

50 years beyond K-means

Authored By: Anil K. Jain

What's After

ans

BIRCH

Presented by: Zheng, Pamela

K-Means?

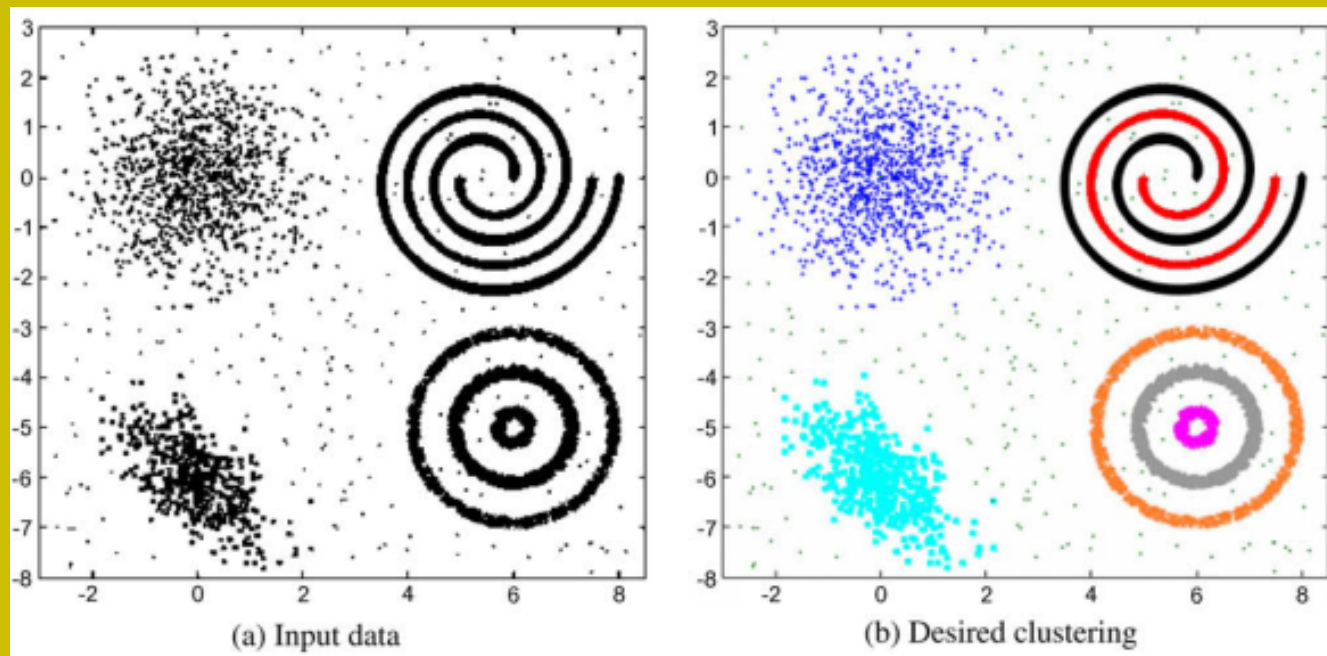
Dynamic



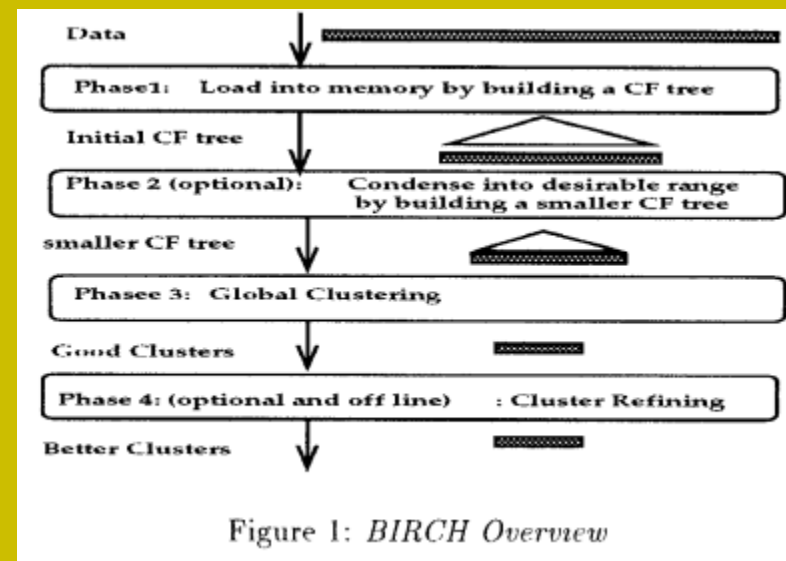
The rising STAR of Texas

Clustering

- Authored By: Anil K. Jain
- Published in Pattern Recognition Letters 31, 2009
- Emphasizes the importance of grouping data.
- Two types of algorithms:
 - Partitional: builds clusters simultaneously
 - Hierarchical: joins or divides clusters.



- Presented by: Zhang, Ramakrishna
- SIGMOD 1996
- Aimed to handle big data within
- Compares data points locally rat
- Can produce good results on a s
- Captures hierarchy of the data



Run Time (sec)

Fig

shnan, and Livny

thin limited memory.

rather than globally.

a single scan of data.

a as a CF Tree.

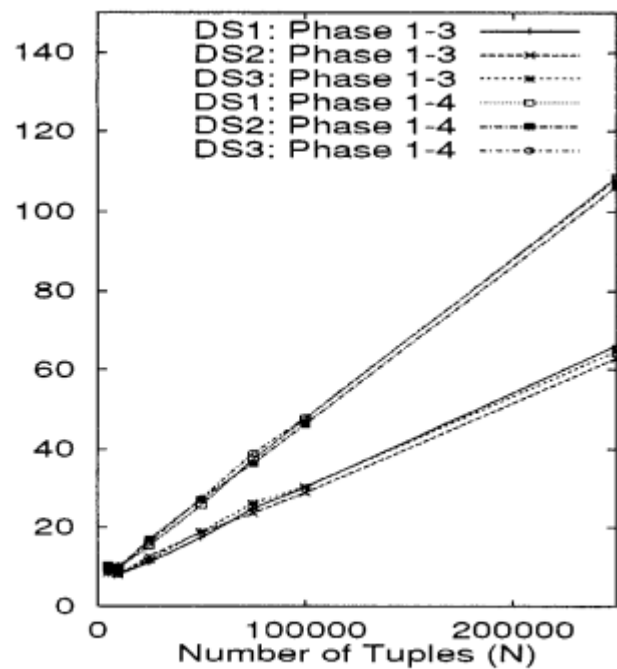
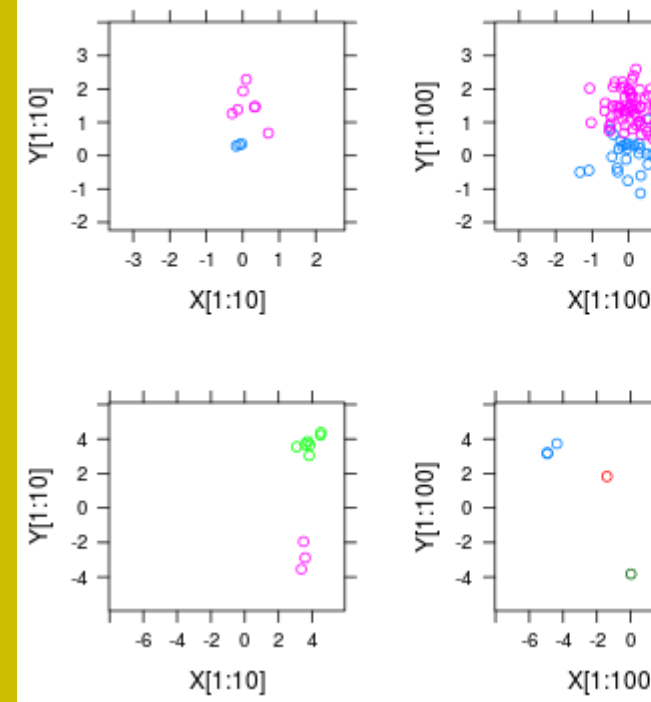


Figure 4: Scalability wrt. Increasing n_l, n_h

- Presented by: Camp
- ANIPS, 2013
- Based on Dirichlet
- Non-parametric: mod
- Clusters can be crea
- Promises high speed



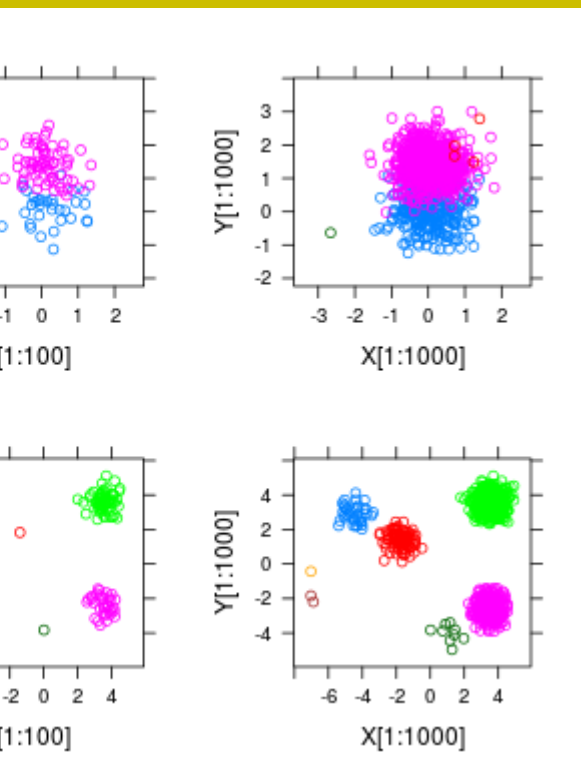
Campbell, Liu, Kulis, How, and Carin.

at Processes

model grows with data.

created, eliminated, and altered.

eed for time-sensitive applications.



Global Kernel K-Means

- Presented by: Tzortzis and Likas
- IEEE Joint Conference on Neural Networks in 2008
- Kernel K-Means finds non linear separable clusters.
- Global K-Means minimizes sensitivity to initialization.
- Comes in “fast” and “near optimal” variety.

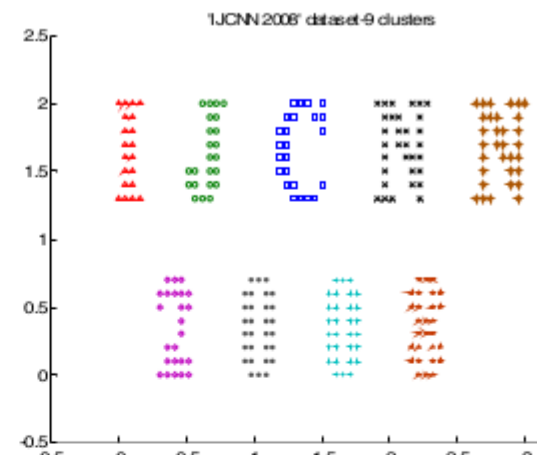
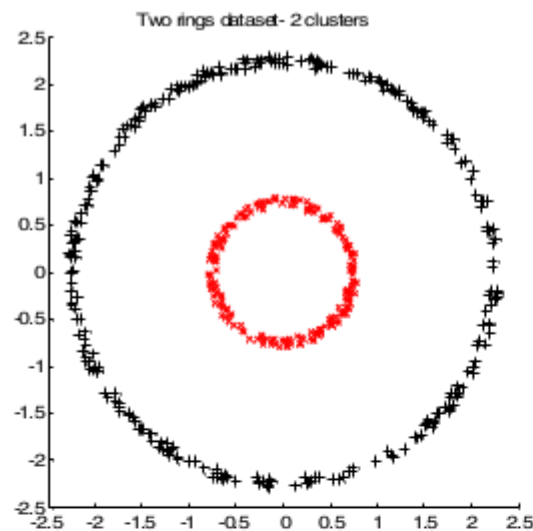
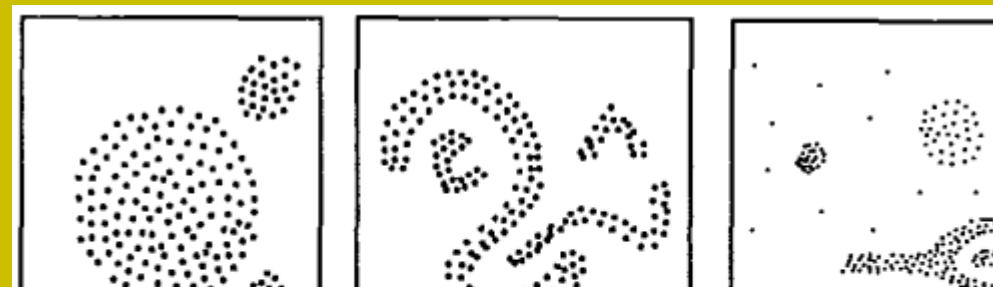


Fig. 2. Global kernel k -means, fast global kernel k -means

A Density Bas

- Presented by: Ester, Kriegel, San
- Association for the Advancement
- Only takes one input parameter.
- Clusters by density of points rat
- Can define clusters of arbitrary
- This is a partitioning algorithm.



ased Algorithm

Sander and Xu.

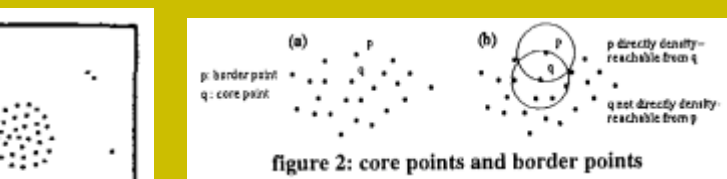
nt of AI in 1996

er.

rather than centroids.

y shapes.

m.



So What

- Focus on Billion-Scale
- Solve the “Curse of D
- Distance measures bre
- Centroids make assum
- Parallelism is a must.
- Fuzzy clustering shoul
- Good algorithms do no



Comes Next?

ale data.

Dimensionality” rather than avoiding.

break down at high dimension

umptions about structure of data.

st.

ould be a possibility.

not make assumptions about data.



Fig. 1. Global kernel k -means, fast global kernel k -means and kernel k -means (run with minimum clustering error) on the two rings dataset.

Fig. 2. Global kernel k -means, fast global kernel k -means and kernel k -means (run with minimum clustering error) on the 'IJCNN 200

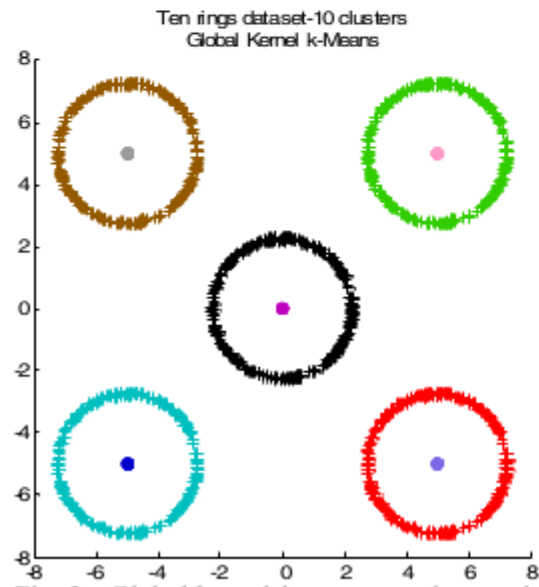


Fig. 3. Global kernel k -means on the ten rings dataset.

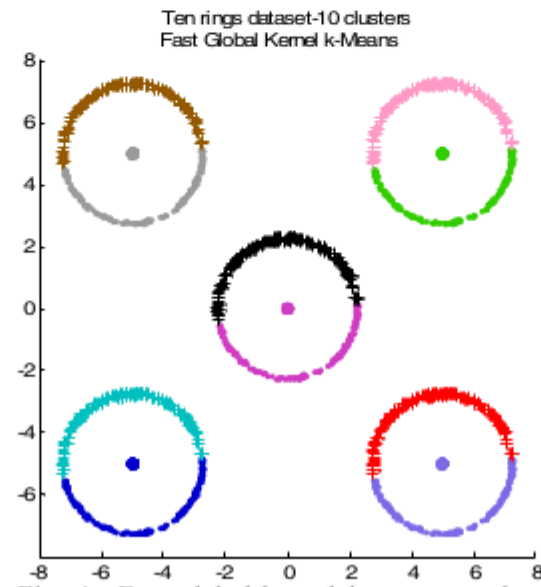


Fig. 4. Fast global kernel k -means on the ten rings dataset.

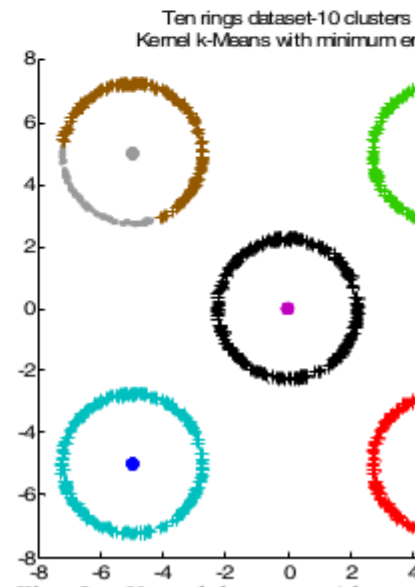
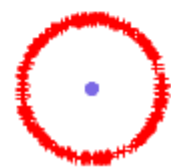
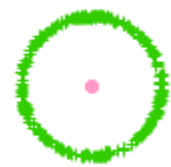


Fig. 5. Kernel k -means (the run with minimum clustering error) on the ten rings dataset.

means and kernel k -
N 2008' logo.

clusters
imum error



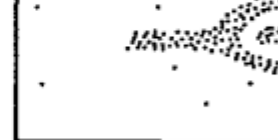
he run with mini-
n rings dataset.



database 1



database 2



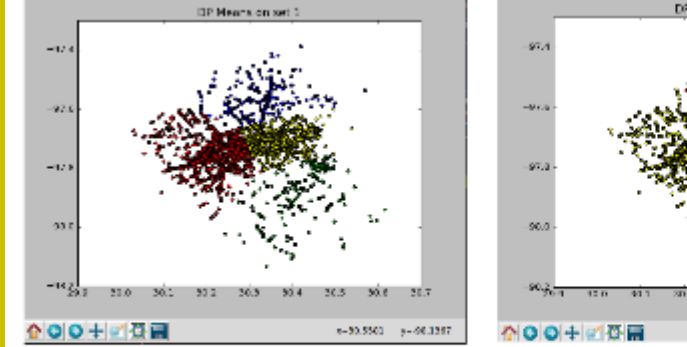
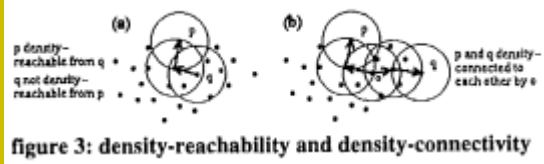
database 3

figure 1: Sample databases

Presenter: Gentry Atkinson

Faculty Mentors: Dr. Tesic and Dr. Tami

e 3

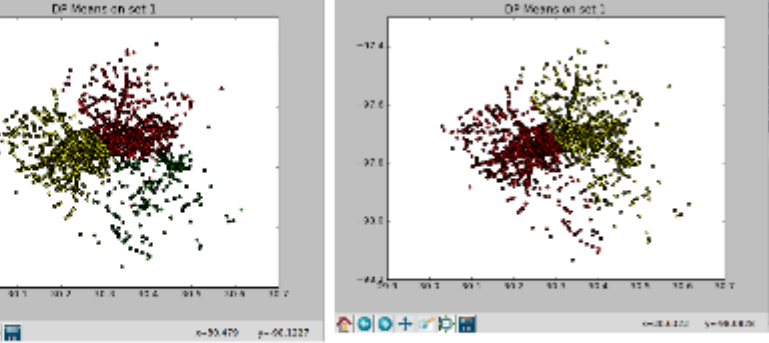


$\Lambda = 0.3$
4 clusters

$\Lambda = 0.4$
3 clusters

A demonstration showing a slight shift in the number of clusters found by varying Λ , finding 2, 3, or 4 clusters on Austin

mir



$s = 0.4$
s

Lambda = 0.5
2 clusters

shift in the DP Means input parameter
stin traffic accidents.