- B-Cubed

Purity Measures (Set measures)

- C = #clusters L = #truth categories
- N = #clustered data points
- Purity = $\sum_{i} \frac{|C_i|}{N} max_j Precision(C_i, L_j)$
- $Precision(C_i, L_j) = \frac{|C_i \cap L_j|}{|C_i|}$

© 2018 Peter V. Henstock

Contingency Tables

- Set of truth labels A,B,C,D
- Set of cluster labels a,b,c,d
 - Might have a,b→A c→C,D d→D

Counts of Pairs	Same Cluster	Different Cluster
Same Truth	SS	DS
Different Truth	SD	DD

- Rand statistic: $\frac{SS+DD}{SS+SD+DS+D}$
- Jaccard Coefficient J = $\frac{SS}{SS+SD+DS}$
- Folkes & Mallows FM = $\sqrt{\frac{SS}{SS + SD}} \frac{SS}{SS + DOS}_{Peter V. Henstock}$

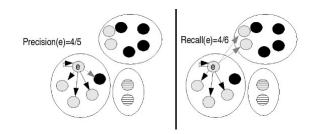
Entropy

$$\text{Entropy} = -\sum_{j} \frac{n_{j}}{n} \sum_{i} P(i, j) \times \log_{2} P(i, j)$$

- P(i,j) = probability of finding element of category I in cluster j
- n_i = # elements in cluster j
- N = # elements overall

© 2018 Peter V. Henstock

B-Cubed



- Compute precision & recall at every point
- Precision ~ purity
 - Fraction of items in with same truth in cluster
- Recall ~ distribution
 - Fraction of total items with same truth in that particular cluster

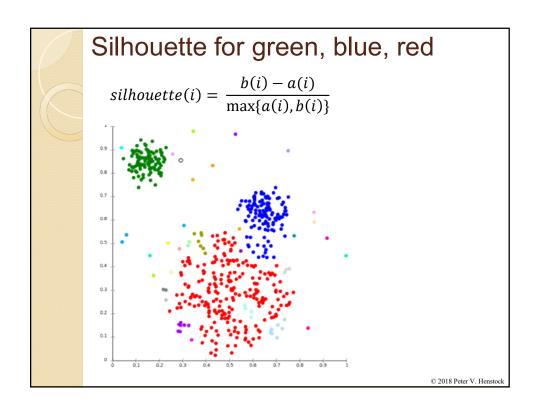
For more information

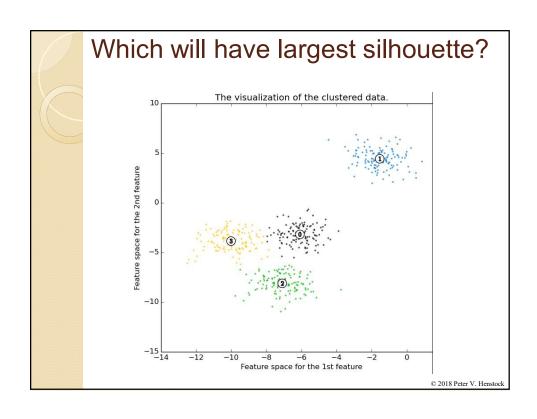
- A comparison of extrinsic clustering evaluation metrics on formal constraints
- Amigo et al. 2009
- http://nlp.uned.es/docs/amigo2007a.pdf

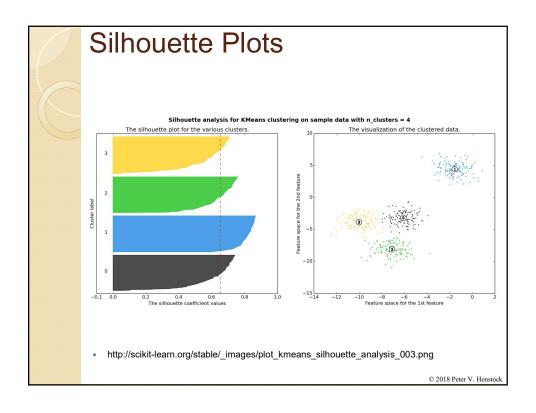
© 2018 Peter V. Henstoc

Unsupervised: Silhouette Measure

- Cluster the data first
- a(i) = average dissimilarity of data i with all other data in the same cluster
 - ~Intra cluster dissimilarity
- b(i) = lowest average dissimilarity of i to any other cluster to which is not a member
 - Dissimilarity to nearest neighbor cluster
- $silhouette(i) = \frac{b(i)-a(i)}{\max\{a(i),b(i)\}}$
 - Near 1 if tight cluster with high separation
 - Near -1 loose cluster with others close fis Peter V. Henstock





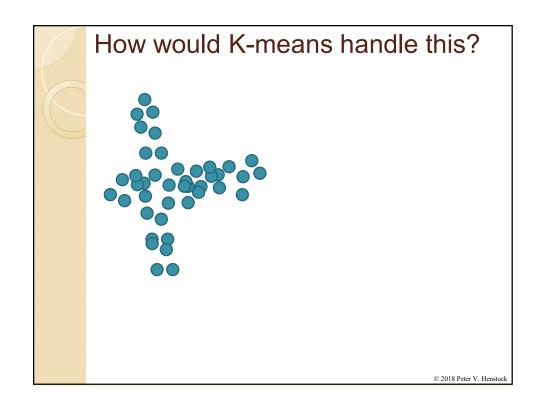


Assessing Number of Clusters

- Empirical \rightarrow k = \sim sqrt(N) N=data pts.
- Elbow/Knee method
- Cross-validation
 - Discuss this in classification
 - Divide data up into pieces and estimate

0

Mixture Models GMM

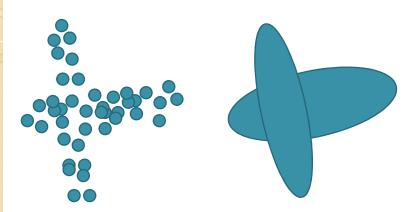


K-Means Issues

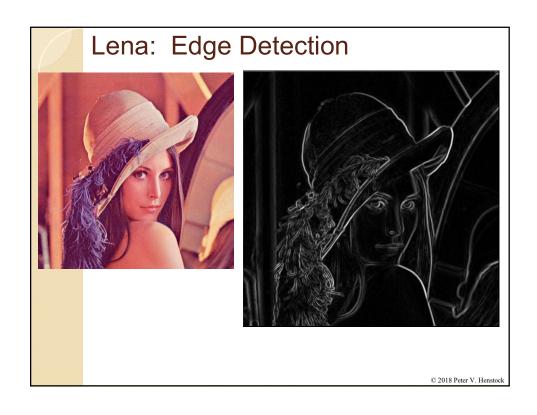
- Likes round/spherical clusters
 - Based on the distance metric
- Hard-limiter:
 - Points are in only one cluster
 - Probability = {0, 1} of membership in cluster
- Statistically, it a centroid and a non-statistical ad hoc radius
- Many fields prefer Gaussian distributions over round 'blob' clusters

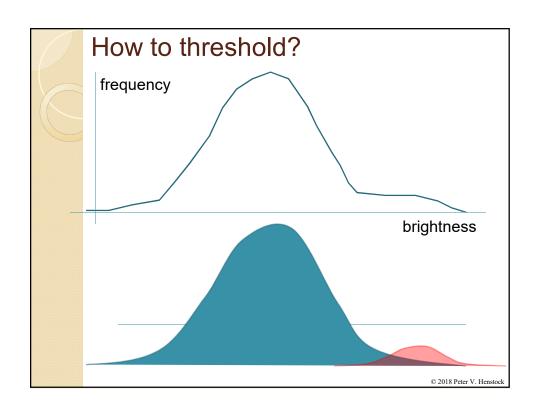
© 2018 Peter V. Henstock

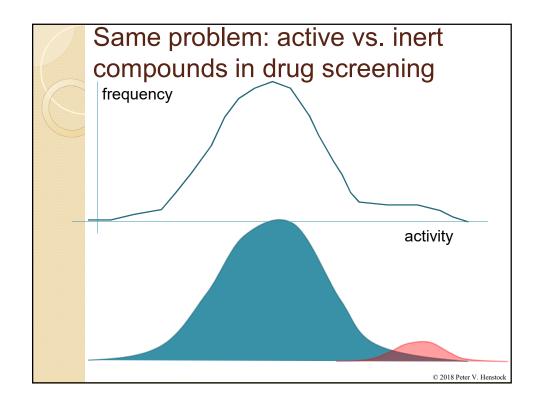
How would K-means handle this?

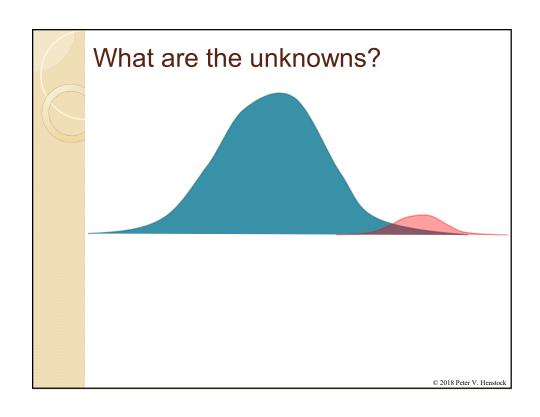


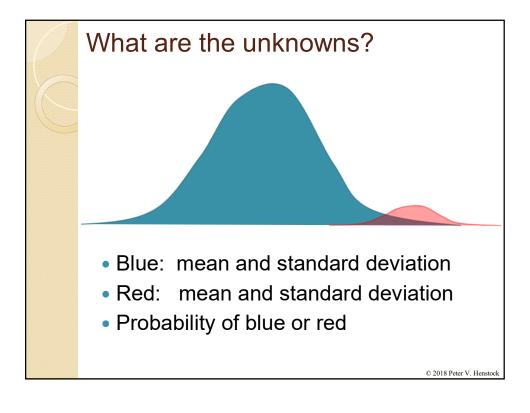
Gaussian distribution is more satisfying

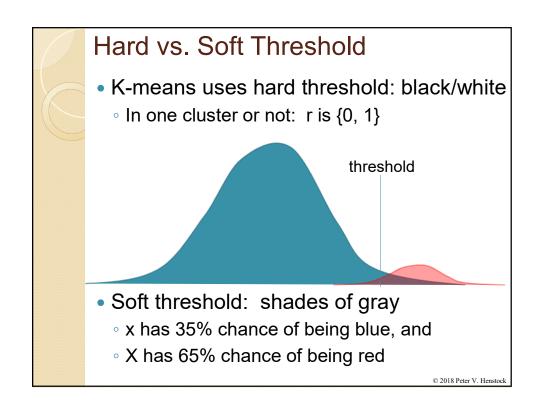












Gaussian Mixture Model



- Actually not mixing anything
- Using multiple
 Gaussians to
 represent data
 instead of just one

https://clipartxtras.com/categories/view/90a208683ac530bb5e65d28d8636d2a5baf35606/mix-drawing.html.

© 2018 Peter V. Henstock

Mixture Models

- $P(x) = \sum_{c=0}^{\#clusters} P(cluster_c)P(x|cluster_c)$
- Generative model for P(x) = one point
 - Pick a cluster at random P(cluster_c)
 - Assume each cluster has a probability
 - Gaussian prior
 - Generate a value X from the Gaussian distribution of cluster_c
 - Sum over values to get P(x) above
 - Also called a "latent" model
- Max likelihood solution:
 - Choose clusters so they are most likely given the data → maximize P(x)

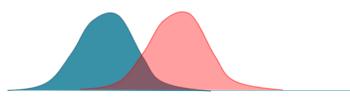
Mixture Models

- Models for p(x|c_i)
 - Normal distributions is standard approach
 - Assume iid so can use a product
- $P(x|cluster_c) = \prod_{j=1}^{N} P(xj|cluster_c)$
- What are we trying to find?

© 2018 Peter V. Henstock

Mixture of Gaussians

$$P(x|\mu_i,\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[\frac{-(x-\mu_i)^2}{2\sigma^2}\right]$$



$$\mathcal{N}(\underline{x}~;~\underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})\right\}$$

How to learn mixture models?

- "Estimation Maximization" algorithm
- = EM algorithm
- Maximizes likelihood in presence of 'missing data' i.e. the clusters
- Defines solution as f(g) and g(f) and iterates the optimization between each
- 2 step algorithm

© 2018 Peter V. Henstock

EM Algorithm

- Initialize parameters
- Loop:
 - E-step:
 - Use current parameters to compute estimates of unobserved variables
 - M-step
 - Update parameter values to maximize the probability of the observed and unobserved data

EM Algorithm

- Initialize parameters (perhaps K-means)
- Loop:
 - E-step—estimating the assignments
 - Use current parameters to compute estimates of unobserved variables
 - Compute probability that point xi maps to cluster_c
 - M-step—maximizing the values
 - Update parameter values to maximize the probability of the observed and unobserved data
 - Compute probability of cluster_c
 - Compute centroid μ_c
 - Compute covariance Σ or stdev σ of cluster c

© 2018 Peter V. Henstock

E-step

- Initialize means (perhaps using K-means)
- For each point x_i: compute

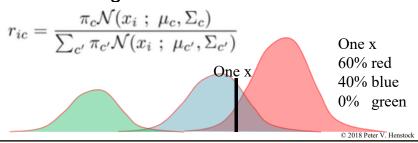
•
$$P(\mu_c|x_j) = \frac{P(x_j|\mu_c)P(\mu_c)}{P(x_j)} = \frac{P(x_j|\mu_c)P(\mu_c)}{\sum_{k=1}^{\#clust}P(\mu_k)P(x_j|\mu_k)}$$

• Recall
$$P(x|\mu_c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left[\frac{-(x-\mu_c)^2}{2\sigma_c^2}\right]$$

 Changing notation, we are computing r just as we did the K-means

E-Step Another way

- $P(\mu_c|x_j) = \frac{P(x_j|\mu_c)P(\mu_c)}{P(x_j)} = \frac{P(\mu_c)P(x_j|\mu_c)}{\sum_{k=1}^{\#clust} P(\mu_k)P(x_j|\mu_k)}$
- In K-means, we computed which x's were assigned to which cluster centroids
- Here, we are using the likelihood values to figure out which cluster it belongs to
- Normalizing in denominator so sums to 1



E-Step Yet another way

- Removing the hard-limit "closest" assignment r of each x_i to cluster c_i
- Adding a probabilistic soft-limit that assigns x_i to each cluster c_i with r in [0,1]
- Each x_i has to belong to some cluster with a probability so we scale them so that require $\Sigma r = 1$
- Sum of assignments of point xi to all clusters sums to 1 over all clusters

M-Step

- M step compute parameters: π_c , μ_c , σ_c
- P(cluster_c) = $\pi_c = \frac{1}{N} \sum_{j=1}^{N} P(\mu_c | x_j)$
 - \circ P(μ_c | x_i) = probability x_i belongs to μ_c

•
$$\mu_c = \frac{\sum_{j=1}^{N} x_j P(\mu_c | x_j)}{\sum_{j=1}^{N} P(\mu_c | x_j)}$$

Weighted average of the values in cluster

•
$$\sigma_c = \sqrt{\frac{\sum_{j=1}^{N} (x_j - \mu_c)^2 P(\mu_c | x_j)}{\sum_{j=1}^{N} P(\mu_c | x_j)}}$$

Weighted average of variance

© 2018 Peter V. Henstock

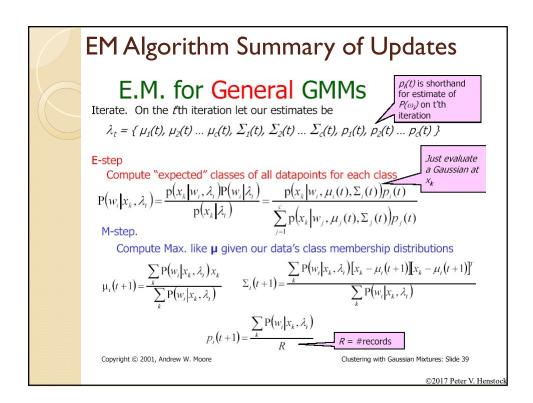
M-Step again: work with 'soft' limit

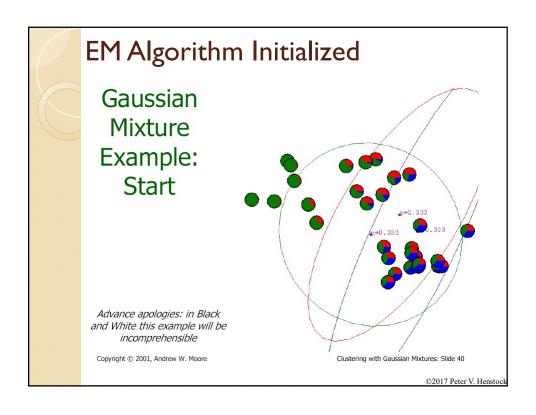
- What % of data are in cluster c?
 - \circ r_{ic} = assignment of xi to cluster c
 - $_{\circ}$ $\Sigma_{\rm i}$ $\rm r_{ic}$ = sum of all xi mapping to cluster c = $\rm m_c$
 - $_{\circ}$ Need to normalize to [0, 1] so divide by Σ_{c} m_{c}
 - \circ π_c = $m_c/$ Σ_k m_k = probability of cluster c
- Mean is defined as Σ xp(x) in general
 - \circ Centroid of cluster c = Σ_{i} r_{ic} xi / m_c
 - If xi is strongly in the cluster r large and impacts the centroid
- Covar cluster c= $\frac{1}{m_c} \sum_i r_{ic} [x_i \mu_c]^T [x_i \mu_c]$

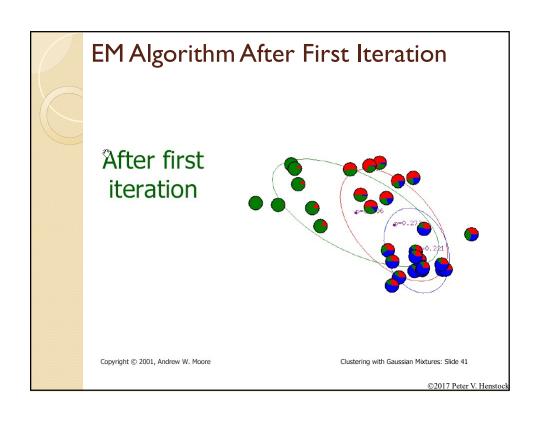
2018 Peter V. Henstock

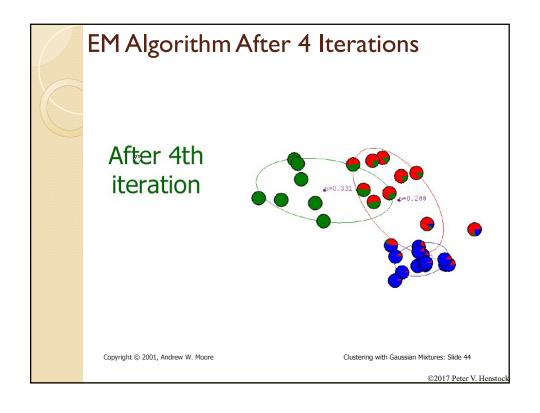
K-Means vs. EM

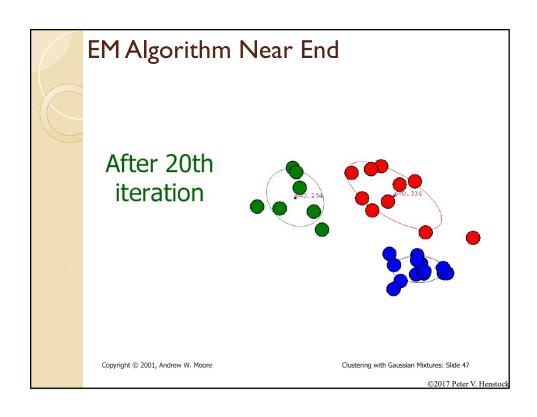
- K-Means is a special case of EM
 - Cluster priors are equal
 - Covariance matrix: identity matrix
 - Means = means
- Assignments in K-Means are not a probability but a hard limit











Good, Bad & Ugly of EM

Good

Bad

Ugly

© 2018 Peter V. Henstock

Good, Bad & Ugly of EM

- Good
 - Statistically accurate and robust to various start points
 - Reasonably fast
 - Ideal for medium to large clusters
 - "Soft" cluster memberships
- Bad
 - Assumes have Gaussian distributions
 - Small isolated clusters work poorly
 - Need to know how many clusters you have
 - Can sometimes be fussy about initialization
- Ugly
 - Quite beautiful and elegant when assumptions apply

2018 Peter V. Henstock

Density Based Techniques

- Mean Shift Clustering
- DBSCAN
- OPTICS