## Urban Informatics

Fall 2018

dr. federica bianco fbianco@nyu.edu



@fedhere





Recap:

- Good practices with data: falsifiability, reproducibility
- Basic data retrieving and munging: APIs, Data formats
- SQL
- Basic statistics: distributions and their moments
- Hypothesis testing: *p*-value, statistical significance
- Statistical and Systematic errors
- Goodness of fit tests
- Likelihood
- OLS
- Topics in (time) series analysis
- Visualizations
- Geospatial analysis
- Topics in time series analysis
- Today: Clusters

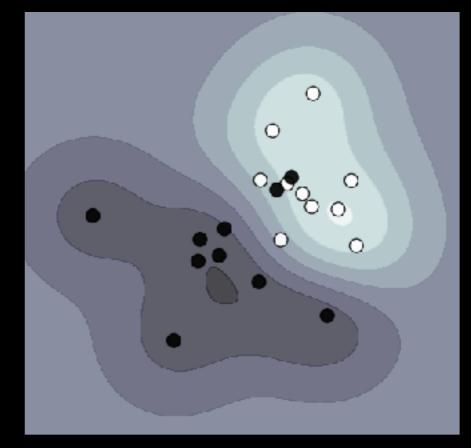
algorithms that can learn from and make predictions on data.





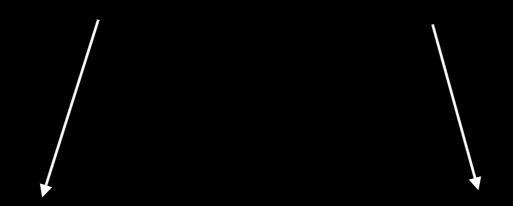
algorithms that can learn from and make predictions on data.

supervised learning extract features and create models that allow prediction where the correct answer is known for a subset of the data





algorithms that can learn from and make predictions on data.



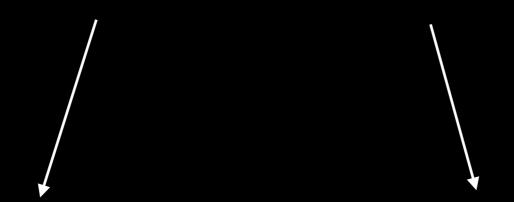
supervised learning extract features and create models that allow prediction where the correct answer is known for a subset of the data

### unsupervised learning

identify features and create models that allow to understand structure in the data



algorithms that can learn from and make predictions on data.



### supervised methods

unsupervised methods

classification

prediction

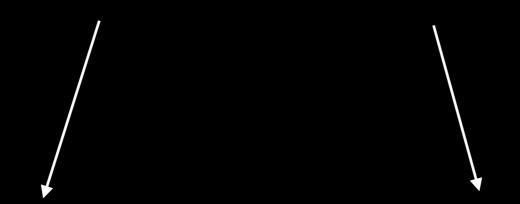
understanding structure

organizing + compressing data

(classification, feature learing)



algorithms that can learn from and make predictions on data.



### supervised methods

unsupervised methods

classification

prediction

understanding structure

organizing + compressing data

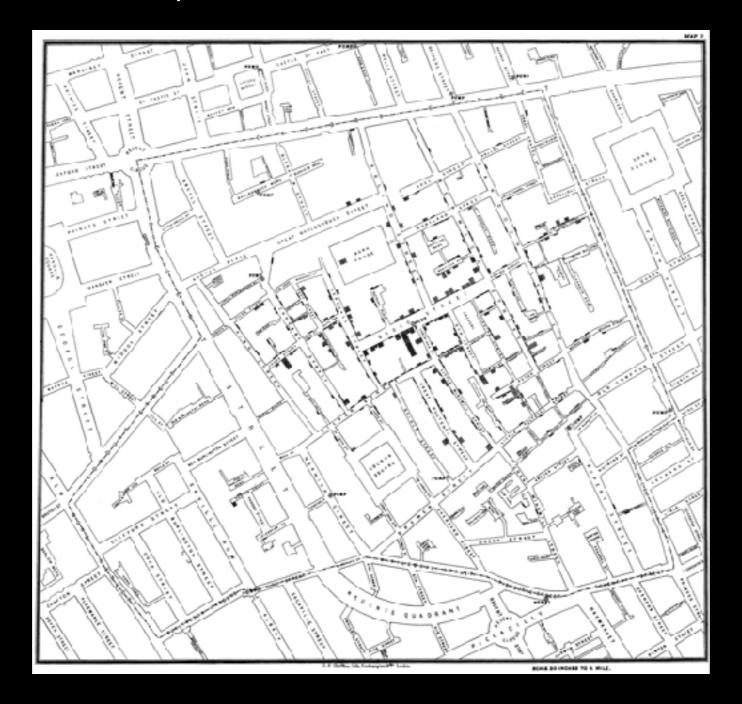
(classification, feature learing)



What is clustering?

XI: Clustering

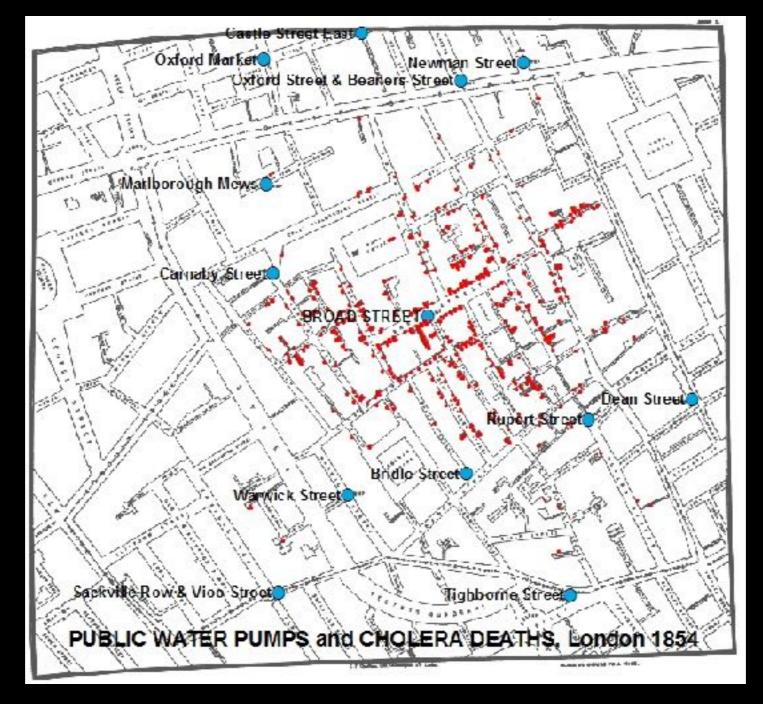
#### Dr. John Snow September 1854 cholera outbreak map



Steven Johnson's 2006 book *The Ghost Map: the Story of London's Most Terrifying Epidemic, and How it Changed Science*.



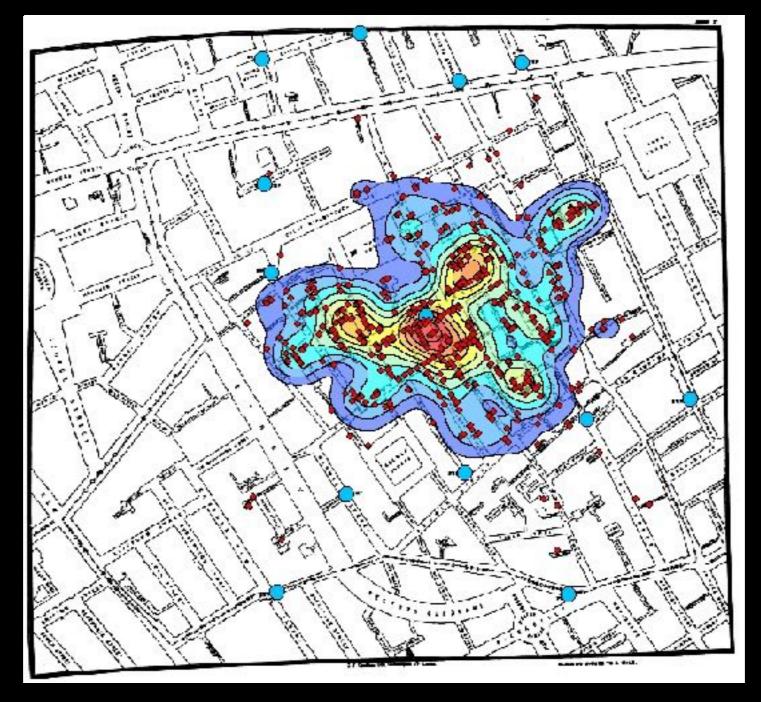
#### Dr. John Snow September 1854 cholera outbreak map



digitization by John Mackenzie, University of Delaware <a href="https://www1.udel.edu/johnmack/frec682/cholera/cholera2.html">https://www1.udel.edu/johnmack/frec682/cholera/cholera2.html</a>

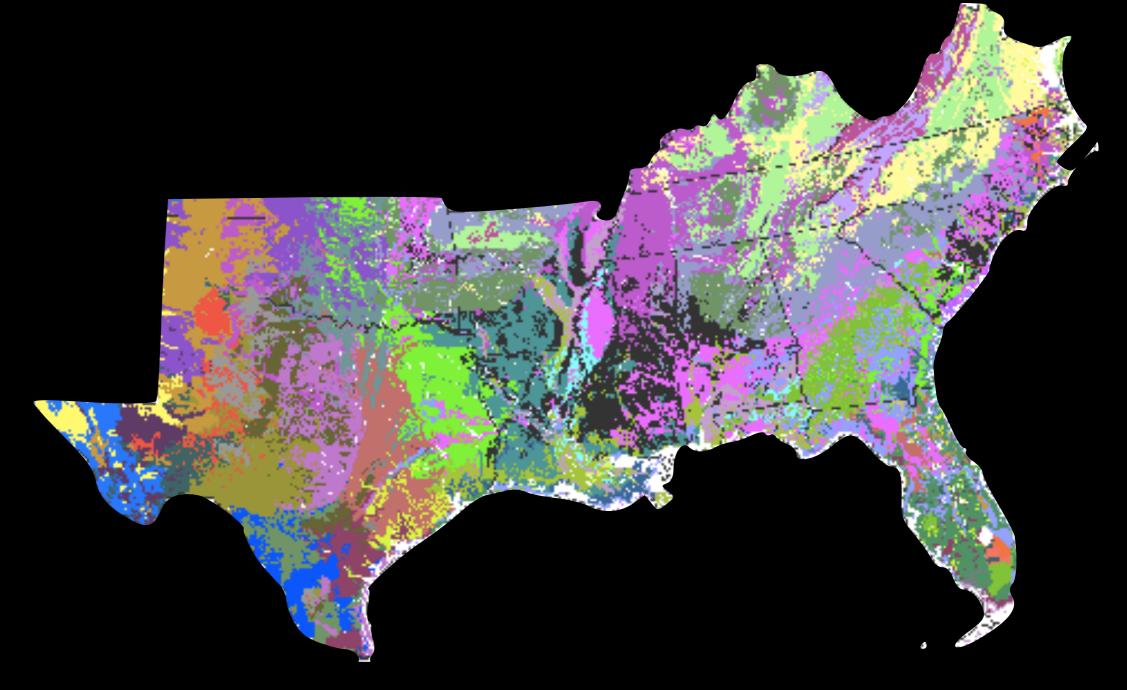


#### Dr. John Snow September 1854 cholera outbreak map

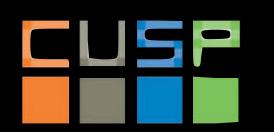


digitization by John Mackenzie, University of Delaware <a href="https://www1.udel.edu/johnmack/frec682/cholera/cholera2.html">https://www1.udel.edu/johnmack/frec682/cholera/cholera2.html</a>



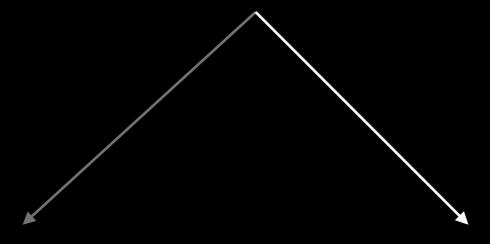


A Spatial Clustering Technique for the Identification of Customizable Ecoregions



William W. Hargrove and Robert J. Luxmoore 50-year mean monthly **temperature**, 50-year mean monthly **precipitation**, **elevation**, total plant-available **water content of soil**, total **organic matter in soil**, and total Kjeldahl **soil nitrogen** 

XI: Clustering



supervised methods

classification prediction

unsupervised methods

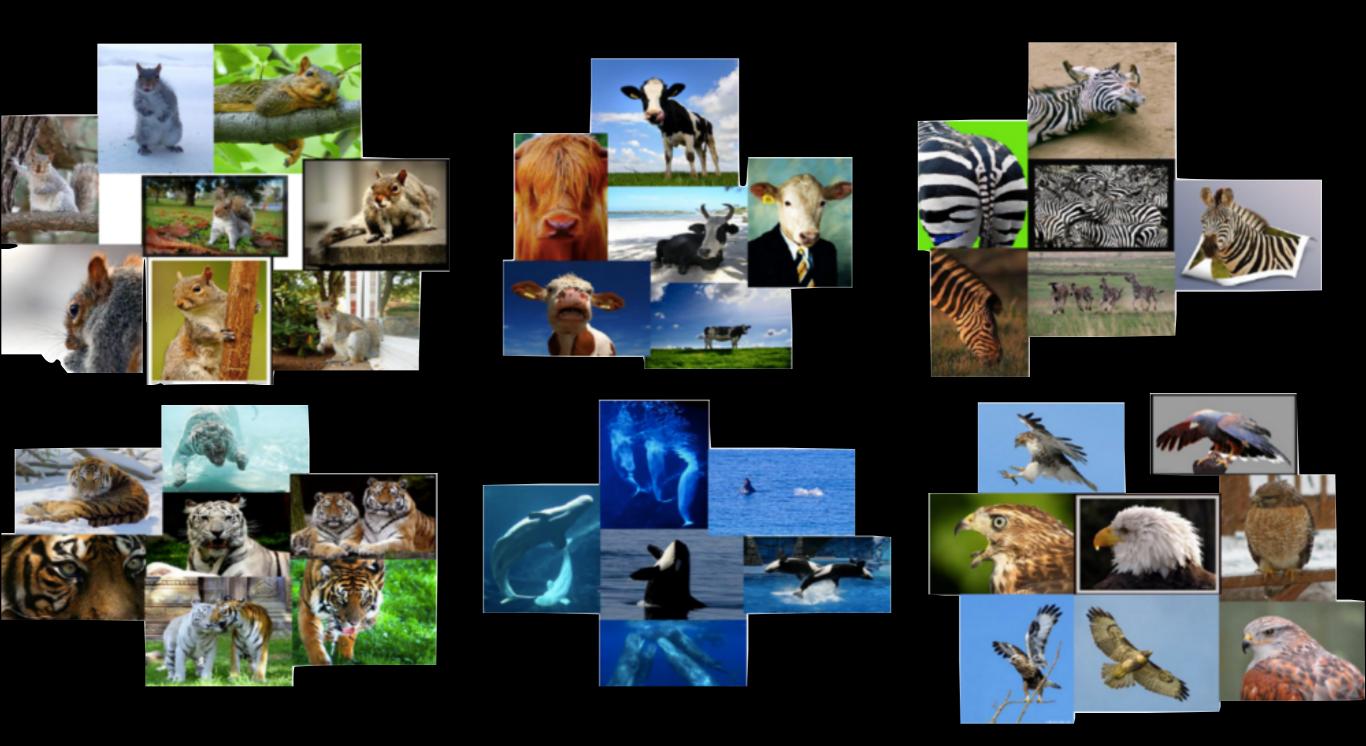
understanding structure

organizing + compressing data

#### **GOAL:**

partitioning data in *maximally homogeneous*, *maximally distinguished* subsets.







#### what is a cluster?

• internal criterion: members of the cluster should be similar to each other (intra cluster compactness)





whales



tigers



eagles

XI: Clustering

#### what is a cluster?

- internal criterion: members of the cluster should be similar to each other
- external criterion: objects outside the cluster should be dissimilar from the objects inside the cluster









#### what is a cluster?

- internal criterion: members of the cluster should be similar to each other
- external criterion: objects outside the cluster should be dissimilar from the objects inside the cluster







https://github.com/fedhere/UInotebooks/blob/master/cluster/ imageProcessingKmeans.ipynb

https://github.com/fedhere/UInotebooks/blob/master/cluster/cluster/

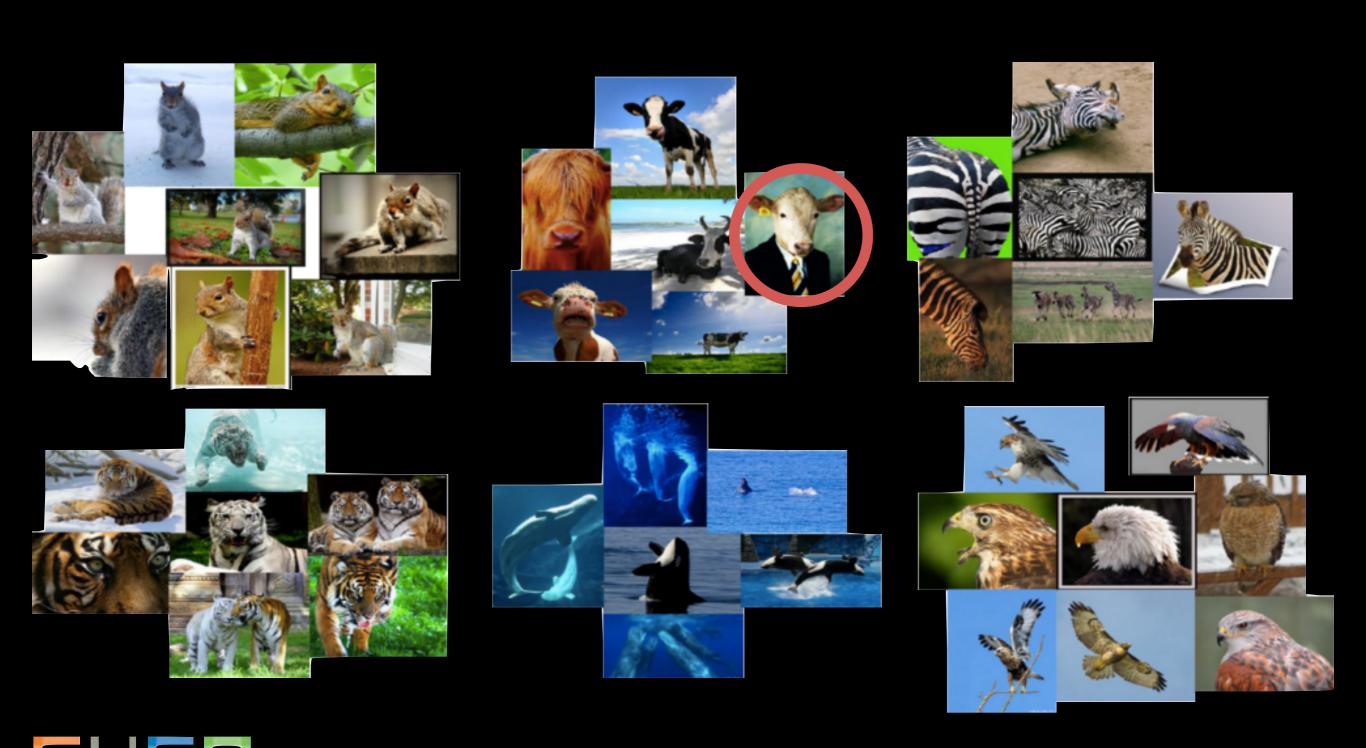


https://github.com/fedhere/UInotebooks/blob/master/cluster/hardVSsoftClustering.ipynb

#### The ideal clustering algorithm:

- Scalability (naive algorithms are Np hard)
- Ability to deal with different types of attributes
- Discovery of clusters with arbitrary shapes
- Minimal requirement for domain knowledge
- Deals with noise and outliers
- Insensitive to order
- Allows incorporation of constraints
- Interpretable





### Defining the distance



# Distance Metrics Continuous variables Minkowski family of distances

$$D(i,j) = p \sqrt{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{iN} - x_{jN}|^p}$$



#### Minkowski family of distances

$$D(i,j) = p / \sum_{k=1}^{N} |x_{ik} - x_{jk}|^p$$
 N features (dimensions)



#### Minkowski family of distances

$$D(i,j) = p \int_{k=1}^{N} |x_{ik} - x_{jk}|^p$$
 N features (dimensions)

$$D(i,i) = 0$$

$$D(i,j) = D(j,i)$$

$$D(i,j) <= D(i,k) + D(k,j)$$



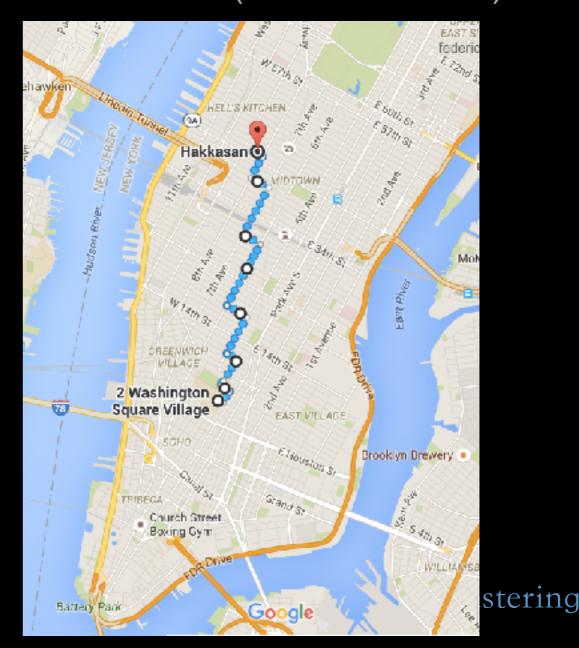
#### Minkowski family of distances

$$D(i,j) = p \left| \sum_{k=1}^{N} |x_{ik} - x_{jk}|^p \right|$$

Manhattan: p = 1

$$D_{Man}(i,j) = \sum_{k=1}^{N} |x_{ik} - x_{jk}|$$

N features (dimensions)



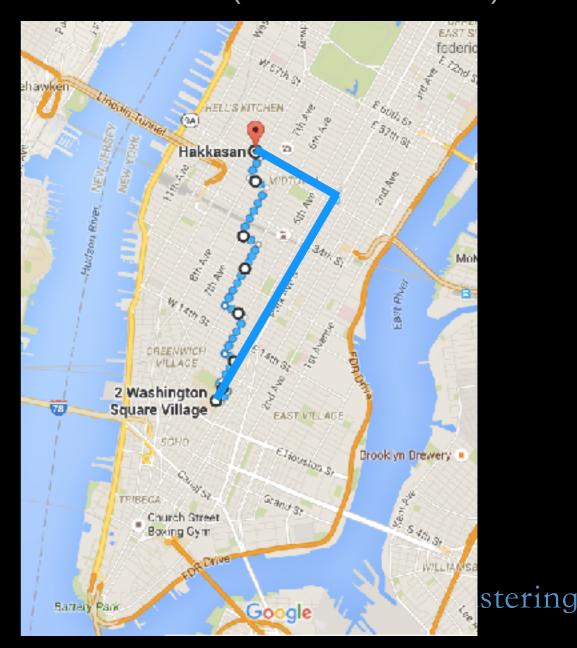
#### Minkowski family of distances

$$D(i,j) = p \left| \sum_{k=1}^{N} |x_{ik} - x_{jk}|^p \right|$$

Manhattan: p = 1

$$D_{Man}(i,j) = \sum_{k=1}^{N} |x_{ik} - x_{jk}|$$

N features (dimensions)



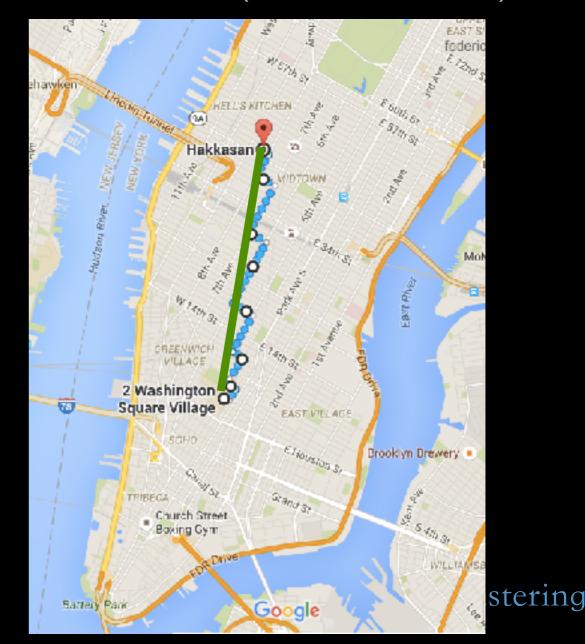
#### Minkowski family of distances

$$D(i,j) = P \int_{k=1}^{N} |x_{ik} - x_{jk}|^p$$

Euclidean: p = 2

$$D_{Euc}(i,j) = \sqrt{\sum_{k=1}^{N} |x_{ik} - x_{jk}|^2}$$

N features (dimensions)



#### Minkowski family of distances

$$D(i,j) = P \left| \sum_{k=1}^{N} |x_{ik} - x_{jk}|^p \right|$$



Great Circle distances:  $\phi_i, \lambda_i, \phi_j, \lambda_j$ 

geographical latitude and longitude

$$D(i,j) = R \arccos(\sin\phi_i \cdot \sin\phi_j + \cos\phi_i \cdot \cos\phi_j \cdot \cos(\Delta\lambda))$$



#### Minkowski family of distances

$$D(i,j) = \sqrt[1/p]{\sum_{k=1}^{N} |x_{ik} - x_{jk}|^p}$$
 N features (dimensions)

Weighted distances:

$$D(i,j) = {}^{p}\sqrt{w_{1}|x_{ik}-x_{jk}|^{p}+w_{2}|x_{i2}-x_{j2}|^{p}+...+w_{N}|x_{iN}-x_{jN}|^{p}}$$

XI: Clustering

#### **Distance Metrics** Binary variables

Uses presence/absence data in two samples

Simple similarity coefficient *SMC* 

$$S_{ij} = \frac{M_{i=0j=0} + M_{i=1j=1}}{M_{00} + M_{01} + M_{10} + M_{11}}$$



#### **Distance Metrics** Binary variables

	1	0	sum	
1	а	b	a+b	
0	С	d	c+d	
sum	a+c	b+d	p	

Uses presence/absence data in two samples

## Simple similarity coefficient *SMC*

$$S_{ij} = \frac{b+c}{a+b+c+d}$$

a = number of items in common,

b = number of items unique to the first set

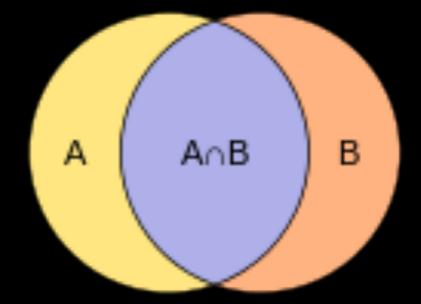


#### **Distance Metrics** Sets variables

Uses presence/absence data

## Jaccard similarity coefficient $S_i$

$$S_j = \frac{a}{a+b+c}$$



a = number of items in common,

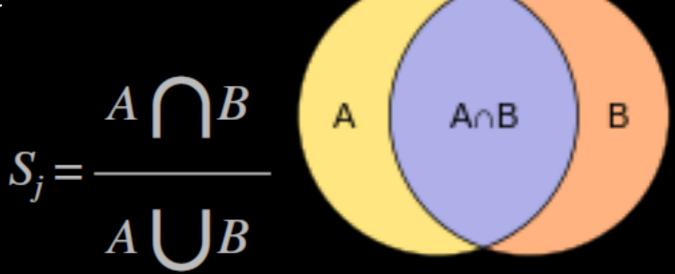
b = number of items unique to the first set



#### **Distance Metrics** Sets variables

Uses presence/absence data

## Jaccard similarity coefficient $S_i$



a = number of items in common,

b = number of items unique to the first set

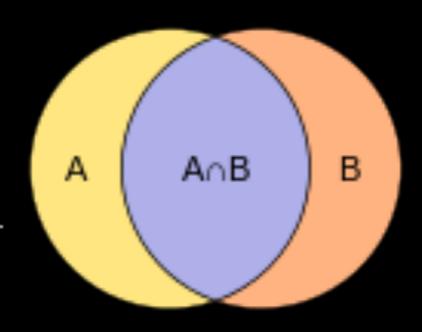


#### **Distance Metrics** Sets variables

Uses presence/absence data

Jaccard distance  $D_i = 1 - S_i$ 

$$S_{j} = \frac{A \bigcap B}{A \mid A \mid B}$$



a = number of items in common,

b = number of items unique to the first set





https://github.com/fedhere/Ulnotebooks/blob/master/cluster/

XI: Clustering



https://github.com/fedhere/Ulnotebooks/blob/master/cluster/

XI: Clustering

## How clustering works



#### **Clustering methods**

Partitioning

#### Hard clustering

K-means (McQueen '67) K-medoids (Kaufman & Rausseeuw '87)

## Soft Clustering Expectation Maximization (Dempster, Laird, Rubin '77)

 Hirarchical agglomerative

devisive

Density based :



DBSCAN (Ester, Kriegel, Sander, Xu'96) A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise (11,000+ citations)

## Clustering methods Partitioning

#### Hard clustering

K-means (McQueen '67) K-medoids (Kaufman & Rausseeuw '87)

Soft Clustering Expectation Maximization (Dempster, Laird, Rubin '77)

Hirarchical agglomerative
 devisive

Density based :



DBSCAN (Ester, Kriegel, Sander, Xu'96) A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise (11,000+ citations)

#### K-means:

- 1. Choose N "centers" guesses: random points in the data space
- 2. Calculate which center each datapoint is closest to: these are the N clusters
- 3. Calculate the new centers as means of the assigned clusters: these are the new N centers
- 4. Iterate 2&3 till convergence: when clusters no longer change



#### K-means:

#### Minimizes the intra cluster variance

Order: #clusters #dimensions #iterations #datapoints O(KdN)

works on minimizing the aggregate distance within the cluster if the distance is Euclidean this is the same amminimizing the variance

Its non-deterministic: the result depends on the (random) starting point

It only works where the mean is defined: alternative is K-medoids which represents the cluster by its central member, rather than by the mean



Must declare the number of clusters upfront

XI: Clustering



https://github.com/fedhere/Ulnotebooks/blob/master/cluster/

XI: Clustering



https://github.com/fedhere/Ulnotebooks/blob/master/cluster/
LUSP kmeans.ipynb

XI: Clustering

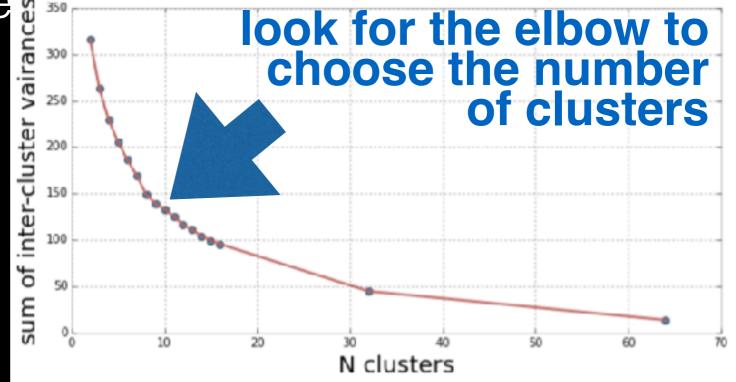
#### K-means:

#### Minimizes the intra cluster variance

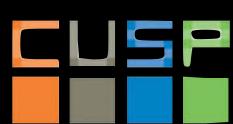
Order: #clusters #dimensions #iterations #datapoints O(KdN)

works on minimizing the aggregate distance within the cluster if the distance is Euclidean this is the same am

minimizing the variance is minimized in the distance is minimized in the d







#### **Clustering methods**

Partitioning

#### Hard clustering

K-means (McQueen '67) K-medoids (Kaufman & Rausseeuw '87)

Soft Clustering Expectation Maximization (Dempster, Laird, Rubin '77)

 Hirarchical agglomerative

devisive

Density based :



DBSCAN (Ester, Kriegel, Sander, Xu'96) A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise (11,000+ citations)

XI: Clustering

#### **Clustering methods**

Partitioning

#### Hard clustering

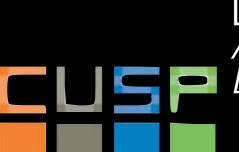
K-means (McQueen '67) K-medoids (Kaufman & Rausseeuw '87)

Soft Clustering Expectation Maximization (Dempster, Laird, Rubin '77)

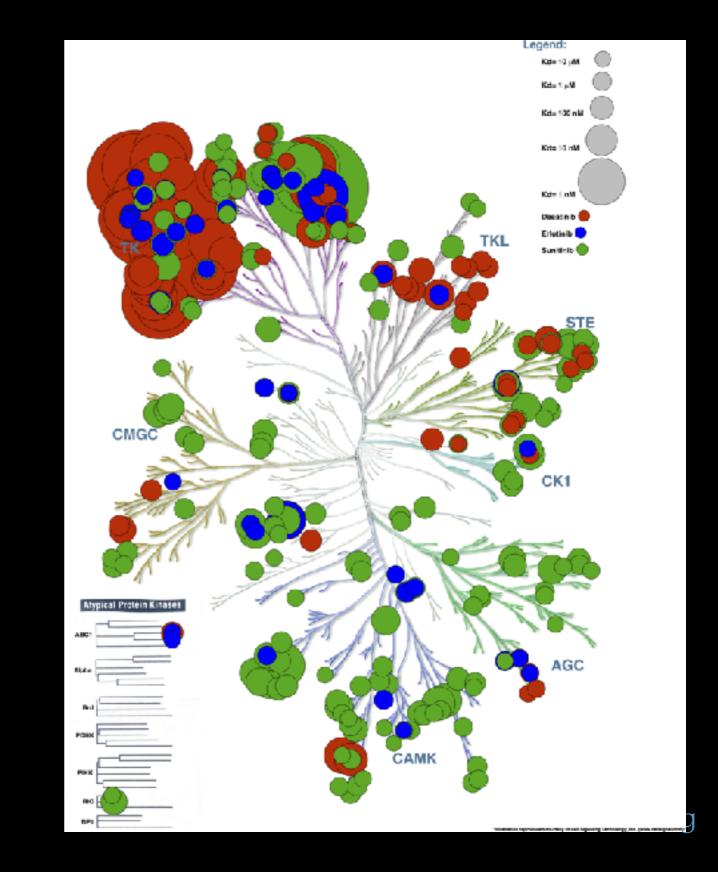
 Hirarchical agglomerative

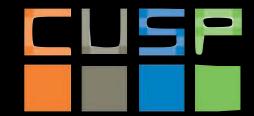
devisive

Density based :

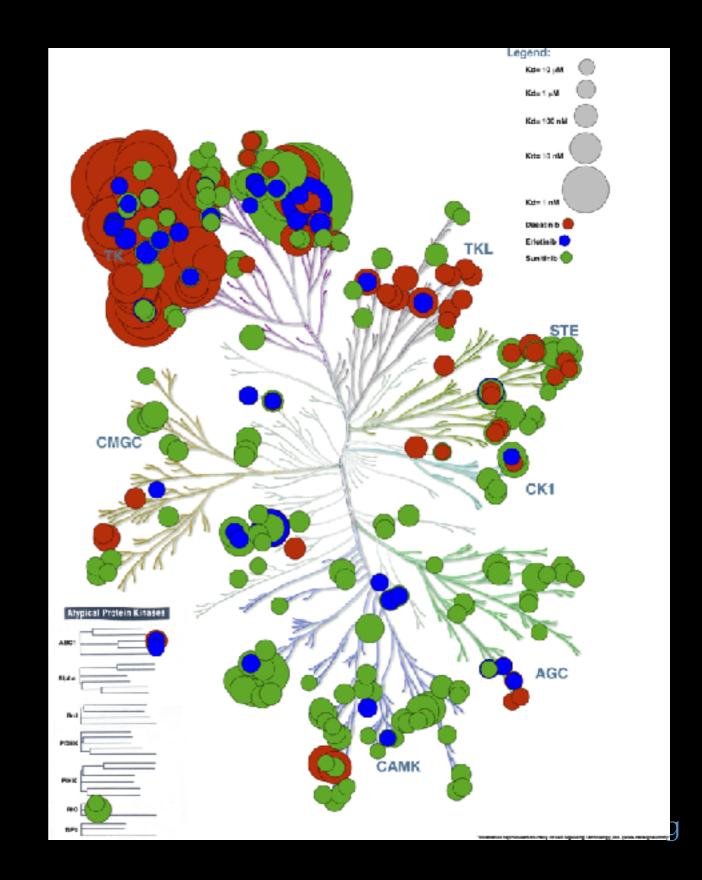


DBSCAN (Ester, Kriegel, Sander, Xu'96) A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise (11,000+ citations)



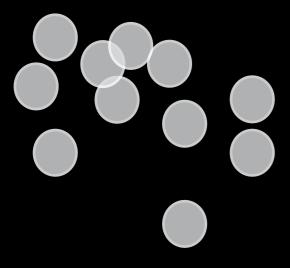


removes the issue of deciding K (number of clusters)



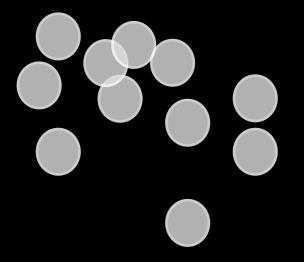


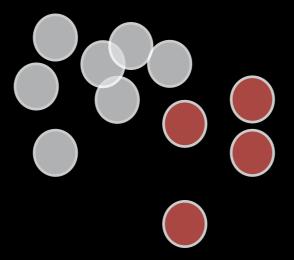
devisive (top-down): e.g. hierarchical k-mean





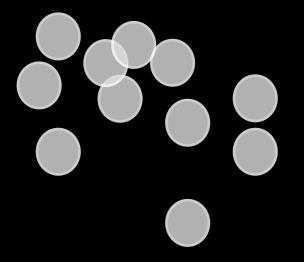
devisive (top-down): e.g. hierarchical k-mean

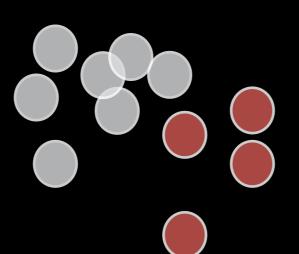


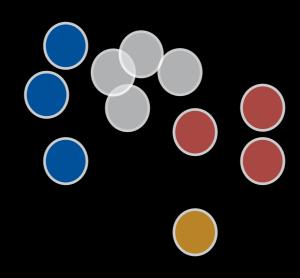




devisive (top-down): e.g. hierarchical k-mean



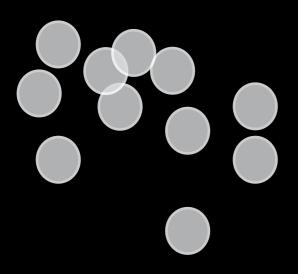


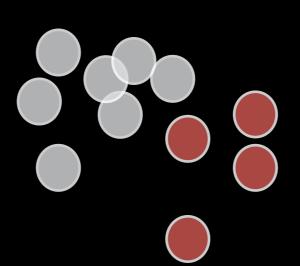


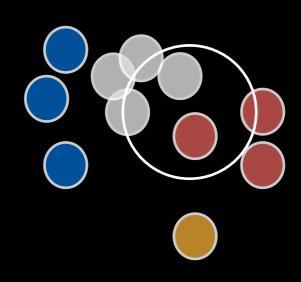


devisive (top-down):

e.g. hierarchical k-mean







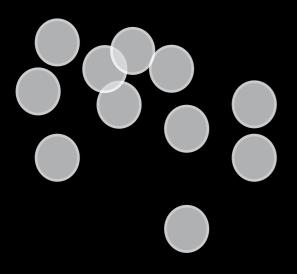
it is non-deterministic

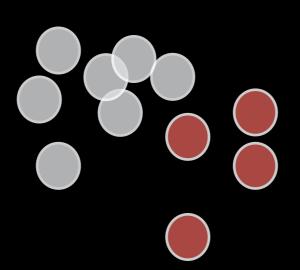
it is *greedy* just as k-means
two nearby points
may end up in
separate clusters

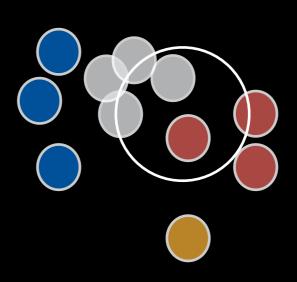


devisive (top-down):

e.g. hierarchical k-mean







it is non-deterministic

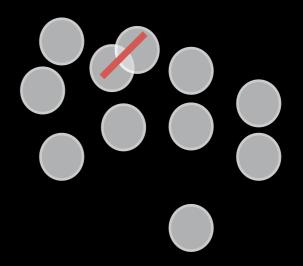
it is *greedy* just as k-means
two nearby points
may end up in
separate clusters

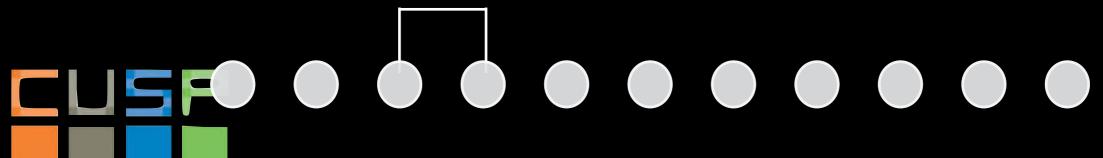
it is simple and fast: complexity

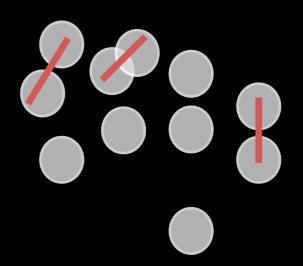
O(NdK log<sub>k</sub>N)

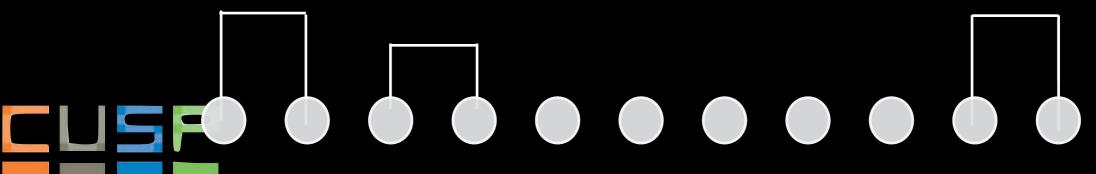


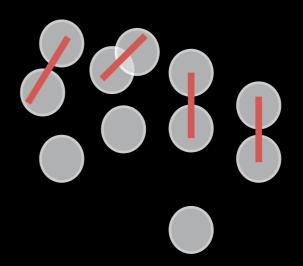
XI: Clustering

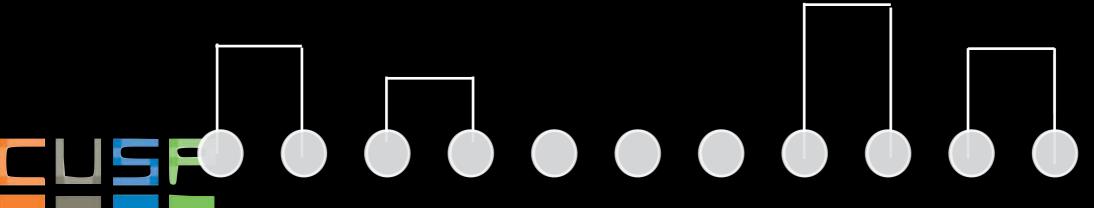


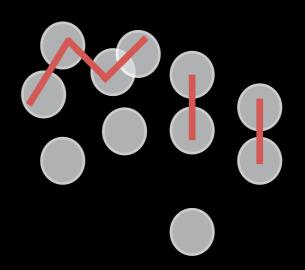


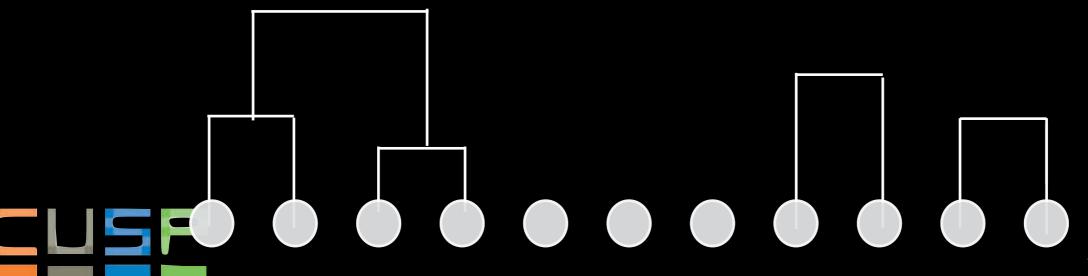


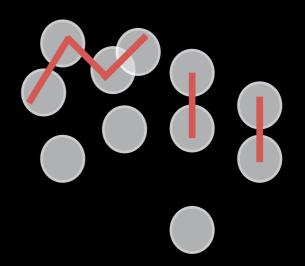


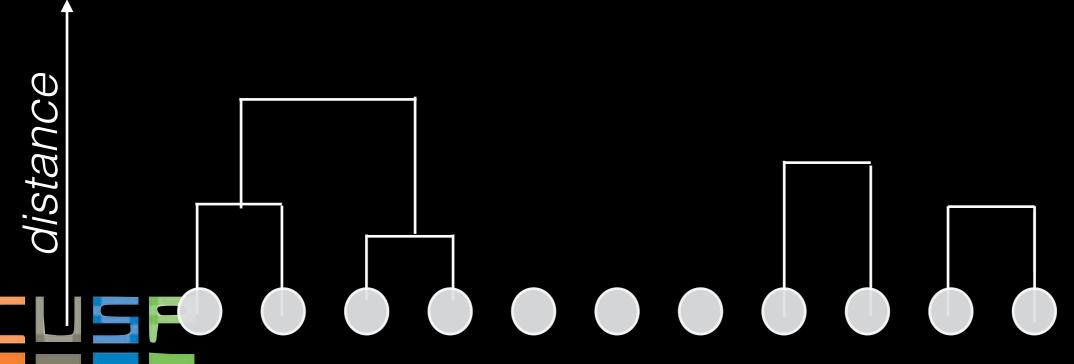




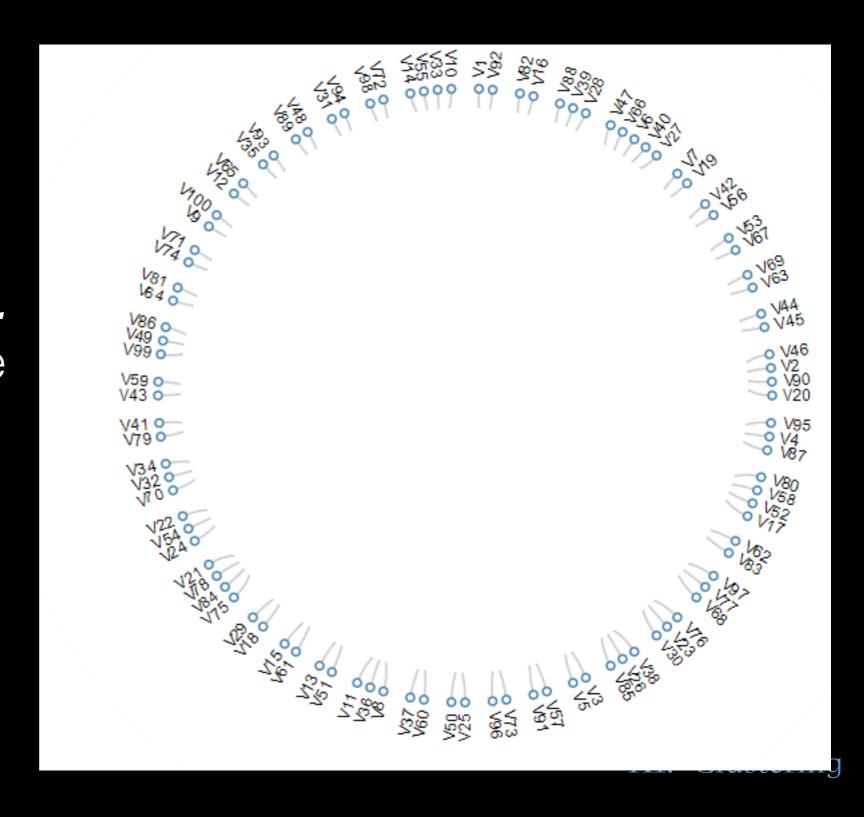






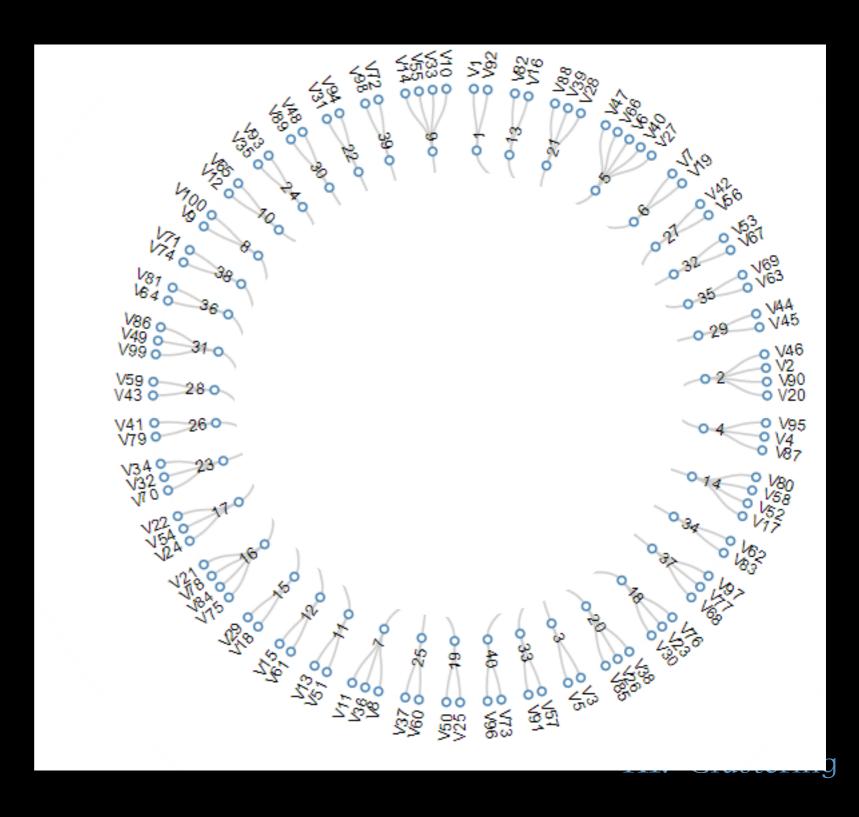


**agglomerative** bottom-up



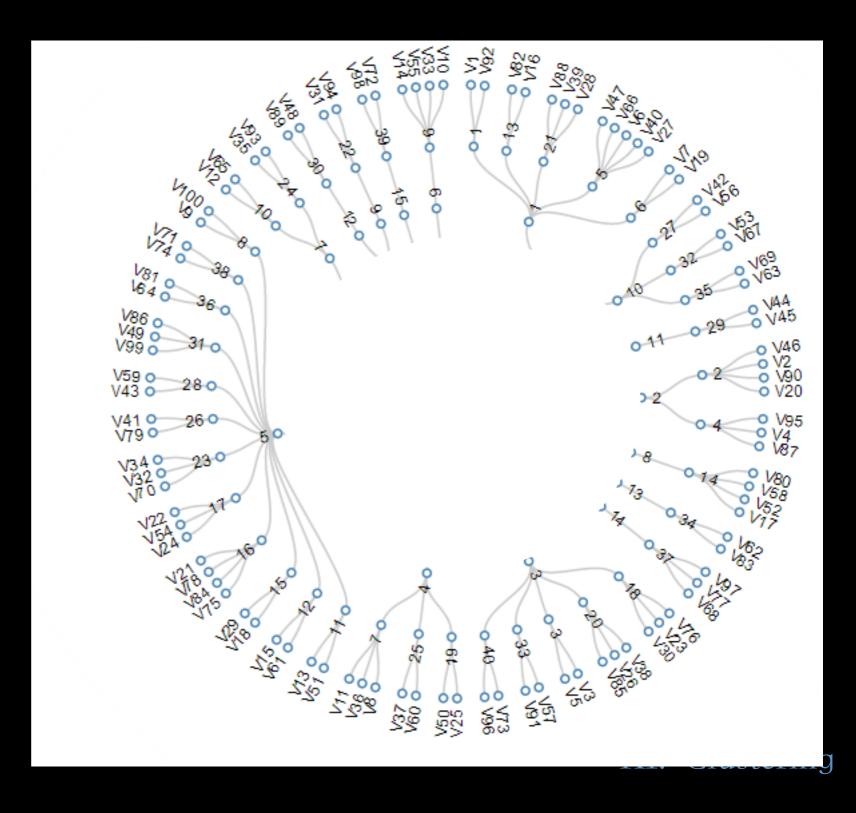


agglomerative bottom-up



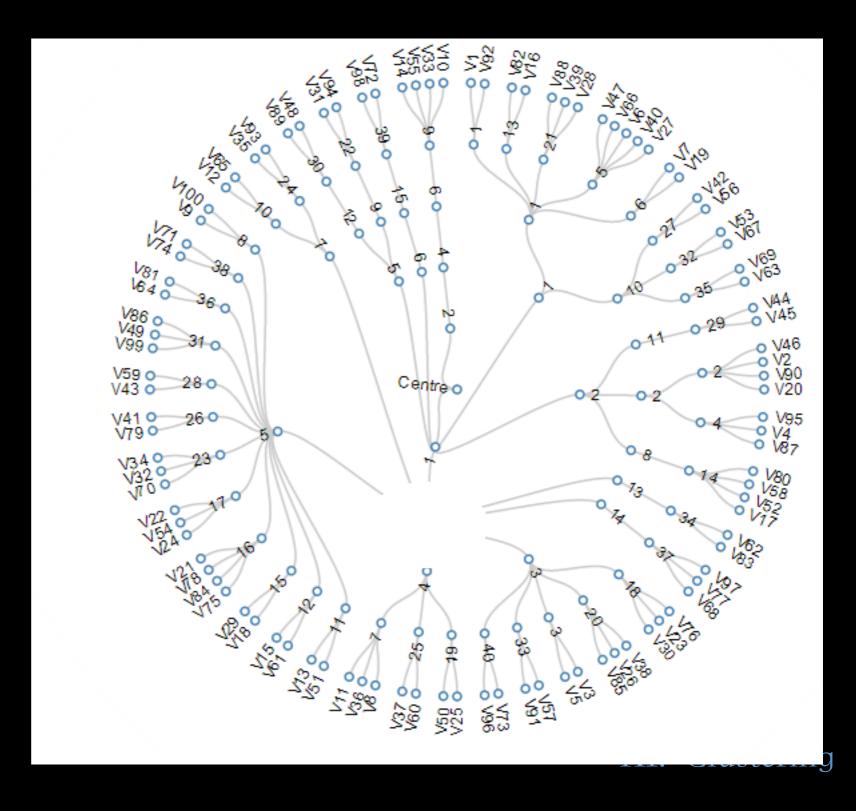


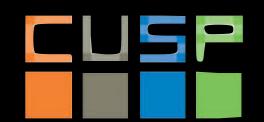
agglomerative bottom-up



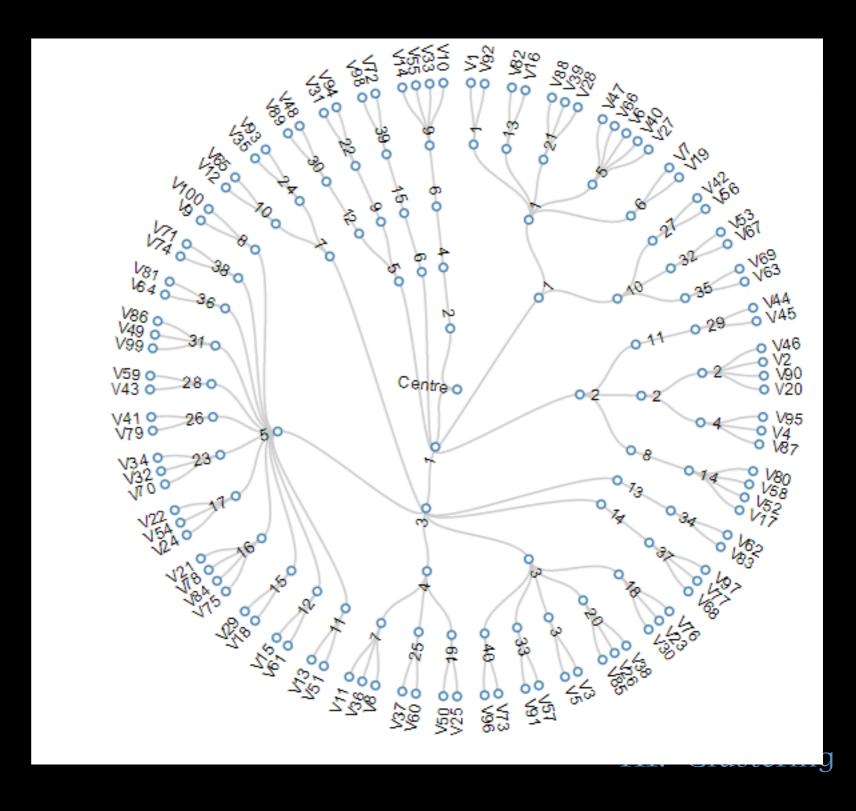


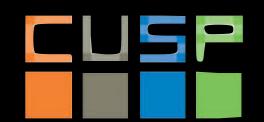
agglomerative bottom-up





agglomerative bottom-up





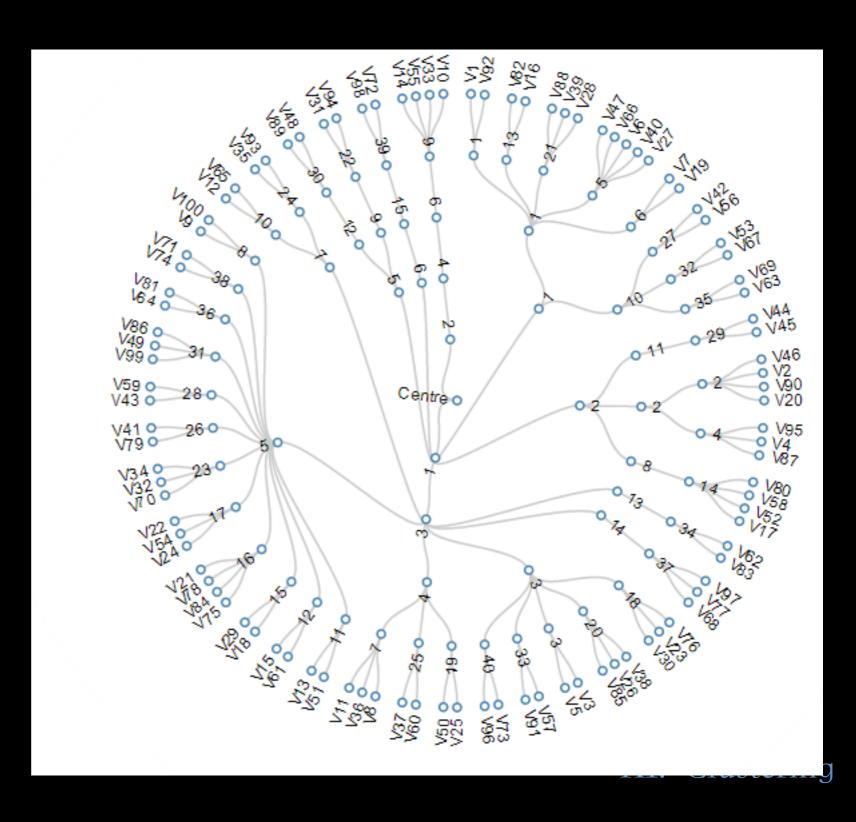
**agglomerative** bottom-up

computationally intense because every *cluster pair* distance has to be calculate

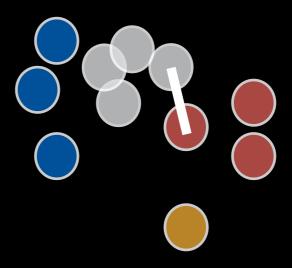
it is slow, though it can be optimize: complexity

 $O(N^2d + N^3)$ 





**agglomerative** bottom-up



## Requires to define

- a distance &
- a "linkage"



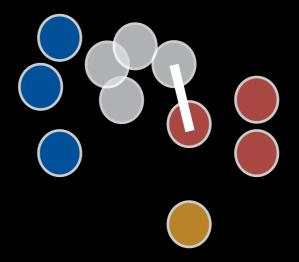
**agglomerative** bottom-up

single link distance

$$D(c1,c2) = min(D(x_{c1}, x_{c2}))$$

complete link distance

$$D(c1,c2) = max(D(x_{c1}, x_{c2}))$$



## Requires to define

- a distance &
- a "linkage"



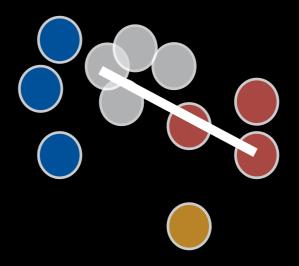
**agglomerative** bottom-up

single link distance

$$D(c1,c2) = min(D(x_{c1}, x_{c2}))$$

complete link distance

$$D(c1,c2) = max(D(x_{c1}, x_{c2}))$$



## Requires to define

- a distance &
- a "linkage"



# **agglomerative** bottom-up

single link distance

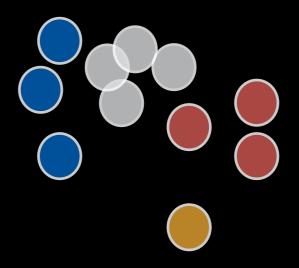
$$D(c1,c2) = min(D(x_{c1}, x_{c2}))$$

complete link distance

$$D(c1,c2) = max(D(x_{c1}, x_{c2}))$$

centroid distance

$$D(c1,c2) = mean(D(x_{c1}, x_{c2}))$$



#### a common distance:

ward distance minimizes variance

$$D_{tot} = \sum_{j} \sum_{i,x_i \in C_j} (x_i - \mu_j)^2$$



https://github.com/fedhere/UInotebooks/blob/master/cluster/

XI: Clustering

#### **Summary and Key concepts**

#### clustering is easy, but interpreting results is tricky

Distane metrics:

Eucledian and other Minchowski metrics geospacial distances metrics for non continuous data

Partitioning methods: inexpensive, typically non deterministic

Hard methods: *K-means, K-medoids* 

Soft (or fuzzy) methods: (i.e. probabilistic approach)

Expectation Maximization Mixture models

Hierarchical methods:

divisive vs agglomerative, dendrograms



#### **RESOURCES:**

#### a comprehensive review of clustering methods

Data Clustering: A Review, Jain, Mutry, Flynn 1999 <a href="https://www.cs.rutgers.edu/~mlittman/courses/lightai03/jain99data.pdf">https://www.cs.rutgers.edu/~mlittman/courses/lightai03/jain99data.pdf</a>

# a blog post on how to generate and interpret a scipy dendrogram by Jörn Hees

https://joernhees.de/blog/2015/08/26/scipy-hierarchical-clustering-and-dendrogram-tutorial/



$$F(\omega) = \frac{1}{2\pi} \int f(t)e^{-i\omega t} dt$$



$$F(\omega) = \frac{1}{2\pi} \int f(t)e^{-i\omega t} dt$$

takes a function in time domain



$$F(\omega) = \frac{1}{2\pi} \int f(t)e^{-i\omega t} dt$$

takes a function in time domain

to a function in frequency domain



$$F(\omega) = \frac{1}{2\pi} \int f(t)e^{-i\omega t} dt$$

takes a function in space domain

to a function in spatial frequency domain

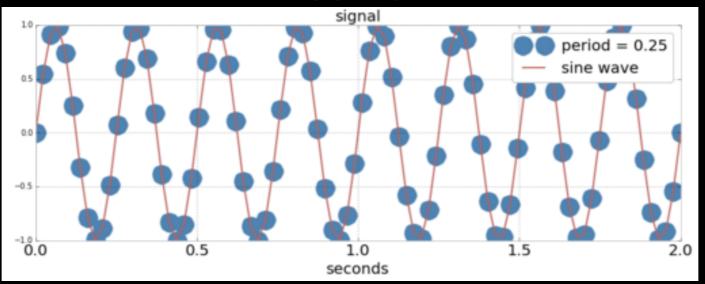


$$F(\omega) = \frac{1}{2\pi} \int f(t)e^{-i\omega t} dt$$

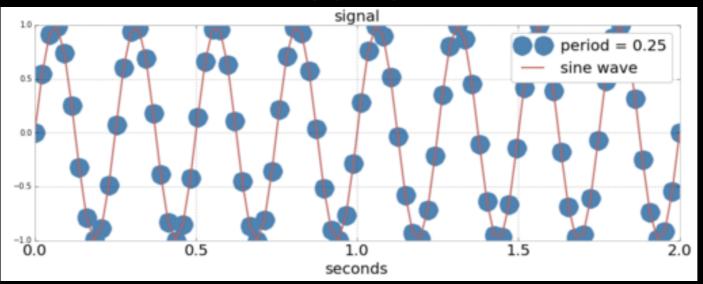
takes a function in space domain f(t) is measured in seconds

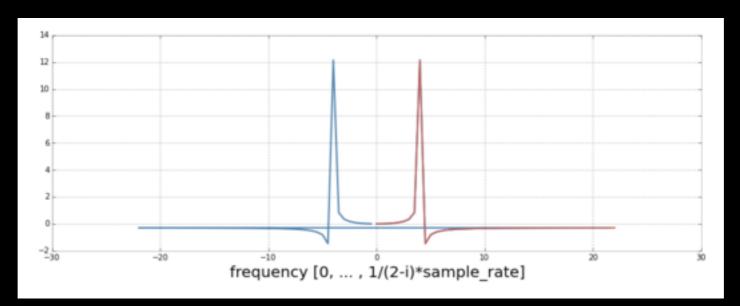
to a function in spatial frequency domain f(t) is measured in 1/seconds or Hz

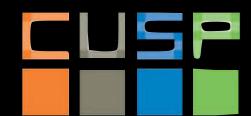








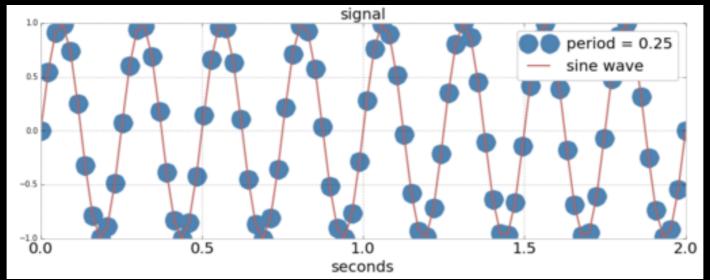


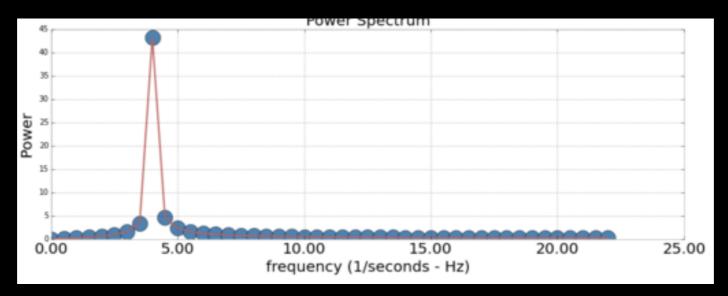


The absolute value of the square of the Fourier transform this is called a Power spectrum.

High value of the power spectrum indicate periodicity at the corresponding frequency

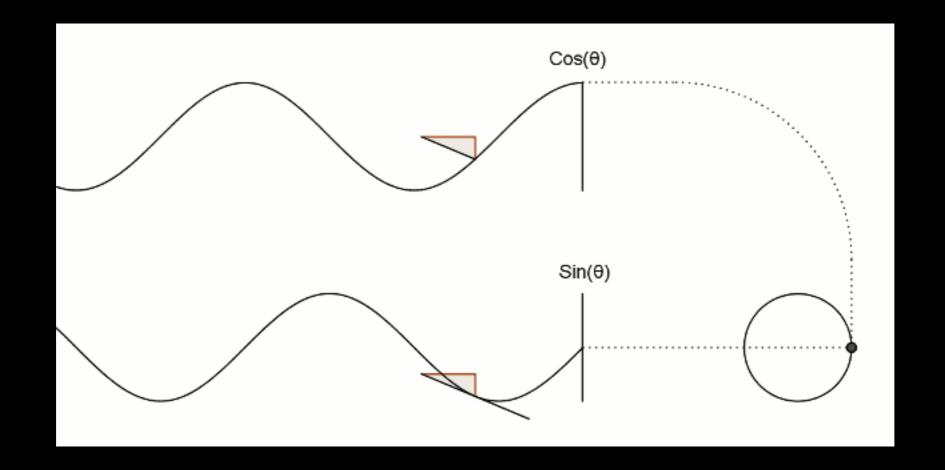




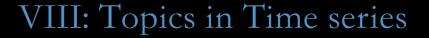




#### Cosine and Sine... just in case



http://www.businessinsider.com/7-gifs-trigonometry-sine-cosine-2013-5





https://github.com/fedhere/Ulnotebooks/blob/master/fourier.ipynb



#### Homework:

Reading: an excellent analysis of time series by Jake Vander Plas (UW e-science center)

https://jakevdp.github.io/blog/2014/06/10/is-seattle-really-seeing-an-uptick-in-cycling/



#### **READING:**

your data aint that big...

https://www.chrisstucchio.com/blog/2013/hadoop\_hatred.html

#### **HW: cluster NYC business history**

- 1. cluster the economic trends in NYC using 2 methods: use K-Means and another method of your choice (e.g. DBscan, agglomerative clustering): use the time behavior of the number of establishments per zip code as your feature space
- 2. see if the clusters based on the time behavior also form spatial clusters.



map the time-based clusters (e.g. with geopanda as a heat map). attempt an interpretation.

XI: Clustering

#### **HW: cluster NYC business history**

Use census data for NYC businesse:

number of establishments per zip code for ~20 years since 1994

you can get the zip code info (list of NYC zip codes and shape files for plotting) here:

http://data.nycprepared.org/dataset/nyc-zip-code-tabulation-areas/resource/0c0e14e9-78e1-404e-97b0-c2fabceb3981

this is the link to the census business data <a href="http://www.census.gov/econ/cbp/download/">http://www.census.gov/econ/cbp/download/</a>

you can download manually, which is labour intensive, or on the terminal via ftp, which requires some wrangling, but i did that for you! (see below)

for ((y=93; y<=99; y+=1)); do wget zbp\$y\totals.zip; done for ((y=0; y<=9; y+=1)); do wget ftp://ftp.census.gov/econ200\$y\\ CBP\_CSV/zbp0\$y\totals.zip; done



for ((y=10; y<=15; y+=1)); do wget ftp://ftp.census.gov/econ20\$y\\ CBP\_CSV/zbp\$y\totals.zip; done XI: Clustering