



CSCI E-82

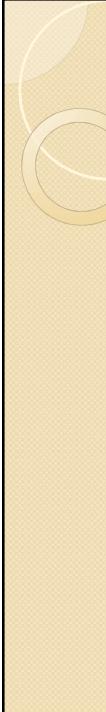
Advanced Machine Learning, Data Mining & Artificial Intelligence

Lecture 2

Peter V. Henstock

Fall 2018

© 2018 Peter V. Henstock



Statistics vs. Machine Learning

- How are they the same?
- How are they different?

© 2018 Peter V. Henstock

Statistics vs. Machine Learning

- “Data Science without statistics is possible, even desirable”
 - <https://www.datasciencecentral.com/profiles/blogs/data-science-without-statistics-is-possible-even-desirable>
- Statistics:
 - Make assumptions
 - Derive formulae
 - Assert properties like “significance”
- Machine learning:
 - Set up an optimal criteria
 - Figure out a computational way of solving it

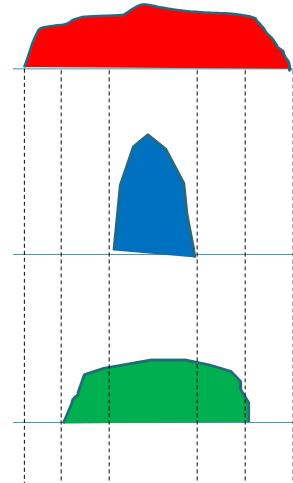
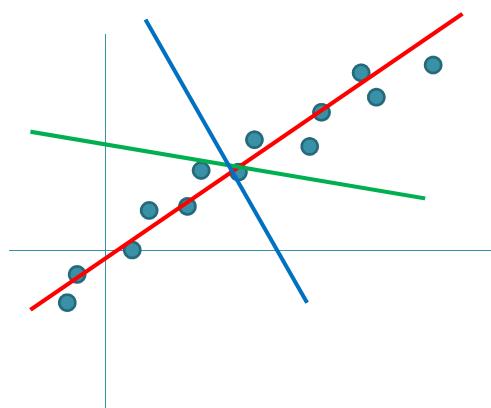
© 2018 Peter V. Henstock

Principal Component Analysis

© 2018 Peter V. Henstock

1st Principal Component Axis

- Another Idea: maximize variance of the distribution projected onto the line



© 2018 Peter V. Henstock

PCA Solution

- Trying to find axes of maximize variance
- New axes represent a “subspace”
- SVD
 - Breaks matrix into USV^T such that
 - S is a diagonal matrix of eigenvalues
 - Columns of U are eigenvectors
 - Eigenvectors of largest eigenvalues are the axes corresponding to maximum variance

© 2018 Peter V. Henstock

Optimization of PCA

- $A \sim USV^T$ where A is covariance matrix
- Minimizing $\|A - USV^T\|_F^2$
 - Frobenius Norm
 - Technically written $\|A - suv^T\|_F^2$
 - Unique minimum when $u^Tu=1$, $v^Tv=1$, $s>0$
- Achieve minimum using SVD for different numbers of eigenvalues (or eigenvectors)

© 2018 Peter V. Henstock

Purge Eigenvalues

U

-0.5930	0.3879	-0.0341	-0.7048
-0.4779	-0.8362	-0.2654	-0.0453
-0.4527	-0.0221	0.8286	0.3287
0.4638	-0.3871	0.4918	-0.6271

S

0.1491	0	0	0
0	0.0820	0	0
0	0	0.0496	0
0	0	0	0.0128

V'

-0.5930	-0.4779	-0.4527	0.4638
0.3879	-0.8362	-0.0221	-0.3871
-0.0341	-0.2654	0.8286	0.4918
-0.7048	-0.0453	0.3287	-0.6271

Original Cov

0.0712	0.0165	0.0350	-0.0485
0.0165	0.0949	0.0227	-0.0126
0.0350	0.0227	0.0660	-0.0130
-0.0485	-0.0126	-0.0130	0.0614

Cov Approximation 1 eigenvalues

0.0648	0.0161	0.0379	-0.0541
0.0161	0.0949	0.0229	-0.0130
0.0379	0.0229	0.0646	-0.0104
-0.0541	-0.0130	-0.0104	0.0563

Cov with 2 eigenvalues

0.0648	0.0156	0.0393	-0.0533
0.0156	0.0914	0.0338	-0.0065
0.0393	0.0338	0.0306	-0.0306
-0.0533	-0.0065	-0.0306	0.0444

Cov with 3 eigenvalues

0.0524	0.0422	0.0400	-0.0410
0.0422	0.0340	0.0322	-0.0330
0.0400	0.0322	0.0305	-0.0313
-0.0410	-0.0330	-0.0313	0.0321

© 2018 Peter V. Henstock

Process Check

$X^*U[1:k] \rightarrow \text{projection}$
If $k = 2$, what dimensionality will we have?

X vectors with 4 features

0.0344	0.4387	0.3816	0.7655
0.7952	0.1869	0.4898	0.4456
0.6463	0.7094	0.7547	0.2760
0.6797	0.6551	0.1626	0.1190
0.4984	0.9597	0.3404	0.5853
0.2238	0.7513	0.2551	0.5060
0.6991	0.8909	0.9593	0.5472
0.1386	0.1493	0.2575	0.8407
0.2543	0.8143	0.2435	0.9293
0.3500	0.1966	0.2511	0.6160

SVD of (cov(A)) $\rightarrow U, S, V$
U columns are eigenvectors

-0.5930	0.3879	-0.0341	-0.7048
-0.4779	-0.8362	-0.2654	-0.0453
-0.4527	-0.0221	0.8286	0.3287
0.4638	-0.3871	0.4918	-0.6271

X matrix 10 x 4 * U matrix is 4x2 after reduction

© 2018 Peter V. Henstock

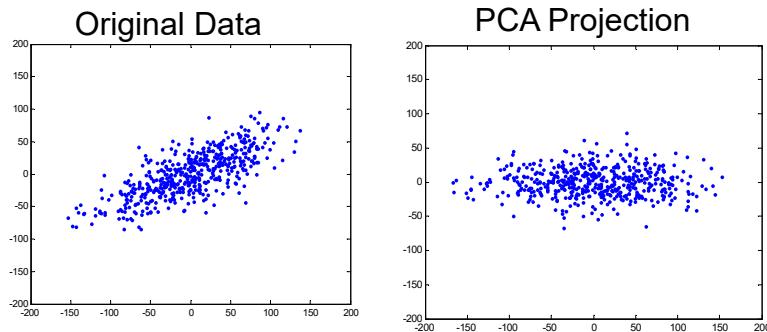
Principal Components

- Create a weighted sum of features
- Use all the features together in one shot
- Each feature' = $\sum_{\text{feature}_j} \text{weight}(j) * \text{feature}(j)$
- Weights are not equal
 - First several are much stronger than rest
 - Eliminate the lowest ranked transformations

© 2018 Peter V. Henstock

Projection

- Generated a random blob and rotated it by $\pi/6$ (Left)
- 1st PC X-axis; 2nd PC is Y-axis
- Dominant PC axis captures the variation



© 2018 Peter V. Henstock

PCA for original feature selection?

- We have lots of features
- We can use a PCA to map to a lower dimensional space
- Can we use PCA for feature selection?
- Can we select high-value features?
- Can we purge low-value features?

© 2018 Peter V. Henstock

PCA for original feature selection?

- We have lots of features
- We can use a PCA to map to a lower dimensional space
- Can we infer which of the original features are less meaningful?
- $\text{PCA1} = w_{11} * \text{feat1} + w_{12} * \text{feat2} + \dots$
- $\text{PCA2} = w_{21} * \text{feat1} + w_{22} * \text{feat2} + \dots$
- $\text{PCA3} = w_{31} * \text{feat1} + w_{32} * \text{feat2} + \dots$

© 2018 Peter V. Henstock

Latent Semantic Indexing

- Document matrix: documents x words
- Consider query dot-product $\text{doc} =$ similarity measure
- Use SVD to create low rank space
- Project query into PC space
- Assumption that need reduced space
- Document Query_k = $\Sigma_k U_k^T \text{query}$
 - Find similar documents
- Term Query_k = $\text{query}^T \Sigma_k V_k^T$
 - Find similar terms

© 2018 Peter V. Henstock

What about outliers?

- Using mean and covariance
- Both are sensitive to outliers
- PCA focuses on preserves pairwise distances—particularly large ones

© 2018 Peter V. Henstock

Robust PCA

- Instead of using mean and variance, it likely uses median and robust standard deviation
- Penalizes the PCA in a way that avoids the huge quadratic cost for outliers
- More complicated solution

© 2018 Peter V. Henstock

Non-Negative Matrix Factorization

- Same idea as PCA but all weights ≥ 0
- Usually uses a different factorization
 - $X = WH$ with same $\|X - WH\|_F^2$
- Face = FacialFeatures * FeatureWeight
- Document = Topic * TopicImportance
- Image = MemberSignal*MemberAbundance
- GeneExp = Specimens*SpecimenWeights
- BrainMeasurements = NeuroProc*Weights
- NP-Hard problem so no SVD

© 2018 Peter V. Henstock

Independent Component Analysis

- Unsupervised learning
- $X_1 \dots X_N$ are hidden variables
- $Y_1 \dots Y_k$ are the observables ($k \geq N$)
 - Y are independent $\text{Info}(Y_i, Y_j) = 0$
 - Maximize Y to X so $\text{Info}(Y_i, X)$ is large
 - Y is linear function of weighted X so $Y = AX$
- “Cocktail Party”
- “Blind Source Separation”
- Breaks world down into components
 - PCA of faces finds average faces
 - ICA of faces finds noses, ears, hair

© 2018 Peter V. Henstock

Independent Component Analysis

- Two sources of ambiguity
 - Can't figure out which X_i is which
 - Can't figure out if it's X_i or $-X_i$
- Algorithm will not work if data is Gaussian
- Math is based on computing cumulative density distributions for the X

© 2018 Peter V. Henstock



- 1967 Western
- Set in civil war
- Name has become an idiom

© 2018 Peter V. Henstock

Good, Bad & Ugly of PCA

- Good?

- Bad?

- Ugly?

© 2018 Peter V. Henstock

Good, Bad & Ugly of PCA

- Good

- Single fast global transform—no parameters
- Widely used so requires less explanation
- Select the extent of the approximation (#eig)

- Bad

- Often don't get the clear delineations since it is a global transform

- Ugly

- Linear transform which sometimes doesn't do much beyond some kind of rotation depending on data
- Ignores probabilistic distribution: variance only

© 2018 Peter V. Henstock



Alternative Methods for Dimensionality Reduction

© 2018 Peter V. Henstock



Multidimensional Scaling

- PCA is linear technique
- Linear transform of the points (rotation)
- Similar approach that is a linear transform
- Goal:
 - Map higher dimensional → lower dimensional
 - Find projection that preserves the pairwise distances in lower dimensional space

© 2018 Peter V. Henstock

MDS

- Optimization function
 - $E(y) = \sum_{i \neq j} \frac{(dy_{ij} - dx_{ij})^2}{dy_{ij}}$
 - dy = pairwise difference of projection
 - dx = pairwise difference of original high dimensional points
- Note that researchers have come up with lots of different stress functions
- How to optimize this?

© 2018 Peter V. Henstock

MDS

- Optimization function or “stress function”
 - $E(y) = \sum_{i \neq j} \frac{(dy_{ij} - dx_{ij})^2}{dy_{ij}}$
 - dy = pairwise difference of projection
 - dx = pairwise difference of original high dimensional points
- Note that researchers have come up with lots of different stress functions
- How to optimize this?

© 2018 Peter V. Henstock

MDS

- Optimization function or “stress function”
 - $E(y) = \sum_{i \neq j} \frac{(dy_{ij} - dx_{ij})^2}{dy_{ij}}$
 - dy = pairwise difference of projection
 - dx = pairwise difference of original high dimensional points
- Note that researchers have come up with lots of different stress functions
- How to optimize this?
 - Gradient descent

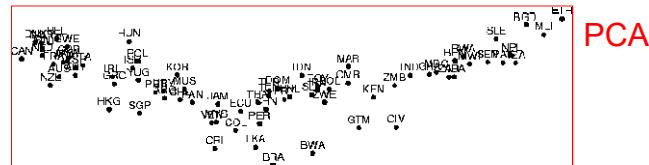
© 2018 Peter V. Henstock

Sammon's Algorithm

- Optimization function
 - $E(y) = \sum_{i \neq j} \frac{(dy_{ij} - dx_{ij})^2}{dx_{ij}}$ ←
 - dy = pairwise difference of projection
 - dx = pairwise difference between original high dimensional points
- MDS reminder
 - $E(y) = \sum_{i \neq j} \frac{(dy_{ij} - dx_{ij})^2}{dy_{ij}}$ ←

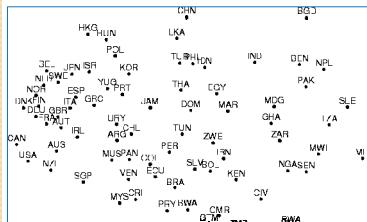
© 2018 Peter V. Henstock

World Bank Data Comparison



PCA

Sammons
Preserves small distances



MDS

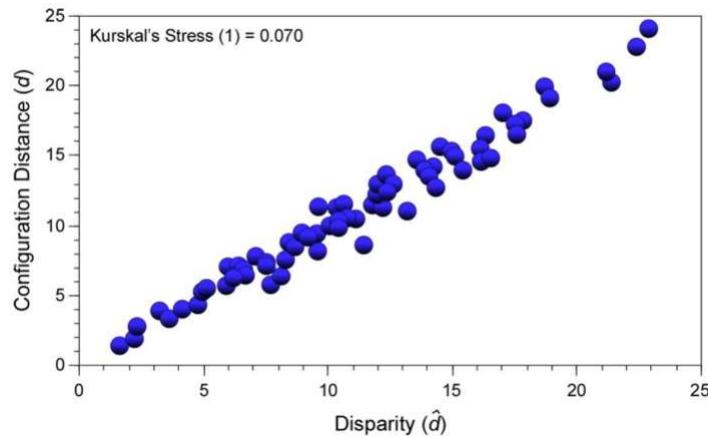
© 2018 Peter V. Henstock

How would you assess these?

© 2018 Peter V. Henstock

How would you assess these?

Shepard Diagram



- http://www.palass.org/modules.php?name=palaeo_math&page=20

© 2018 Peter V. Henstock

Where do these fail?

- PCA is linear technique
 - Linear transform of the points (rotation)
 - Removes noise on hyperplane
- MDS/Sammons focus on pairwise distances in space
- What if data doesn't follow a hyperplane?
 - Curve
 - Manifold
 - ?

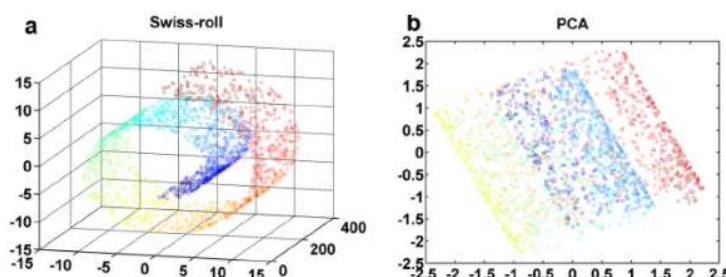
© 2018 Peter V. Henstock

Computational Speed?

- $O(N^2)$
- Typically used for up to 1000s of points
- Computational speed is an issue
- Optimizing 1,000,000 pairwise distances is challenging

© 2018 Peter V. Henstock

Nonlinear manifolds



Luo, Lijia. (2014). Process Monitoring with Global–Local Preserving Projections. Industrial & Engineering Chemistry Research. 53. 7696–7705.

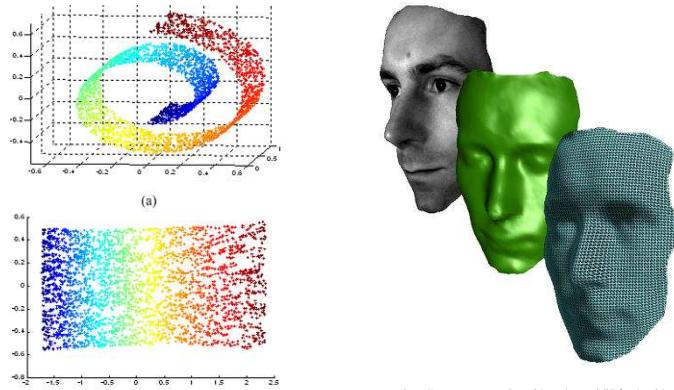


http://scikit-learn.org/stable/auto_examples/manifold/plot_swissroll.html

© 2018 Peter V. Henstock

Isomap

- Tenenbaum, de Silva and Langford
- Science 290 (5500): 2319-2323, 22 Dec. 2000
- What if have nonlinear manifold?



© 2018 Peter V. Henstock

Isomap

- How do you construct distances?

© 2018 Peter V. Henstock

Isomap

- How do you construct distances?
- Same process as the connectivity for the clustering
- Construct neighborhood graph
- Draw edge between all pairs of points with distance less than threshold
- Compute shortest distances along graph between all pairs of points
- Apply MDS

© 2018 Peter V. Henstock

Floyd-Warshall Algorithm

- Computes shortest paths in [weighted] graphs with + and – edges
- Will not work for negative cycles
- $O(V^3)$
- Vertices 1 to N
- Trying to find shortest path from i to j
- Only use k vertices at a point in time
- Goal: find shortest path for $k+1$ vertices

© 2018 Peter V. Henstock

Floyd-Warshall contd.

- For each pair of vertices i and j
- Two possibilities of shortest path
 - 1) Uses vertices 1 to k
 - 2) Using path from $i \rightarrow k+1 \rightarrow j$

Algorithm:

initialize with edge weights

for $k = 1$ to $|V|$

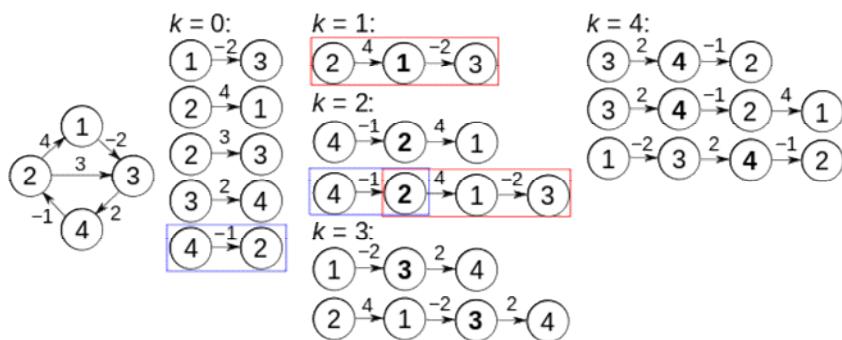
 for each pair i, j

$$\text{dist}(i, j) = \min[\text{dist}(i, j), \text{dist}(i, k) + \text{dist}(k, j)]$$

© 2018 Peter V. Henstock

Example from Wikipedia

- $K = 0$ is original edges
- $K = 1$ is paths through node {1}
- $K = 2$ is path through {1,2}



© 2018 Peter V. Henstock

Local Linear Embedding

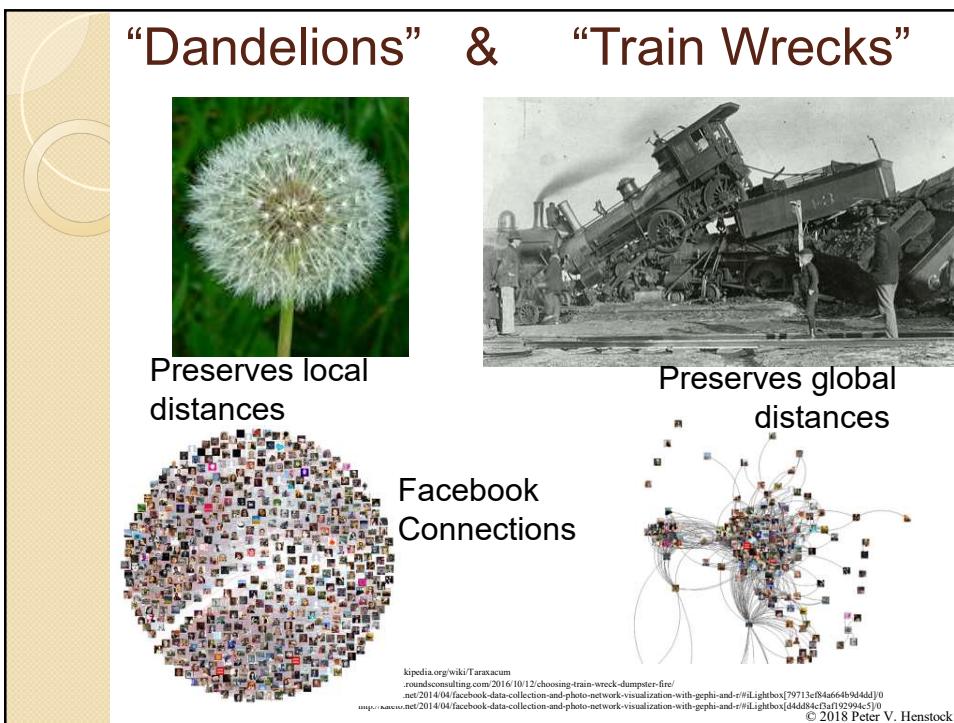
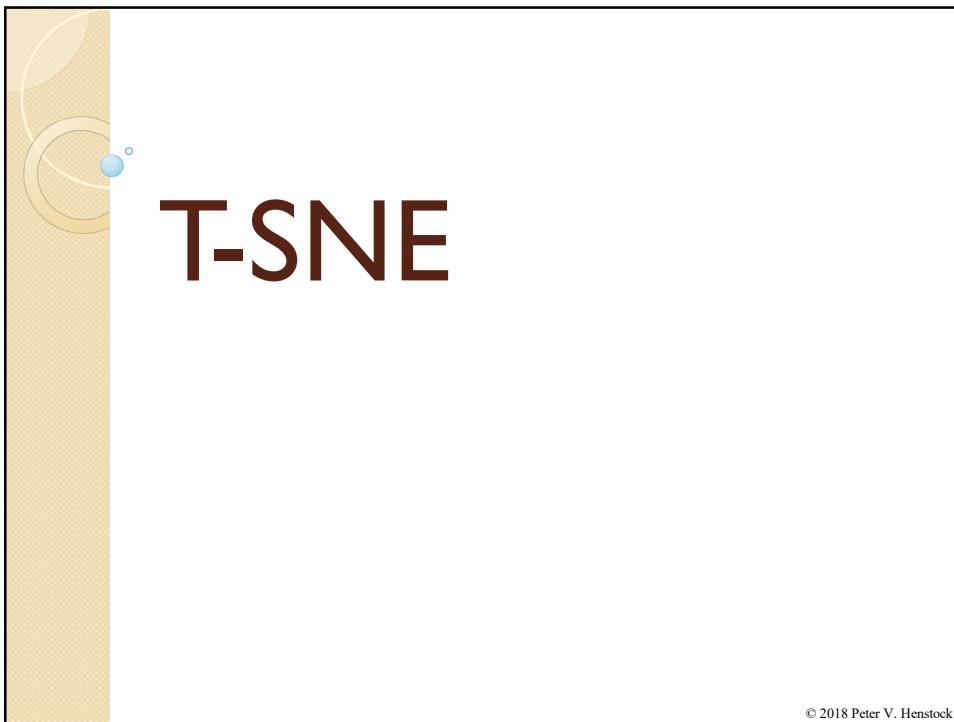
- Nonlinear Dimensionality Reduction by Locally Linear Embedding, Roweis & Saul 2000
- Embedding = mapping to lower dimensional space
- Motivation:
 - Most methods fail on nonlinear spaces
 - Take patches of local geometries together
 - Map patches to plane to preserve global geometry
- Assumptions:
 - Continuous manifold without holes and not noisy data
 - Other versions overcome these
- Approach:
 - Compute KNN of each point X_i ($K=$ only parameter)
 - Minimize $\sum_i [X_i - \sum_j W_{ij} X_j]^2$ $W_{ij} > 0$ for neighbors j ; ensure $\sum_j W_{ij} = 1$
 - Create $\Phi(Y) = \sum_i [Y_i - \sum_j W_{ij} Y_j]^2$ s.t. $\sum_i Y_i = 0$ and $\sum_i Y_i Y_i^T = N I$

© 2018 Peter V. Henstock

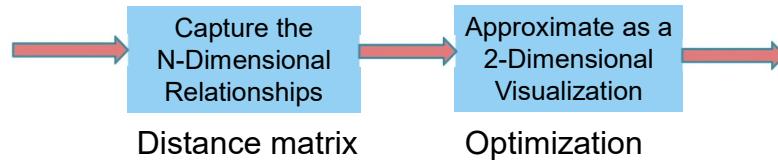
Local Linear Embedding

- Optimization of Y is equivalent to finding eigenvectors
- Space: $N \times K$ weight matrix
- Time: $O(DN^2)$ for D input dimensions
 - $O(DN^2)$ for KNN worst case
 - $O(DNK^3)$ to obtain W weight matrix
 - $O(DN^2)$ for the Y embedding

© 2018 Peter V. Henstock

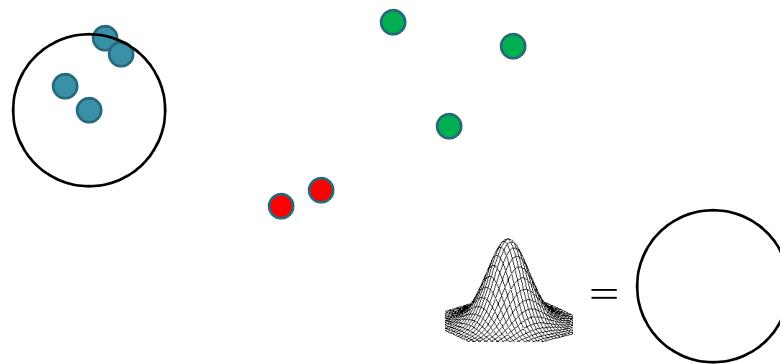


How t-SNE class of solutions works



© 2018 Peter V. Henstock

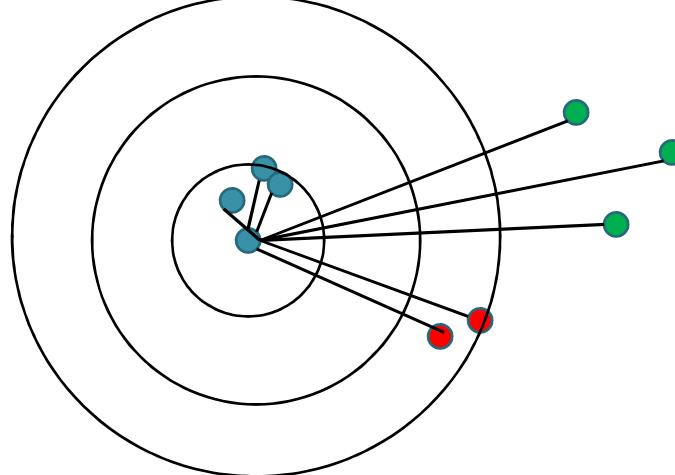
T-SNE Computing output similarity



- Compute square similarity matrix
- Similarities defined on Gaussian
- 1.0 is perfect match and 0.0 is distant

© 2018 Peter V. Henstock

T-SNE: Computing input similarity



- After computing similarities
- Normalize scores to sum of all similarities
- Why? Balance spread on clusters

© 2018 Peter V. Henstock

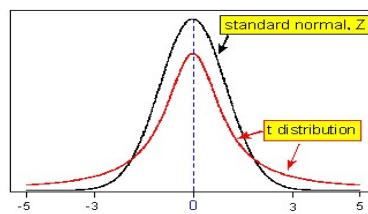
About this Gaussian variance

- Assumed that we had this magical distribution on the blue cluster
- Relates to the perplexity parameter
- Optimal calculation based on neighborhood density of each point
- Are all densities the same? No!
 - $\text{Sim}(\text{blue } j, \text{ red } k) \neq \text{Sim}(\text{red } k, \text{ blue } j)$
 - So...we average them
 - $p(j,k) = \text{Sim}(j,k) / 2 + \text{Sim}(k,j) / 2$

© 2018 Peter V. Henstock

T-SNE Computing output similarity

- Trying to reproduce the distances on the output of course since want similar things to be close
- Setting up the same idea of similarity
- But to compute “q” output distance:
 - Use t-distribution instead of Normal



© 2018 Peter V. Henstock

Now the optimization

- We can compute the input similarity
- We can compute the output similarity
- Now we need something to optimize
- Kullback-Leibler distance
 - $KL(P|Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$
 - KL is based on an entropy based equation
 - Do we want this small or large?
- If p is similar then q should be....
- If p is distant then q should be ...

© 2018 Peter V. Henstock

Gradient Descent

- What is it?

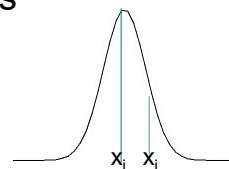
© 2018 Peter V. Henstock

Gradient Descent

- What is it?
- How do we start?

© 2018 Peter V. Henstock

t-SNE

- T-Dist. Stochastic Neighbor Embedding
 - Strategy:
 - Focus on local neighborhoods
 - Use normalized Gaussian probabilities of similarity between local points
- $$p_{ij} = \frac{\exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})}{\sum_k \sum_{l \neq k} \exp(-\frac{\|x_k - x_l\|^2}{2\sigma^2})}$$
- 
- Note that the variance σ^2 is global here but that's not good enough

© 2018 Peter V. Henstock

t-SNE

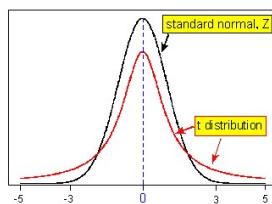
- $p_{ij} = \frac{\exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})}{\sum_k \sum_{l \neq k} \exp(-\frac{\|x_k - x_l\|^2}{2\sigma^2})}$
- Replace global variance with local variance
- $p_{j|i} = \frac{\exp(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2})}{\sum_{k \neq i} \exp(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2})}$
- $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}$ kluge to restore symmetry
- $q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|y_k - y_l\|^2)^{-1}}$ in low dim space
- Gradient descent optimize Kullback-Leibler distance $KL(P||Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$

© 2018 Peter V. Henstock

P & Q have different similarities?

- $p_{j|i} = \frac{\exp(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2})}{\sum_{k \neq i} \exp(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2})}$: Gaussian $f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

- $q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|y_k - y_l\|^2)^{-1}}$: Student-T $\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$



<http://election.princeton.edu/2014/09/18/how-our-predictions-work-continued/>

- Heavy tails allow for distant points to be further apart than will usually happen

© 2018 Peter V. Henstock

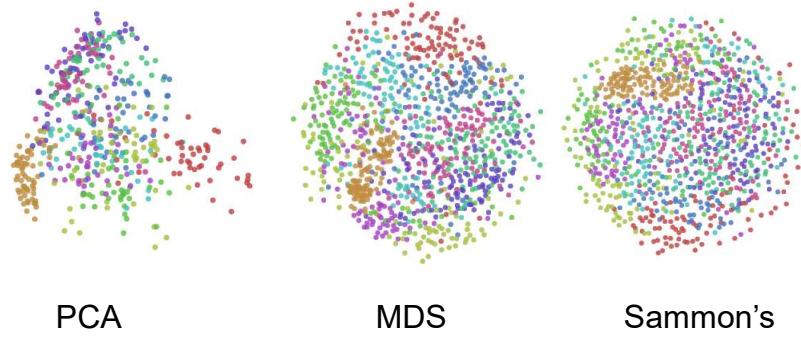
Where is the perplexity parameter?

- One of the more cryptic parts of t-SNE
- It's buried within the σ_i
- Create probability distribution P_i for given σ_i
- Perplexity(P_i) = $2^{-\sum_j p_{j|i} \log_2(p_{j|i})}$
- Exponent part is the entropy of $p_{j|i}$
- t-SNE searches for a σ_i that achieves perplexity requested by the user

© 2018 Peter V. Henstock

Comparison on MNIST

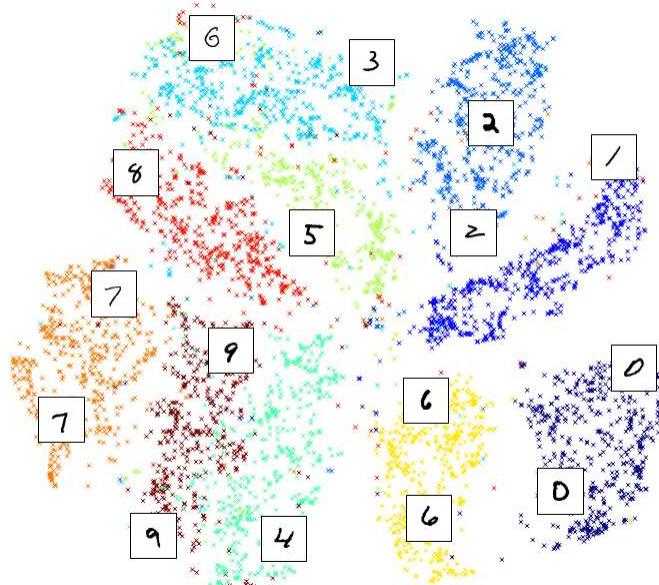
- <http://colah.github.io/posts/2014-10-Visualizing-MNIST/>
- 10K or 60K handwritten numbers



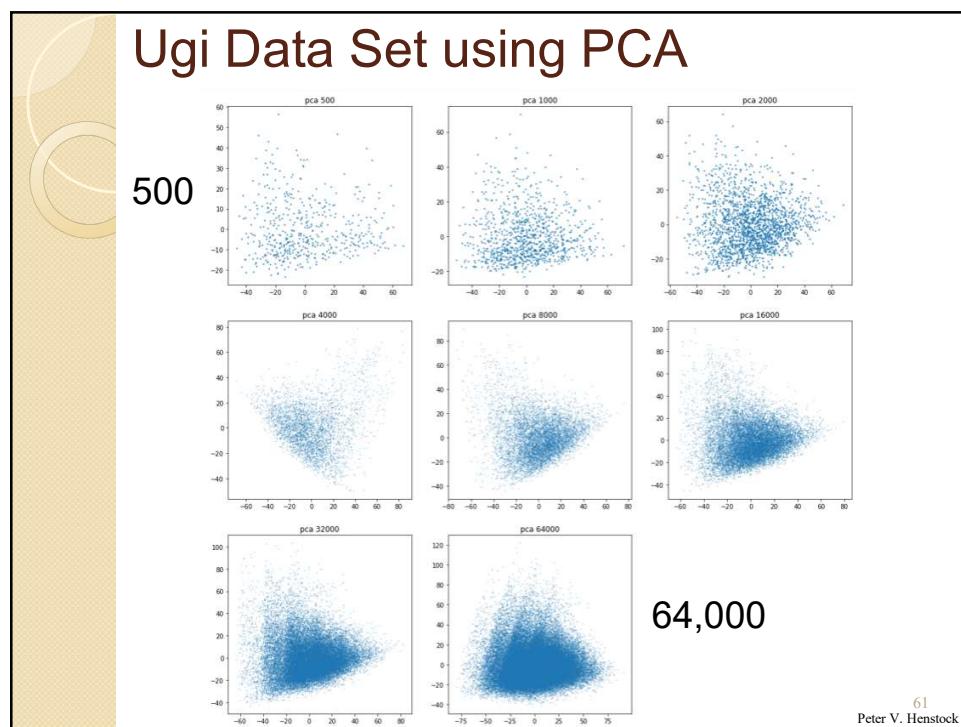
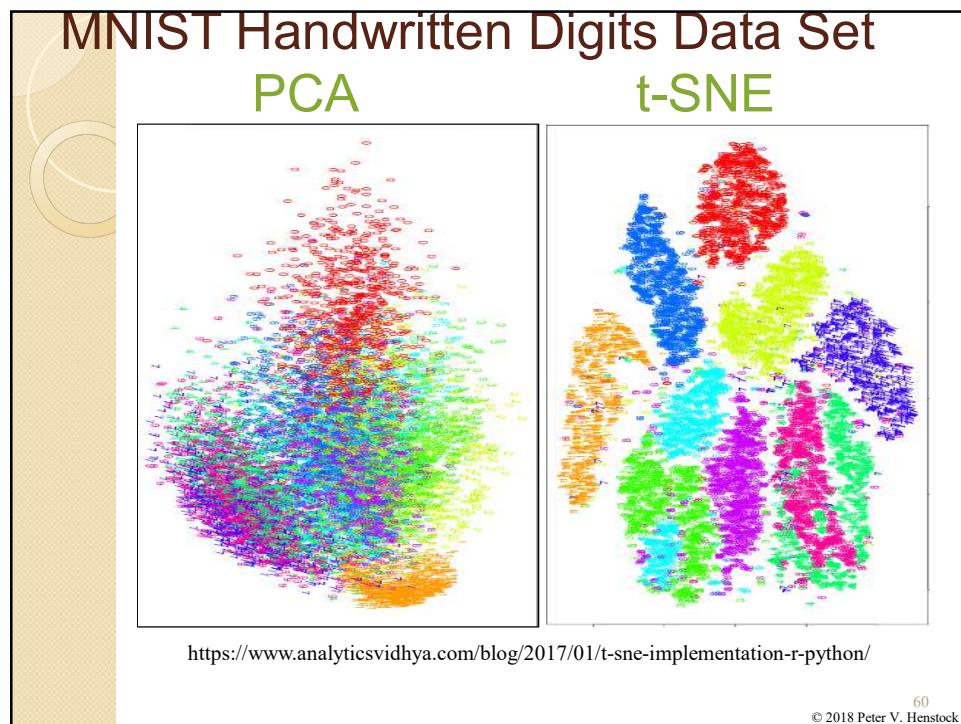
© 2018 Peter V. Henstock

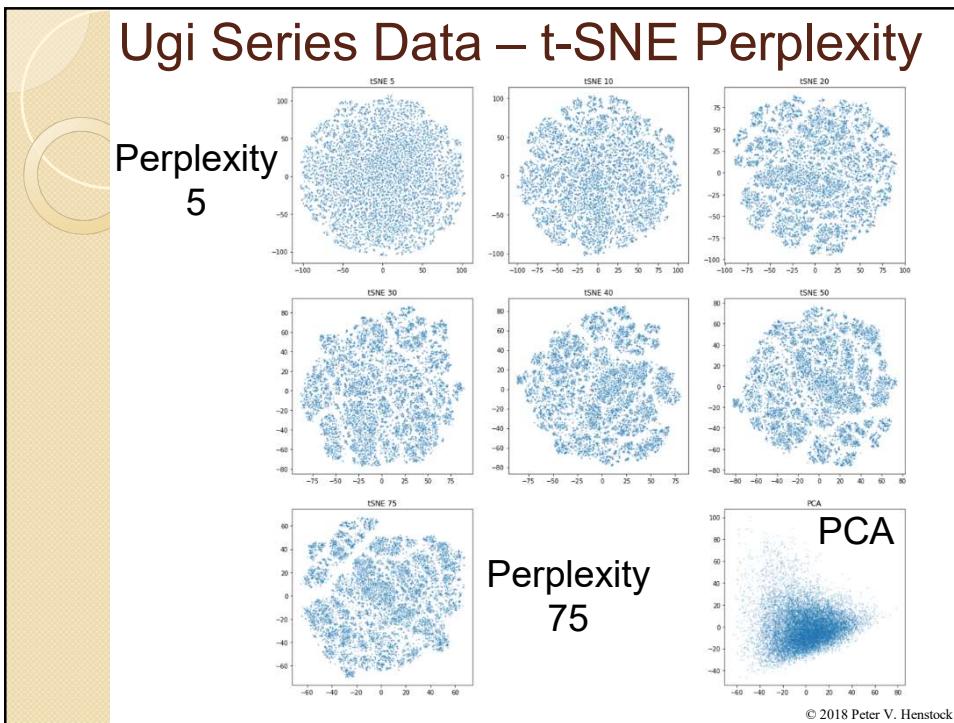
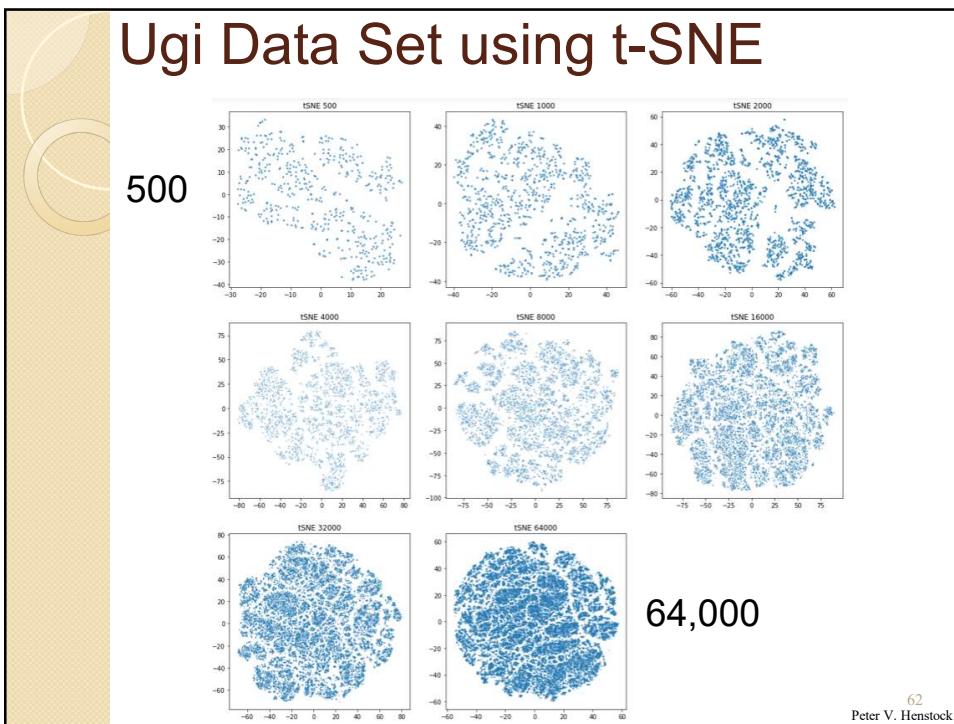
T-SNE on MNIST from

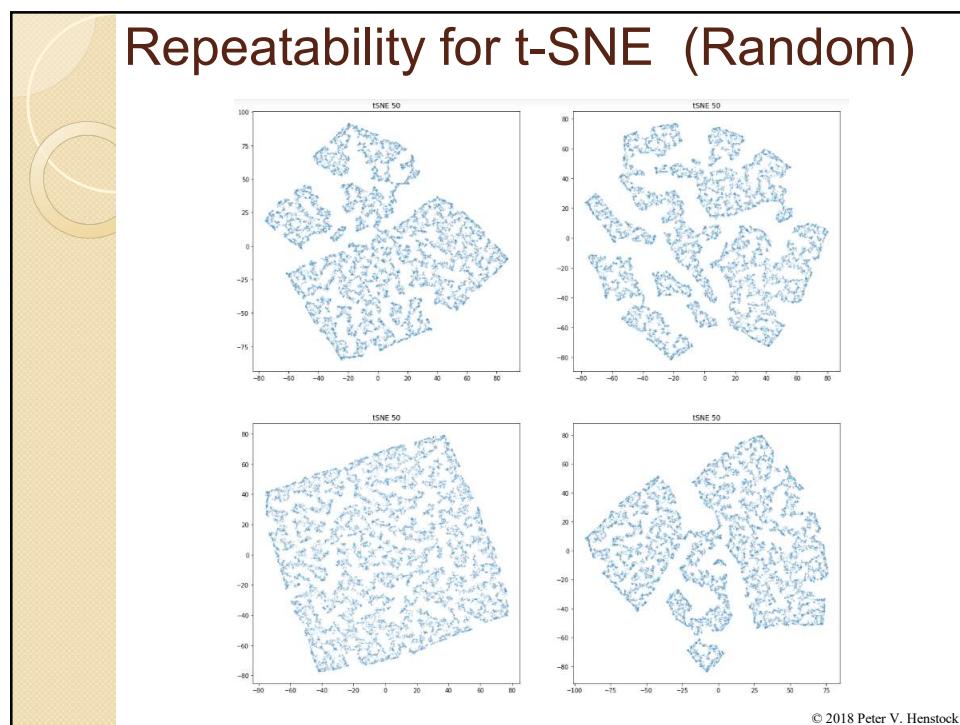
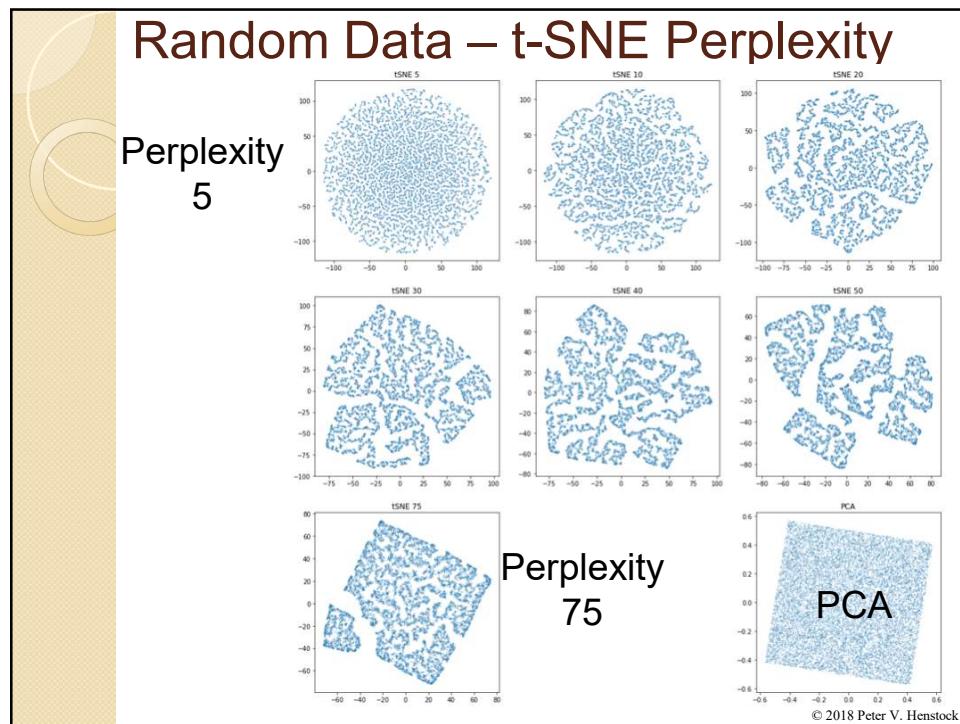
- <http://alexanderfabisch.github.io/t-sne-in-scikit-learn.html>

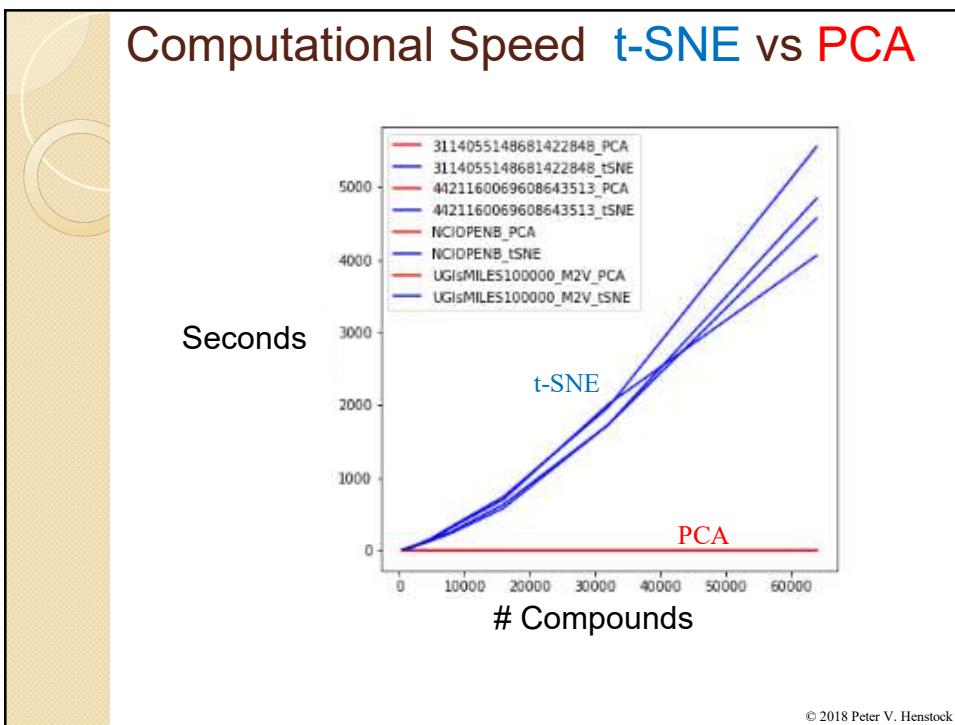
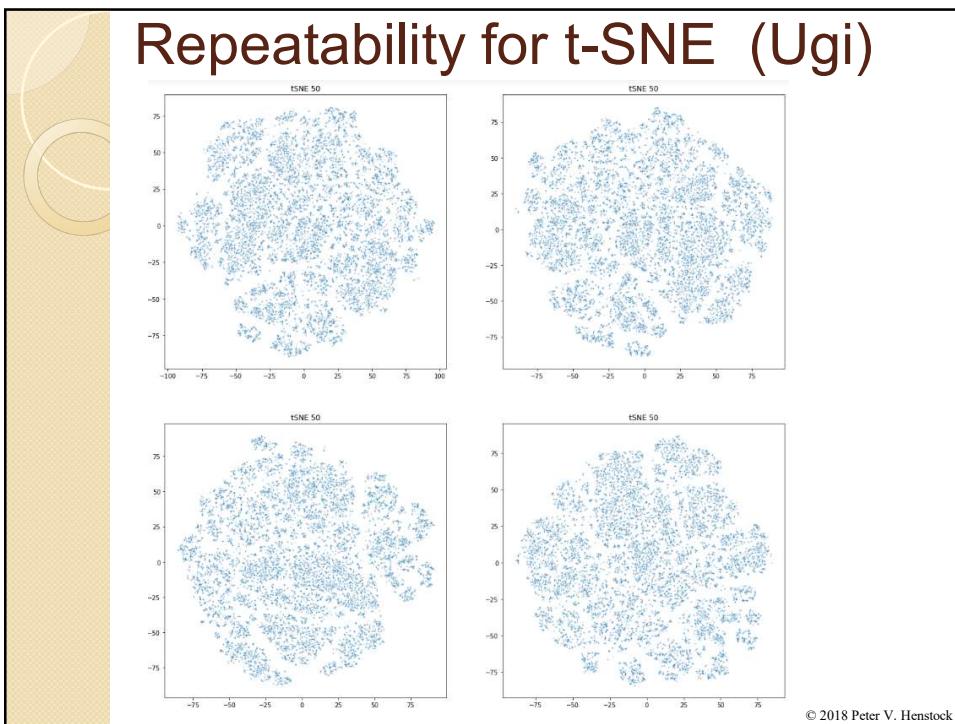


© 2018 Peter V. Henstock

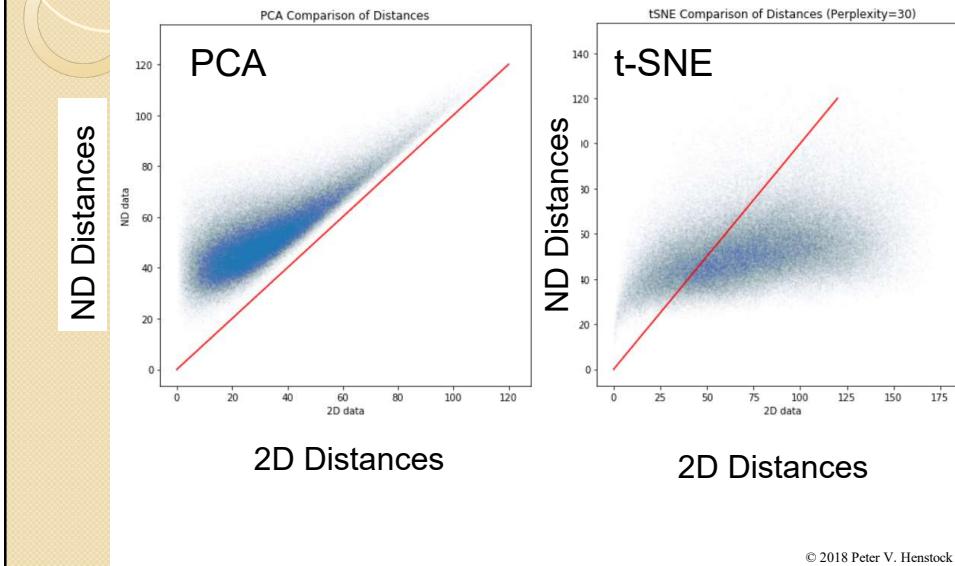




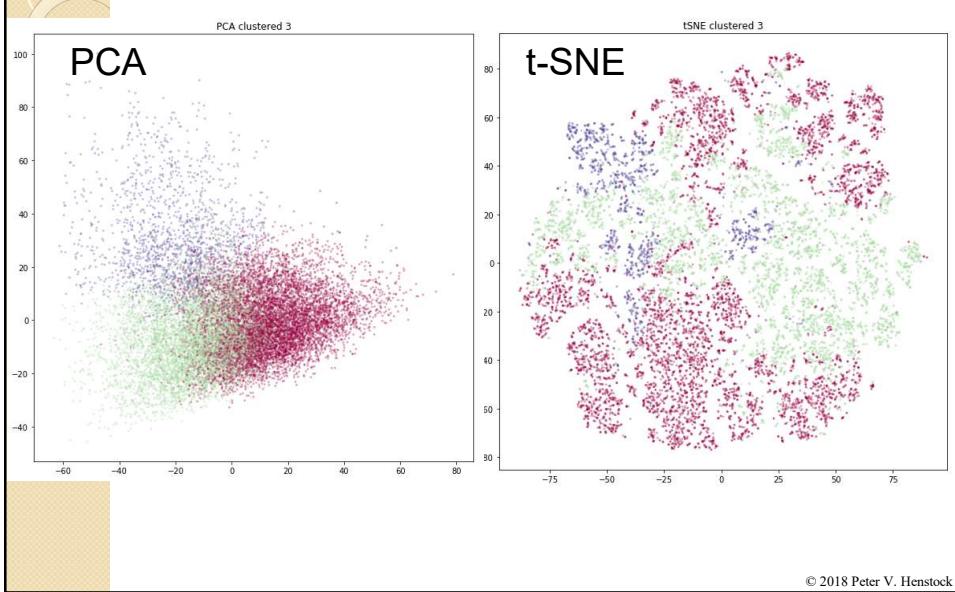


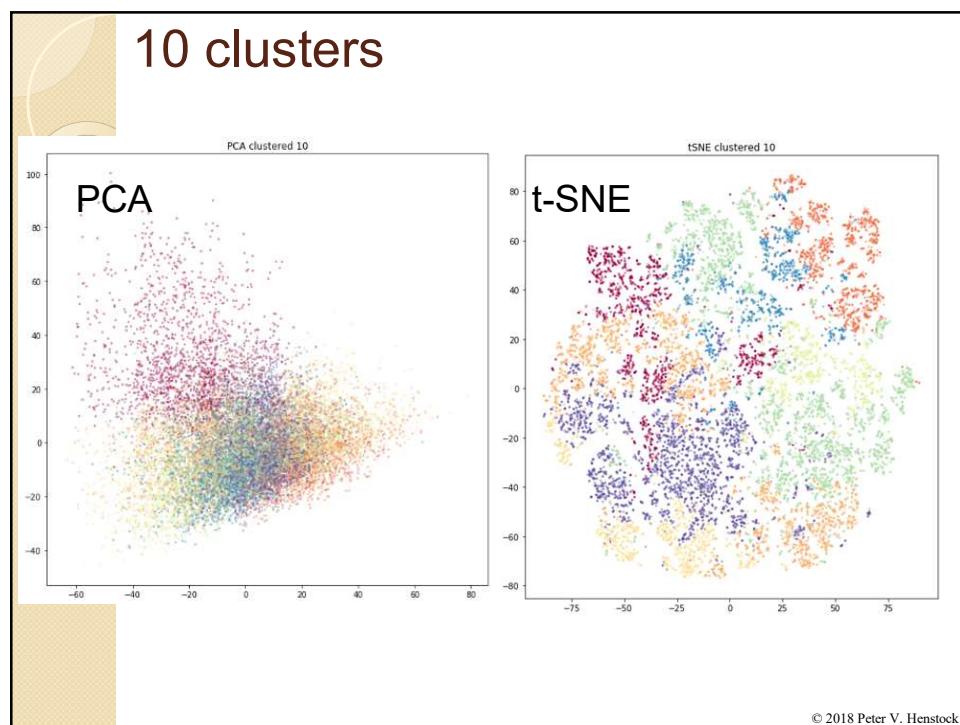
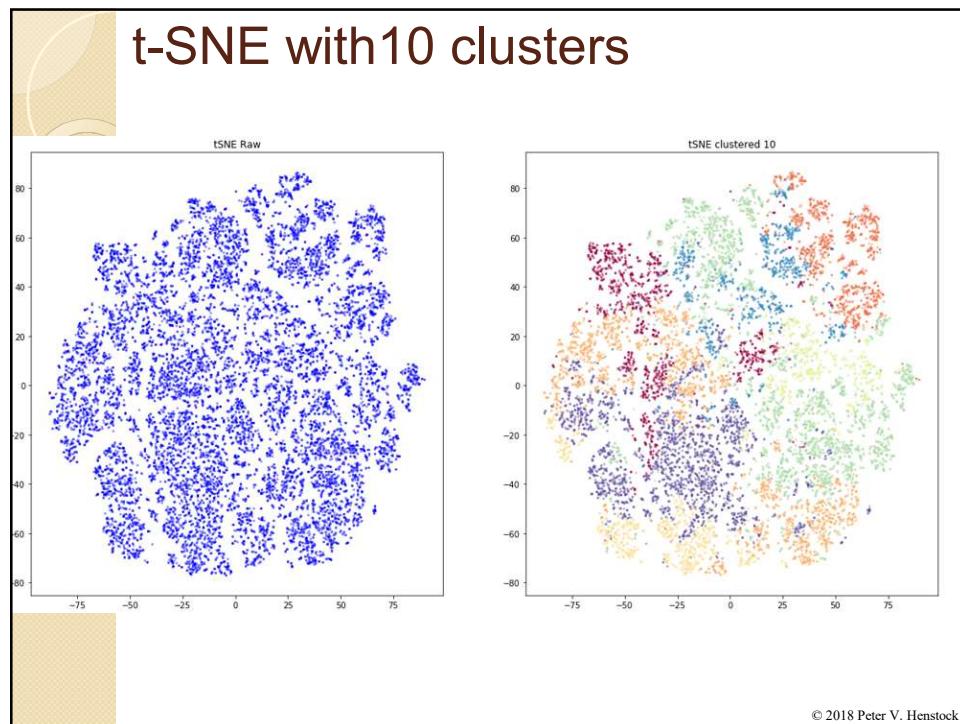


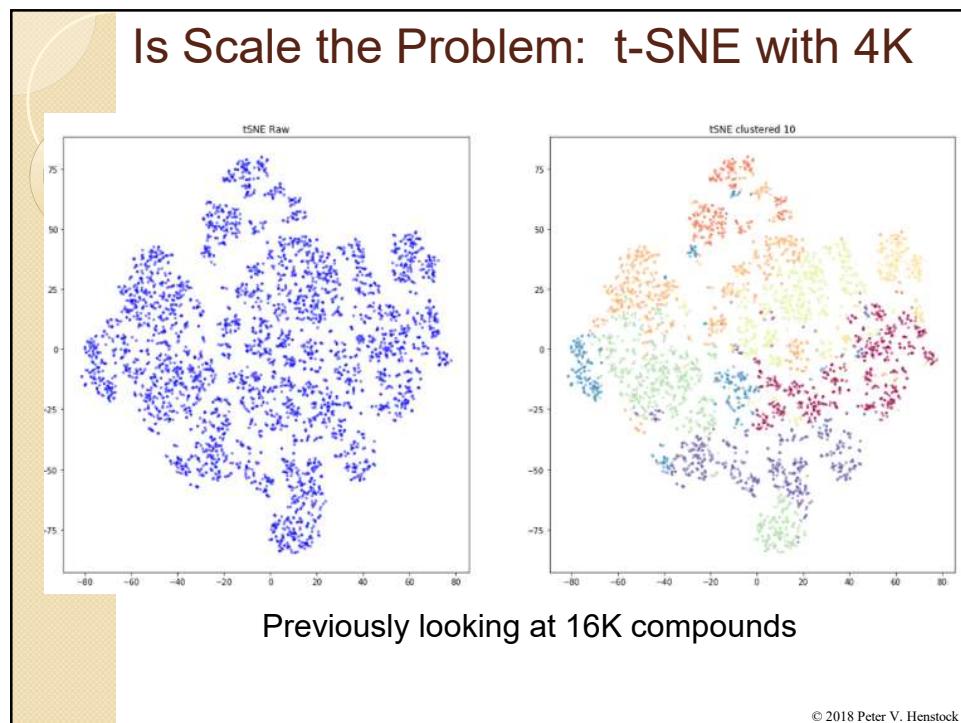
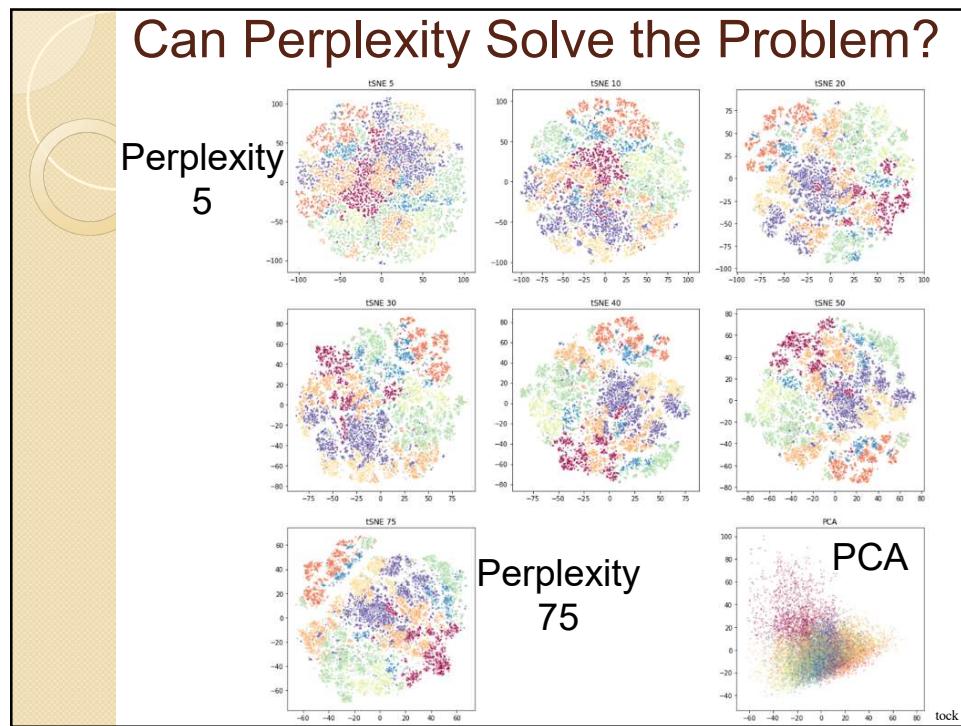
Comparing Distances: ND vs. 2D



Clustering Data with 3 clusters



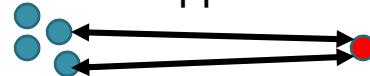




Barnes-Hut

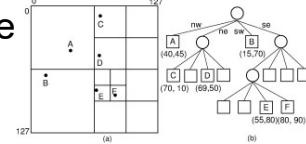
- Computational time for t-SNE
 - Compute the distance matrix so N^2
 - If have 1 million points → zzz

- Barnes-Hut approximation:



- Distance between red and blue are all fairly similar so approximate them together
- Implementation is quadtree

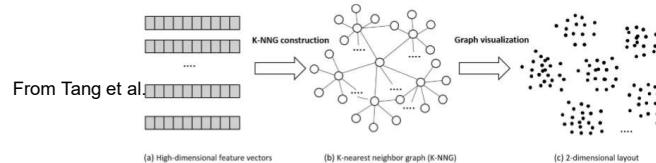
- Run time?



© 2018 Peter V. Henstock

LargeVis

- “Visualizing Large-Scale and High-Dimensional Data” Tang et al. 2016
- Approach is $O(sMN)$ $s=2,3$ $M=\# \text{neg samp}$
 - Computes k-NN graph & maps
 - Layouts graph to lower dimensional space



From Tang et al. Figure 1: A typical pipeline of data visualization by first constructing a K-nearest neighbor graph and then projecting the graph into a low-dimensional space.

- Speed 30x > t-SNE in KNN, 7x for layout
- Better data structures & fewer parameters

© 2018 Peter V. Henstock

Fast KNN graph

- Typical is too slow $O(N^2d)$
- Builds tree by randomly selecting points and dividing up space evenly
- Points in same node: NN candidates
- Standard for high accuracy: many trees
 - LargeVis avoids bottleneck: uses a few trees
 - Takes neighbors of neighbors for KNN
- Uses same t-SNE equation for p_{ij}

© 2018 Peter V. Henstock

Embedding

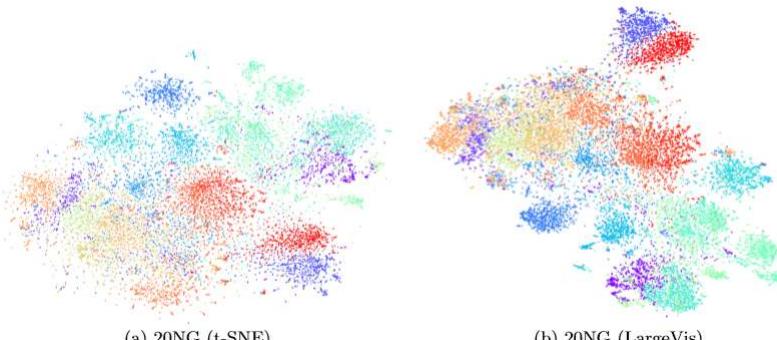
- $P(\text{edge from } i \text{ to } j) = P(e_{ij}=1) = f(||y_i - y_j||)$
 - High when edges are close
- $P(\text{weighted edge } i-j) = P(e_{ij}=w_{ij}) = P(e_{ij}=1)^{w_{ij}}$
- $P(\text{graph}) = \prod P(e_{ij}=1)^{w_{ij}} \prod (1-P(e_{ij}=1))^\gamma$
 - γ = weight for negative (i.e. no) edges
 - $O(N^2)$ operation so avoid this

$$\sum_{ij \in E} w_{ij} \log p(e_{ij} \text{ is 1}) + \sum_{k=1}^M E_{jk \sim P_n(j)} \gamma \log(1 - p(e_{ij} = 1))$$
 - Randomly sample vertices and estimate $P_n(j)$ as negative edges
- Asynchronous gradient descent on samples

© 2018 Peter V. Henstock

T-SNE vs. LargeVis

- 20K documents from 20 different newsgroups
- From LargeVis paper



(a) 20NG (t-SNE)

(b) 20NG (LargeVis)

© 2018 Peter V. Henstock

Good, Bad & Ugly of t-SNE

- Good?
- Bad?
- Ugly?

© 2018 Peter V. Henstock

Good, Bad & Ugly of t-SNE

- Good

- Currently best in separating high dimensional classes
- Gaining traction as a valuable tool

- Bad

- Stochastic part → often can't get same output
- Not a perfect clustering even without Barnes-Hut

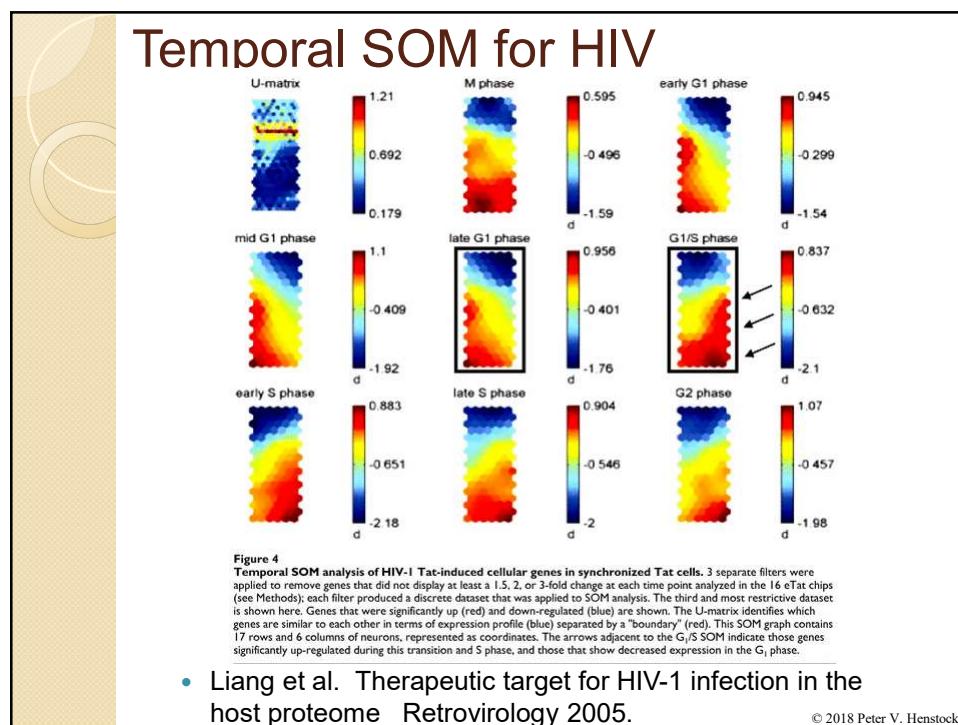
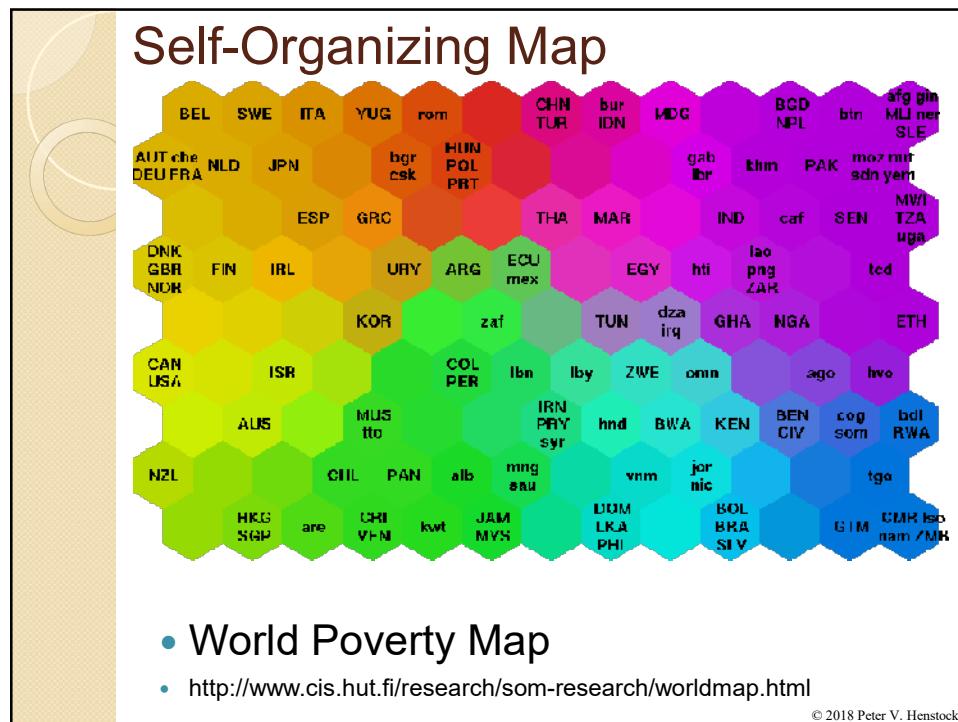
- Ugly

- Perplexity over perplexity: artifacts and tweaking
- Control allowed with balance of local vs. global

© 2018 Peter V. Henstock

An Introduction to Self- Organizing Maps

© 2018 Peter V. Henstock

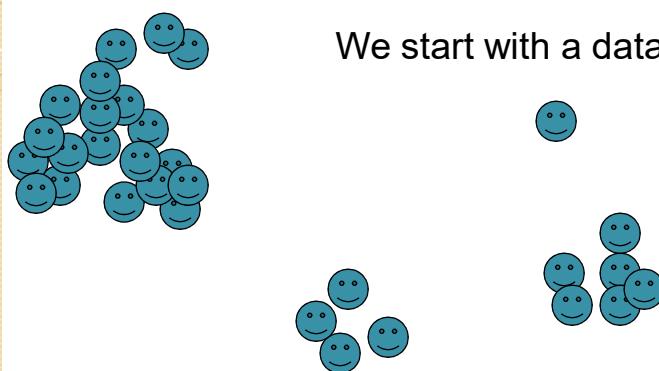


Competitive Learning & Clustering

- Class of clustering algorithms based around neural networks
- Notion is that clusters compete to represent a given data point
- Winning cluster definition is updated or “learns” to better represent that point in the future
- Examples:
 - Adaptive Resonance Theory
 - Grossman @Boston Univ)
 - Self-Organizing Map
 - Kohonen @Helsinki Univ., Finland)

© 2018 Peter V. Henstock

Visual Competitive Learning

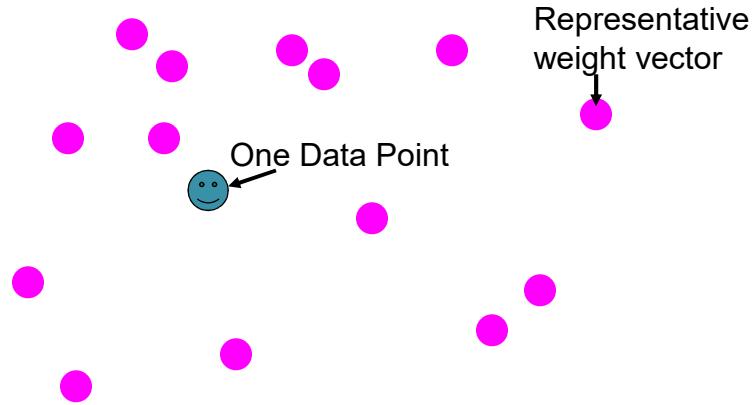


- Each is a data point.
 - Here, they are points in a 2D space.
- We want to find the best representatives.

© 2018 Peter V. Henstock

Visual Competitive Learning

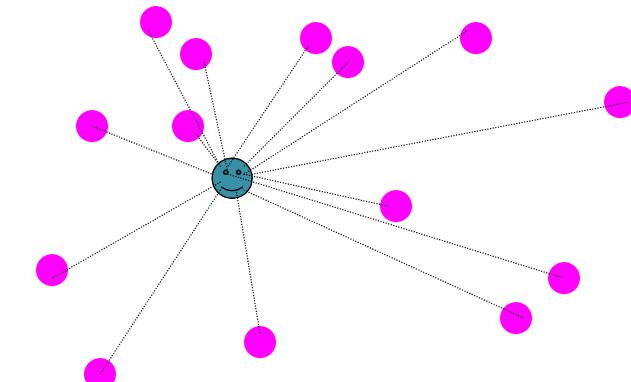
We guess some representatives...



- Each data point is processed sequentially, independent from other data points

© 2018 Peter V. Henstock

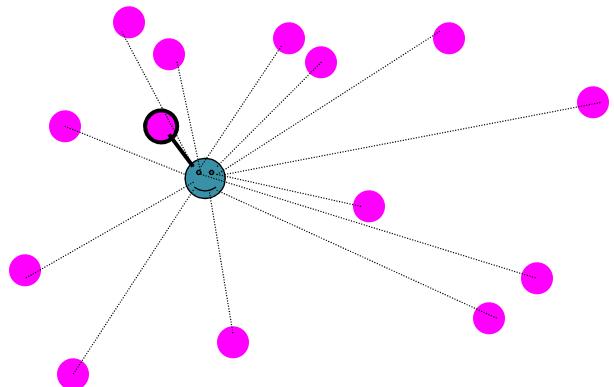
Visual Competitive Learning



- Calculate distances to all weight/reps

© 2018 Peter V. Henstock

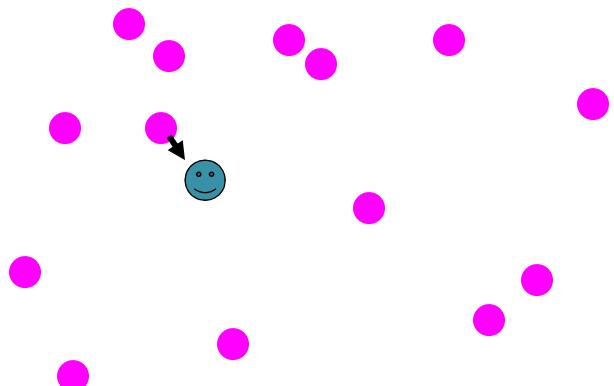
Visual Competitive Learning



- Find the closest weight
 - = best representative

© 2018 Peter V. Henstock

Visual Competitive Learning



- Move weight α distance toward data point

© 2018 Peter V. Henstock

Visual Competitive Learning

- Repeat for next data point
- Repeat for N epochs over set of all data points

© 2018 Peter V. Henstock

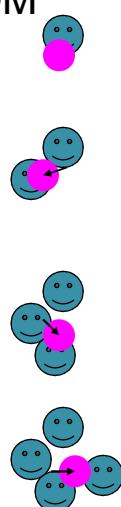
Clustering with Competitive Learning

- Weights migrate to represent clusters
 - More data points—more representatives
- Some weights might be loners

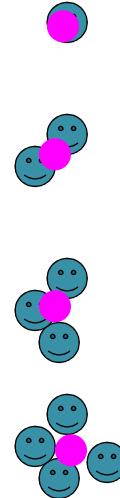
© 2018 Peter V. Henstock

Comparison with K-Means

SOM



K-Means



© 2018 Peter V. Henstock

Maintaining a Topology

- Up until now, we have explained competitive learning as a clustering tool
- Very similar to K-means or other approaches
- Why are SOMs different?
 - Ability to represent N-dimensional data in 2D
 - Similar areas of the high dimensional space map to similar areas of the 2D space

© 2018 Peter V. Henstock

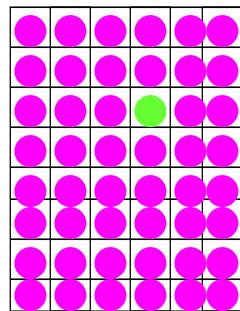
Maintaining a Topology

- Create a 2D array of weight vectors (representatives)
 - Initialize each representative to a random location in the compound space
 - For i=1 to E epochs {
 - For J = 1 to D data points {
 - Find closest representative to dataPoint[J]
 - Move representative a distance α towards dataPoint[J]
 - **Move the K neighbors of the closest representative in 2D array a distance α in direction of dataPoint[J]**
 - }
- Reduce learning rate and neighborhood.

© 2018 Peter V. Henstock

2D SOM Map (often hexagons)

2D map space

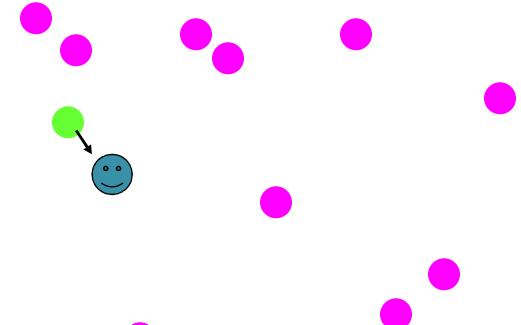


- Representative weights are our 2D visualization of the space
- It's a 2D array of pointers
- Each circle is a pointer to a location in ND space
- The pointer in ND is located at the representative of the cluster
- The representative moves
- The displayed position in the grid does not move...ever

© 2018 Peter V. Henstock

Updating the Neighbors

n-dimensional space of the data set

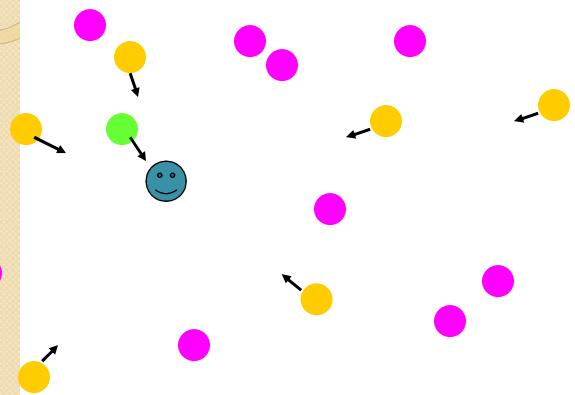


2D map space

© 2018 Peter V. Henstock

Updating the Neighbors

n-dimensional space of the data set

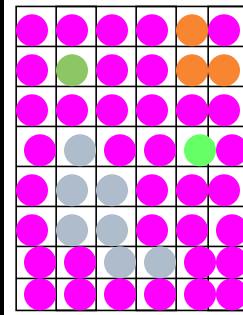
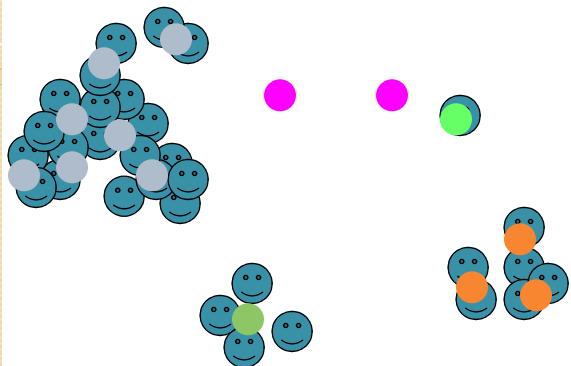


2D map space

- Neighbors of best weight in 2D array learn

© 2018 Peter V. Henstock

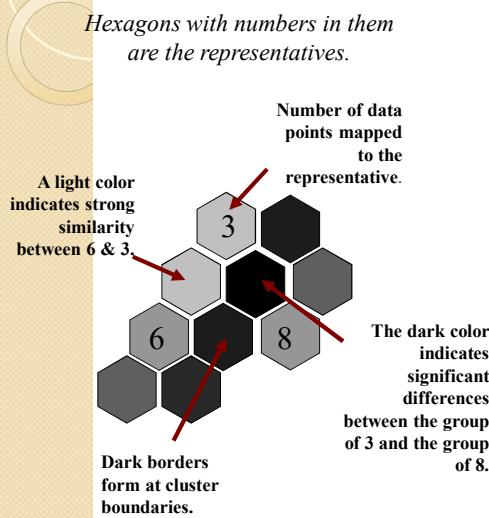
Effect of Preserving Topology



- **Similar representatives map to similar areas**
- Relative distances between clusters will be represented
 - if you are careful about neighborhood parameters
 - We are limited by the reduction of dimensionality

© 2018 Peter V. Henstock

U-Matrix Coloring Algorithm

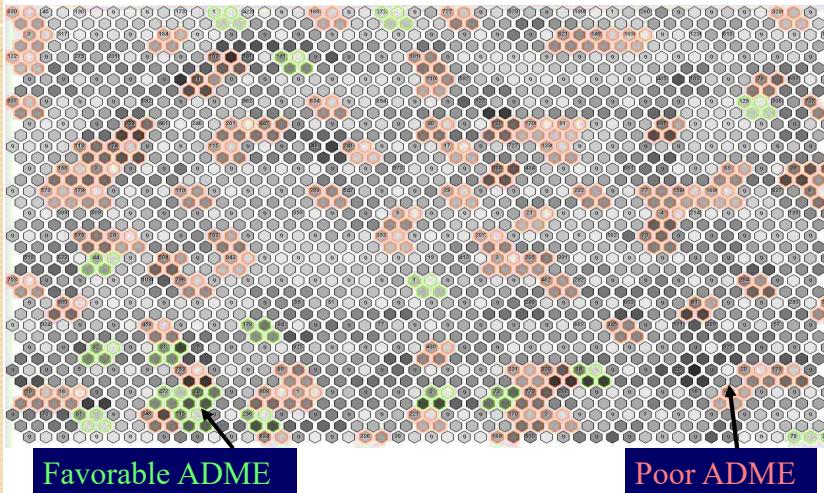


- Hexagons are convenient to represent pairwise distances
- Create histogram of all pairwise distances and map to color table
- Dark = Dissimilar
 - Far distance in n-dimensional space
- Light = Similar
 - Short distance in n-dimensional space

© 2018 Peter V. Henstock

Scalable Self-Organizing Map

- Only known tool that can visualize a large virtual library of 100,000+ compounds and show ADME



©2017 Peter V. Henstock

Good, Bad & Ugly of SOM

- Good?
- Bad?
- Ugly?

© 2018 Peter V. Henstock

Good, Bad & Ugly of SOM

- Good
 - Scalable and accurate with relatively few parameters
 - Provides a visualization along with the clustering
- Bad
 - Difficult to explain to fellow researchers and customers
 - No implicit definition of cluster
 - Sometimes issues with local minima between similar clusters
- Ugly
 - Appreciating the visualization requires time

© 2018 Peter V. Henstock

Text Mining

© 2018 Peter V. Henstock

Records discovery could solve 50-year-old Birmingham civil rights era bombing case

By Joseph D. Bryant | [jbryant@al.com](#)
[Email the author](#) | [Follow on Twitter](#)
 on September 19, 2012 at 7:30 PM, updated September 19, 2012 at 11:39 PM

[Print](#)



BIRMINGHAM, Alabama -- A discovery among Birmingham civil rights era police records and documents could solve a 50-year-old bombing cold case.

Birmingham city records analyst, Bruce Wright, looks over boxes of 50-year-old police reports and documents at City Hall Wednesday. The latest discovery of a package of documents might solve a cold case civil rights era bombing. (Birmingham News Tamika Moore)

Birmingham lawyer Doug Jones, the former U.S. attorney who successfully prosecuted two men in the 1963 Sixteenth Street Baptist Church bombing, said papers stuffed in an envelope and long forgotten in city files might provide details needed to finally close the case.

Jones, in an interview with The Birmingham News this afternoon, would not name the case, but expressed excitement over the finding. Jones was recently called to City Hall to see the historic files.

"The city had just pulled out all these old police records," he said. "I just looked at the file and said, 'We may be able to solve it.'"

http://blog.al.com/spotnews/2012/09/records_discovery_could_solve.html

© 2018 Peter V. Henstock

Enron

- Dot-com energy company
- Executives fudged the accounting
- Bankruptcy



- Email evidence used to convict senior leadership for securities fraud
- Now a public data set

© 2018 Peter V. Henstock

Power of Language

- Horse
- The horse

© 2018 Peter V. Henstock

Power of Language

- Horse
- The horse
- The horse raced.

© 2018 Peter V. Henstock

Power of Language

- Horse
- The horse
- The horse raced.
- The horse raced past.

© 2018 Peter V. Henstock

Power of Language

- Horse
- The horse
- The horse raced.
- The horse raced past.
- The horse raced past the barn.

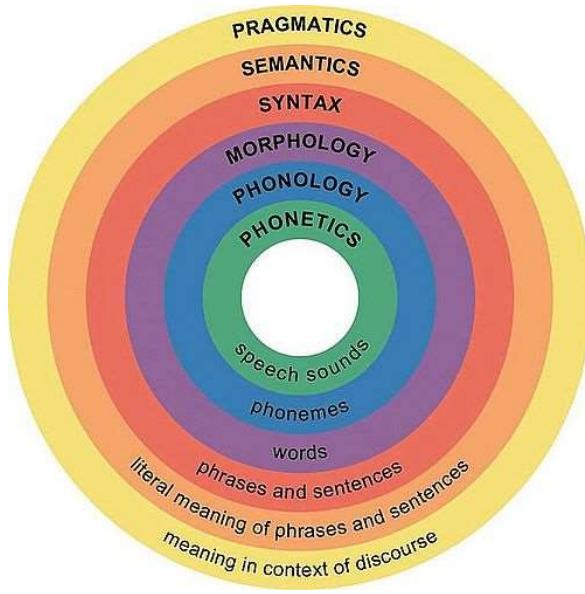
© 2018 Peter V. Henstock

Power of Language

- Horse
- The horse
- The horse raced.
- The horse raced past.
- The horse raced past the barn.
- The horse raced past the barn fell.

© 2018 Peter V. Henstock

Linguistics Hierarchy

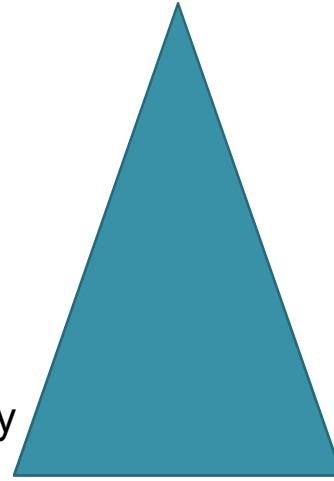


<http://www.sanjaymeena.net/introduction-to-linguistics.html>

© 2018 Peter V. Henstock

Linguistics Pyramid

- Pragmatics
- Sociolinguistics
- Semantics
- Syntax
- Morphology
- Lexical
- Phonology / Orthography



© 2018 Peter V. Henstock

Linguistics Pyramid

- Pragmatics
- Sociolinguistics
- Semantics
- Syntax
- Morphology
- Lexical
- Phonology

© 2018 Peter V. Henstock

Levels of Linguistics

thehorseracedpastthebarn
 The horse raced past the barn.
 article noun verb adverb article noun
 [Noun phrase] [Adverbial phrase]
 [sentence]

The diagram illustrates the levels of linguistics for the sentence "The horse raced past the barn." It shows the breakdown into words, then into Noun phrases ("horse", "barn"), then into an Adverbial phrase ("raced past"), and finally into a sentence. Below the words are corresponding icons: a horse for "animal", a horse race for "movement", and a barn for "location".

animal movement location

Move(agent:animal, where: past barn)

© 2018 Peter V. Henstock

Introduction to Text Mining

The diagram shows a stack of books labeled "Forms" and a stack of emails labeled "Emails". Below them is a stack of colorful books labeled "Tweets". A small figure sits on top of the email stack. An arrow points from the bottom right towards a large red question mark, symbolizing the goal of text mining.

Emails

Forms

Tweets

- <http://www.turnafrownaround.org/volunteer/documents/>
- <http://www.thebooksamaritan.com/>
- http://onclkads.net/?auction_id=74e2bffd191e9427&zoneid=302895&pbk2=6060e265141b149e46c75dc98df186da6201044410040393632&r=%2Foc%2Fhan2Flomb

© 2018 Peter V. Henstock

Introduction to Text Mining

Emails

Forms

Tweets

Features

- <http://www.turnafrownaround.org/volunteer/documents/>
- <http://www.thebooksamaritan.com/>
- http://onclickads.net/?auction_id=74e2bffd191e9427&zoneid=302895&pbk2=6060e265141b149e46c75dc98df186da6201044410040393632&r=%2Foc%2Fham2fomb

© 2018 Peter V. Henstock

What are the rows and columns?

Data feature

Data instance

© 2018 Peter V. Henstock

What are the rows and columns?

	Words					
Sentence						
Paragraph						
Page						
Email						
Act						
Scene						
Book						
Tweet						

© 2018 Peter V. Henstock

Approaches to Text Mining

- Bag of words
- N-gram
- Shallow parsing
- Full parsing

© 2018 Peter V. Henstock

Approaches to Text Mining

- **Bag of words**

- Randomly ordered collection of words
- {hello, my, name, is, peter}

- **N-gram**

- Order within each set of N adjacent units but random N-grams
- {hello my, my name, name is, is peter}

- **Shallow parsing**

- Recognize subset of subject-object-verb or other relationships
 - Subject=Name
 - Verb=is
 - Predicate Nominative=Peter

© 2018 Peter V. Henstock

Full parsing

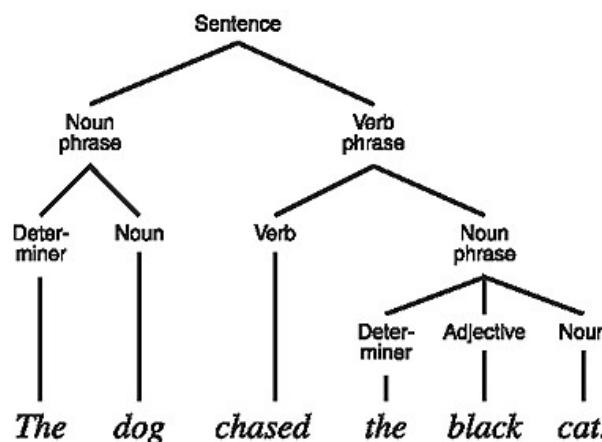
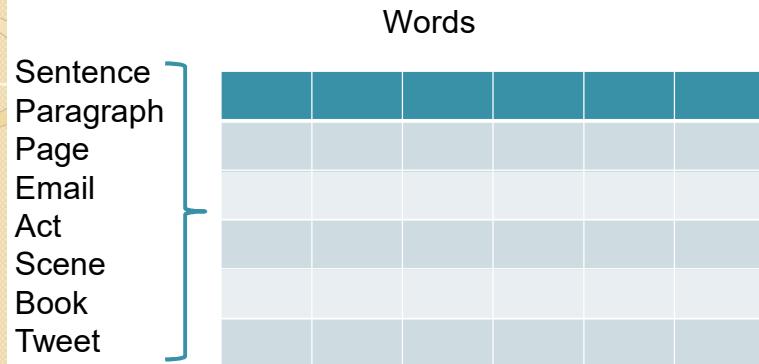


FIGURE 192. Parsing: structure of a sentence

- <http://www.suggestkeyword.com/cGFyc2luZyAgbWVhbmluZw/>

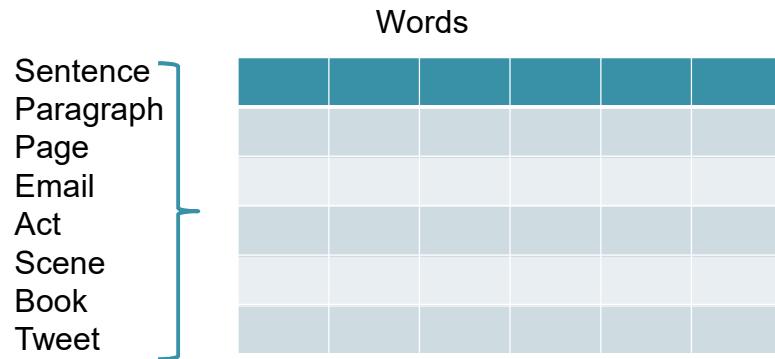
© 2018 Peter V. Henstock

What are the rows and columns?



© 2018 Peter V. Henstock

Thinking about clustering text...



- How would you characterize the distance between two sentences?

© 2018 Peter V. Henstock

Workflow for Bag of Words

- Extract the sentence/paragraph/page
- Remove the punctuation
 - Remove capitalization?
- Extract the words
- Apply [Brill] Part of Speech tagger?
- Remove morphological endings
- Remove stop words

© 2018 Peter V. Henstock

Brill Tagger

- Developed Eric Brill in 1992
- Supervised learning approach to label each word with a part of speech
- Uses a dictionary
- Uses the word structure to guess if it's not in the dictionary
- Set of rules
 - IN NN WDPREVTAG DT while
 - IN=preposition, NN = noun DT=determiner
 - Word "while" should be changed from preposition to noun if previous word tag is DT
 - "in a while" should have while be noun

© 2018 Peter V. Henstock

Stemmers or Lemmatization

- Morphology = modification of a word
- Swim → swims, swimming, swam, etc.
- Peach → peachy, peaches
- Amylase → β -amylase
- Goal: swimming, swims → swim
- Porter Stemmer is one of most common
 - Arguably pretty bad
- Some stemmers require part of speech
 - Noun plurals are good for nouns only

© 2018 Peter V. Henstock

All words are not created equal

- Clustering documents are not likely to be influenced by filler words like “a”, “the”, “in”, etc.
- Useless words are usually removed
- Called “stop words”

© 2018 Peter V. Henstock