Analytics2011 CONFERENCE SERIES

Dynamical Data Clustering

Carl Meyer, North Carolina State University

Copyright © 2011, SAS Institute Inc. All rights reserved

#analytics2011

Cluster Analytics

Classify objects from a set into clusters

Objects in the same cluster are more similar to each other than to those in other clusters.

Detect, reveal, and analyze hidden patterns

A Primary Exploratory Task

- Data and text mining
- Pattern recognition
- Image analysis

- Network analysis
- Information retrieval
- Bioinformatics

Fundamental Theorem Of Clustering

Nothing Works Always!

So Many Choices

- Hierarchical & agglomerative clustering
- k-means and its many variations and derivatives
- Nonnegative matrix factorization
- Spectral and subspace clustering
- Graph partitioning and min-cut techniques
- PDDP and PCA based partitioning algorithms
- Self organizing maps & neural network methods
- Gaussian mixture models & generalizations
- Nearest neighbor implementations
- Hard vs. fuzzy
- · · · more · · · and more · · · and more · · ·

Static Data vs. Dynamic Data

Hidden patterns are more difficult define in static data

But most clustering algorithms are built for the analysis of static data



Is It Relevant How The Fish Moves?

By "fish intellegence?"

- Complicated system of PDE's ?
- Strictly at random?

Is It Relevant How The Fish Moves?

By "fish intellegence?"

Complicated system of PDE's ?

Strictly at random?

Doesn't Matter

- As long as the spots on the fish stay attached to fish's body
- Spots move together and relatively faster to slower movement of the background.

Dynamical Clustering

Somehow impart "motion" to static data to reveal clusters moving by observing which data points move in concert relative to background noise.

Dynamical Clustering

Somehow impart "motion" to static data to reveal clusters moving by observing which data points move in concert relative to background noise.

The Strategy

Observe the evolution of a differentiated time-scale stochastic process imposed on the static data.

Dynamical Clustering

Somehow impart "motion" to static data to reveal clusters moving by observing which data points move in concert relative to background noise.

The Strategy

Observe the evolution of a differentiated time-scale stochastic process imposed on the static data.

How To Do It?

Reverse Simon–Ando process

Simon–Ando Theory

- Herbert Simon (1916–2001)
- Carnegie Mellon University



- Nobel Prize in economics in 1978

- Albert Ando (1930–2002)
- University of Pennsylvania



"Aggregation of variables in dynamic systems," *Econometrica*, Vol. 29, No. 2 (Apr., 1961), pp. 111-138.

The Goal Of Simon–Ando

- Analyze long-term economic stability of a macro economy containing closely coupled micro economies by analyzing (or observing) the evolution of the micro economies for a short period of time.
- Small example
 - Nine industries
 - Three closely coupled clusters



Cluster 1 = Manufacturing (steel, machine tools, heavy equipment) Cluster 2 = Entertainment (movies, TV, books & magazines) Cluster 3 = Beverage (sugar, water, packaging)





Modeled by a Markov chain

Modeled by a Markov chain

Individual industries are the states

Modeled by a Markov chain

Individual industries are the states

Flow of capital between industries are the transitions

Modeled by a Markov chain

Individual industries are the states

Flow of capital between industries are the transitions

- Transition flows are row normalized (row sums = 1)
 - $\mathbf{P}_{n \times n}$ row stochastic transition matrix
 - Aperiodic Markov chain

- Modeled by a Markov chain
- Individual industries are the states
- Flow of capital between industries are the transitions
- Transition flows are row normalized (row sums = 1)
- $\mathbf{P}_{n \times n}$ row stochastic transition matrix
- Aperiodic Markov chain
- ► *k* distinct micro economies (or clusters)
 - There is a permutation such that **P** is nearly uncoupled

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} & \dots & \mathbf{P}_{1k} \\ \mathbf{P}_{21} & \mathbf{P}_{22} & \dots & \mathbf{P}_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{P}_{k1} & \mathbf{P}_{k2} & \dots & \mathbf{P}_{kk} \end{bmatrix} \quad \max_{i} ||\mathbf{P}_{i\star}||_{\infty} = \zeta << 1$$

Uncoupling By Censoring



$$\mathbf{S}_1 \equiv \mathbf{P}_{11} + \mathbf{Q}_1 \qquad \mathbf{S}_2 \equiv \mathbf{P}_{22} + \mathbf{Q}_2 \qquad \mathbf{S}_3 \equiv \mathbf{P}_{33} + \mathbf{Q}_3$$





Censored probabilities

$$Q_{1} = [P_{12} \quad P_{13}] \begin{bmatrix} I - P_{22} & P_{23} \\ P_{32} & I - P_{33} \end{bmatrix}^{-1} \begin{bmatrix} P_{21} \\ P_{31} \end{bmatrix}$$
$$= P_{1*} \left(I - \widetilde{P}_{11} \right)^{-1} P_{*1}$$





Censored probabilities

$$Q_{1} = [P_{12} \quad P_{13}] \begin{bmatrix} I - P_{22} & P_{23} \\ P_{32} & I - P_{33} \end{bmatrix}^{-1} \begin{bmatrix} P_{21} \\ P_{31} \end{bmatrix}$$
$$= P_{1*} \left(I - \widetilde{P}_{11} \right)^{-1} P_{*1}$$

Censored transition matrix for censored chain

$$S_1 = P_{11} + Q_1 = P_{11} + P_{1*} \left(I - \widetilde{P}_{11} \right)^{-1} P_{*1}$$





Censored probabilities

$$Q_{1} = [P_{12} \quad P_{13}] \begin{bmatrix} I - P_{22} & P_{23} \\ P_{32} & I - P_{33} \end{bmatrix}^{-1} \begin{bmatrix} P_{21} \\ P_{31} \end{bmatrix}$$
$$= P_{1*} \left(I - \widetilde{P}_{11} \right)^{-1} P_{*1}$$

Censored transition matrix for censored chain

$$S_1 = P_{11} + Q_1 = P_{11} + P_{1*} \left(I - \widetilde{P}_{11} \right)^{-1} P_{*1}$$

► In General,
$$\mathbf{S}_i = \mathbf{P}_{ii} + \mathbf{P}_{i*} \left(\mathbf{I} - \widetilde{\mathbf{P}}_{ii} \right)^{-1} \mathbf{P}_{*i}$$

These are called stochastic complements

Steady States

• σ_i = steady-state distribution of i^{th} censored chain $\sigma_i(0)$ = initial prob dist $\sigma_i(t) = \sigma_i(t-1)\mathbf{S}_i = t$ -step dist $\sigma_i = \lim_{t \to \infty} \sigma_i(t)$ $\sigma_i = \sigma_i \mathbf{S}_i$

Steady States

 σ_i = steady-state distribution of i^{th} censored chain $\sigma_i(0)$ = initial prob dist $\sigma_i(t) = \sigma_i(t-1)\mathbf{S}_i = t$ -step dist $\sigma_i = \lim_{t \to \infty} \sigma_i(t)$ $\sigma_i = \sigma_i \mathbf{S}_i$

 $\pi = \text{steady-state distribution of global chain}$ $\pi(0) = \text{initial prob dist}$ $\pi(t) = \pi(t-1)\mathbf{P} = t\text{-step dist}$ $\pi = \lim_{t \to \infty} \pi(t)$ $\pi = \pi \mathbf{P}$

Steady States

 σ_i = steady-state distribution of i^{th} censored chain $\sigma_i(0)$ = initial prob dist $\sigma_i(t) = \sigma_i(t-1)\mathbf{S}_i = t$ -step dist $\sigma_i = \lim_{t \to \infty} \sigma_i(t)$ $\sigma_i = \sigma_i \mathbf{S}_i$ π = steady-state distribution of global chain $\pi(0)$ = initial prob dist $\pi(t) = \pi(t-1)\mathbf{P} = t$ -step dist $\pi = \lim_{t \to \infty} \pi(t)$ $\pi = \pi \mathbf{P}$

Coupling Theorem

• $\pi = (\xi_1 \sigma_1 \ \xi_2 \sigma_2 \ \xi_3 \sigma_3), \quad \xi_1, \xi_2, \xi_3 \text{ are "coupling" constants}$

Equibrium Phases

Short Run ($a \le t \le b$)

 $\boldsymbol{\pi}(t) \approx \left(\xi_1(t) \boldsymbol{\sigma}_1 \quad \xi_2(t) \boldsymbol{\sigma}_2 \quad \xi_3(t) \boldsymbol{\sigma}_3 \right) \qquad \qquad \xi_i(t) = \left\| \boldsymbol{\pi}_i(0) \right\|_1 = \nu_i$

CONSTANTS

Equlibrium Phases

Short Run $(a \le t \le b)$ $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \xi_i(t) = \|\pi_i(0)\|_1 = \nu_i$

Constants

• Middle Run (t > b)

 $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \qquad \xi_i(t) \text{ Varies With Time}$

Equlibrium Phases

Short Run $(a \le t \le b)$ $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \qquad \xi_i(t) = \|\pi_i(0)\|_1 = \nu_i$ CONSTANTS

Middle Run (t > b)

 $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \xi_i(t)$ Varies with Time

For example, consider states i and j in cluster #1

$$\frac{\pi_i(t)}{\pi_j(t)} \approx \frac{\xi_1(t)[\sigma_1]_i}{\xi_1(t)[\sigma_1]_j} = \frac{[\sigma_1]_i}{[\sigma_1]_j} = \text{ a constant} \qquad (t > a)$$

Equlibrium Phases

Short Run $(a \le t \le b)$ $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \qquad \xi_i(t) = \|\pi_i(0)\|_1 = \nu_i$ CONSTANTS

Middle Run (t > b)

 $\pi(t) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3) \qquad \qquad \xi_i(t) \text{ Varies With Time}$

For example, consider states i and j in cluster #1

$$\frac{\pi_i(t)}{\pi_j(t)} \approx \frac{\xi_1(t)[\sigma_1]_i}{\xi_1(t)[\sigma_1]_j} = \frac{[\sigma_1]_i}{[\sigma_1]_j} = \text{ a constant} \qquad (t > a)$$

▶ Long Run $\xi_i(t) \to \xi_i$ (Constant) as $t \to \infty$

 $\pi(t) \to \pi(\infty) = (\xi_1 \sigma_1 \ \xi_2 \sigma_2 \ \xi_3 \sigma_3) = \pi$ (Global Equibrium)

Simon–Ando Conclusions

Short-run behavior reveals long-run behavior

Simon–Ando Conclusions

Short-run behavior reveals long-run behavior

Long-run equibrium in a macro economy containing clusters of micro economies is determined by the short-run evolution of the micros.

Simon–Ando Conclusions

Short-run behavior reveals long-run behavior

Long-run equibrium in a macro economy containing clusters of micro economies is determined by the short-run evolution of the micros.

Longer-term economic predictions can be made from shorterterm observations.

Reverse Simon–Ando

Long-run behavior reveals short-run behavior

$$\pi(\infty) = \pi = (\pi_1 \ \pi_2 \ \pi_3) \Longrightarrow \ \sigma_i = \frac{\pi_i}{\|\pi_i\|_1}$$

Reverse Simon–Ando

Long-run behavior reveals short-run behavior

$$\pi(\infty) = \pi = (\pi_1 \ \pi_2 \ \pi_3) \Longrightarrow \ \sigma_i = \frac{\pi_i}{\|\pi_i\|_1}$$

 $\pi(a \leq t \leq b) \approx (\nu_1 \sigma_1 \ \nu_2 \sigma_2 \ \nu_3 \sigma_3)$ where $\nu_i = \|\pi_i(0)\|_1$ (Short-run stabilization)
Reverse Simon–Ando

Long-run behavior reveals short-run behavior

$$\pi(\infty) = \pi = (\pi_1 \ \pi_2 \ \pi_3) \Longrightarrow \sigma_i = \frac{\pi_i}{\|\pi_i\|_1}$$
 $\pi(a \le t \le b) \approx (\nu_1 \sigma_1 \ \nu_2 \sigma_2 \ \nu_3 \sigma_3) \text{ where } \nu_i = \|\pi_i(0)\|_1$
(Short-run stabilization)

And middle-run behavior

 $\pi(t > b) \approx (\xi_1(t)\sigma_1 \ \xi_2(t)\sigma_2 \ \xi_3(t)\sigma_3)$ where $\xi_i(t) \rightarrow \|\pi_i\|_1$

Suppose that $\pi(\infty)$ is uniform

$$\pi(\infty) = \frac{1}{n} \left(\underbrace{1, 1, \dots, 1}_{n_1}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_3} \right) = \left(\begin{array}{cc} \mathbf{n}_i = \# \text{ states in cluster } i \\ \mathbf{\pi}_1 & \mathbf{\pi}_2 & \mathbf{\pi}_3 \end{array} \right)$$

Suppose that $\pi(\infty)$ is uniform

 $\boldsymbol{\pi}(\infty) = \frac{1}{n} \left(\underbrace{1, 1, \dots, 1}_{n_1}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_3}, \underbrace{n_3}_{n_3} \right) = \left(\boldsymbol{\pi}_1 \quad \boldsymbol{\pi}_2 \quad \boldsymbol{\pi}_3 \right)$

•
$$\sigma_i = \pi_i / \|\pi_i\|_1 = \frac{1}{n_i} (\underbrace{1, 1, ..., 1}^{n_i})$$

Suppose that $\pi(\infty)$ is uniform

 $\boldsymbol{\pi}(\infty) = \frac{1}{n} \left(\underbrace{1, 1, \dots, 1}_{n_1}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_3} \right) = \left(\begin{array}{cc} \boldsymbol{\pi}_1 & \boldsymbol{\pi}_2 & \boldsymbol{\pi}_3 \end{array} \right)$

$$\boldsymbol{\sigma}_i = \boldsymbol{\pi}_i / \left\| \boldsymbol{\pi}_i \right\|_1 = \frac{1}{n_i} \left(\overbrace{1, 1, \dots, 1}^{\mathbf{n}_i} \right)$$

 $\alpha_i = \nu_i / n_i$

 $\pi(a \leq t \leq b) \approx (\nu_1 \sigma_1 \ \nu_2 \sigma_2 \ \nu_3 \sigma_3) = (\alpha_1 \cdots \alpha_1 | \alpha_2 \cdots \alpha_2 | \alpha_3 \cdots \alpha_3)$

Suppose that $\pi(\infty)$ is uniform $\pi(\infty) = \frac{1}{n} \left(\overbrace{1,1,\dots,1}^{\mathbf{n}_{1}}, \overbrace{1,1,\dots,1}^{\mathbf{n}_{2}}, \overbrace{1,1,\dots,1}^{\mathbf{n}_{3}} \right) = \left(\pi_{1} \quad \pi_{2} \quad \pi_{3} \right)$ $\sigma_{i} = \pi_{i} / \left\| \pi_{i} \right\|_{1} = \frac{1}{n_{i}} \left(\overbrace{1,1,\dots,1}^{\mathbf{n}_{i}} \right)$ $\sigma_{i} = \nu_{i} / n_{i}$

 $\boldsymbol{\pi}(a \leq t \leq b) \approx (\nu_1 \boldsymbol{\sigma}_1 \quad \nu_2 \boldsymbol{\sigma}_2 \quad \nu_3 \boldsymbol{\sigma}_3) = (\alpha_1 \cdots \alpha_1 \mid \alpha_2 \cdots \alpha_2 \mid \alpha_3 \cdots \alpha_3)$

 $\pi(t > b) \approx (\xi_1(t)\sigma_1 \quad \xi_2(t)\sigma_2 \quad \xi_3(t)\sigma_3)$ $= (\beta_1(t)\cdots\beta_1(t) \mid \beta_2(t)\cdots\beta_2(t) \mid \beta_3(t)\cdots\beta_3(t)) \qquad \beta_{i}(t) \to 1/n$

Suppose that $\pi(\infty)$ is uniform

 $\boldsymbol{\pi}(\infty) = \frac{1}{n} \left(\underbrace{1, 1, \dots, 1}_{n_1}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_2}, \underbrace{1, 1, \dots, 1}_{n_3} \right) = \left(\begin{array}{cc} \mathbf{n}_i = \# \text{ states in cluster } i \\ \mathbf{\pi}_1 & \mathbf{\pi}_2 & \mathbf{\pi}_3 \end{array} \right)$

$$\boldsymbol{\sigma}_i = \boldsymbol{\pi}_i / \|\boldsymbol{\pi}_i\|_1 = \frac{1}{n_i} \left(\underbrace{1, 1, ..., 1}_{n_i} \right)$$

 $\pi(a \leq t \leq b) \approx (\nu_1 \sigma_1 \quad \nu_2 \sigma_2 \quad \nu_3 \sigma_3) = (\alpha_1 \cdots \alpha_1 \mid \alpha_2 \cdots \alpha_2 \mid \alpha_3 \cdots \alpha_3)$

 $\boldsymbol{\pi}(t > b) \approx \left(\xi_1(t) \boldsymbol{\sigma}_1 \quad \xi_2(t) \boldsymbol{\sigma}_2 \quad \xi_3(t) \boldsymbol{\sigma}_3 \right)$

 $= (\beta_1(t) \cdots \beta_1(t) | \beta_2(t) \cdots \beta_2(t) | \beta_3(t) \cdots \beta_3(t)) \qquad \beta_i(t) \to 1/n$

Nearly equal entries in $\pi(t > a)$ belong to the same cluster





Answer

Simply force P to be doubly stochastic

- i.e., force all row sums and all column sums = 1



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows ightarrow scale columns ightarrow scale rows



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns \rightarrow scale rows \rightarrow scale columns



Answer

Simply force P to be doubly stochastic

- i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns \rightarrow scale rows \rightarrow scale columns \rightarrow etc.



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns \rightarrow scale rows \rightarrow scale columns \rightarrow etc.

Converges (usually) — if not, it can be forced



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns \rightarrow scale rows \rightarrow scale columns \rightarrow etc.

Converges (usually) — if not, it can be forced

Sinkhorn–Knopp preserves cluster integrity



Answer

Simply force **P** to be doubly stochastic — i.e., force all row sums and all column sums = 1

Sinkhorn–Knopp Procedure

Scale rows \rightarrow scale columns \rightarrow scale rows \rightarrow scale columns \rightarrow etc.

Converges (usually) — if not, it can be forced

Sinkhorn–Knopp preserves cluster integrity

Sinkhorn–Knopp preserves symmetry

Putting Things Together

Raw Data

► Step 1. $\mathbf{A}_{mn} = [\mathbf{a}_1 | \mathbf{a}_2 | \cdots | \mathbf{a}_n] \rightarrow \mathbf{C}_{nn}$ (A Similarity Matrix)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ $d(\star, \star) = \text{similarity metric of your choice}$

Putting Things Together

Raw Data

► Step 1. $\mathbf{A}_{mn} = [\mathbf{a}_1 | \mathbf{a}_2 | \cdots | \mathbf{a}_n] \rightarrow \mathbf{C}_{nn}$ (A SIMILARITY MATRIX)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ $d(\star, \star) = \text{similarity metric of your choice}$

- Euclidean distance
- Minkowski norms
- Correlation
- Angular distance
- Gaussian metrics
- Hamming distance or variations
- Gabriel graph
- Delaunay triangulation
- Mean first passage time
- Ensemble (consensus) metrics
- etc

Putting Things Together

Raw Data

► Step 1. $\mathbf{A}_{mn} = [\mathbf{a}_1 | \mathbf{a}_2 | \cdots | \mathbf{a}_n] \rightarrow \mathbf{C}_{nn}$ (A SIMILARITY MATRIX)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ $d(\star, \star) = \text{similarity metric of your choice}$

- Applications usually dictates choice
- e.g., text vs. numeric
- $\mathbf{C} = \mathbf{C}^T$ (symmetric matrix—usually)

- Euclidean distance
- Minkowski norms
- Correlation
- Angular distance
- Gaussian metrics
- Hamming distance or variations
- Gabriel graph
- Delaunay triangulation
- Mean first passage time
- Ensemble (consensus) metrics
- etc

Step 2. Sinkhorn–Knopp Procedure

- $\mathbf{C}_{nn} \rightarrow \mathbf{P}_{nn}$ (by successive row & column scaling)
- Similarity matrix (Symmetric) \rightarrow doubly stochastic (Symmetric)

Step 2. Sinkhorn–Knopp Procedure

- $\mathbf{C}_{nn} \rightarrow \mathbf{P}_{nn}$ (by successive row & column scaling)
- Similarity matrix (Symmetric) \rightarrow doubly stochastic (Symmetric)

Step 3. Initialize Markov Chain

- Pick $\pi(0)$ to be significantly different than uniform
- To see where state (or data point) #i clusters, pick

$$\pi(0) = \mathbf{e}_i = (0, 0, ..., \underbrace{\mathbf{1}}_{i}, 0, ..., 0)$$

Step 2. Sinkhorn–Knopp Procedure

- $\mathbf{C}_{nn} \rightarrow \mathbf{P}_{nn}$ (by successive row & column scaling)
- Similarity matrix (Symmetric) \rightarrow doubly stochastic (Symmetric)

Step 3. Initialize Markov Chain

- Pick $\pi(0)$ to be significantly different than uniform
- To see where state (or data point) #*i* clusters, pick $\pi(0) = \mathbf{e}_i = (0, 0, ..., \underbrace{1}_{i}, 0, ..., 0)$

Step 4. Observe Markov Chain For A Few Steps

- $\pi(t+1) = \pi(t)\mathbf{P}$ t = 1, 2, ..., a (until short-run stabilization)
- Order entries in $\pi(t > a)$ to identify gaps
- Nearly equal entries belong to the same cluster

Leukemia Exmple

 AML — Acute myeloid leukemia (Most common in adults)



 ALL — Acute lymphoblastic leukemia (Most common in children)



In AML, myeloid stem cells develop into immature abnormal white blood cells, myeloblasts, that don't become healthy white blood cells

In ALL, too many stem cells develop into lymphoblasts that don't mature into lymphocytes, the white blood cells required to fight infections.



DNA Microarray



Broad Institute (MIT/Harvard)

- "Molecular classification of cancer: class discovery and class prediction by gene expression monitoring," T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasen-beek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfeld, and E. S. Lander, *Science*, October 1999.
- "Metagenes and molecular pattern discovery using matrix factorization," J. P. Brunet, P. Tamayo, T. Golub, J. Mesirov, *Proceedings of the National Academy of Sciences*, March, 2004.
- S8 cancer patients gene expression data from bone marrow samples 5000 genes
- ► Patient Diagonsis: #1–19 = ALL(B), #20–27 = ALL(T), #28–38 = ALM
- Good test case since clusters are known
- Similarity matrix was build by a consensus method

The Iterations (Chuck














































































Determining The Number Of Clusters

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} & \mathbf{P}_{13} \\ \mathbf{P}_{21} & \mathbf{P}_{22} & \mathbf{P}_{23} \\ \mathbf{P}_{31} & \mathbf{P}_{32} & \mathbf{P}_{33} \end{bmatrix} \approx \begin{bmatrix} \mathbf{S}_1 & & \\ & \mathbf{S}_2 & \\ & & \mathbf{S}_3 \end{bmatrix}$$

Determining The Number Of Clusters



Determining The Number Of Clusters



k = # clusters = # eigenvalues near $\lambda = 1$ (determined by largest gap)

Eigenvalues For Leukemia Data



▶ Step 1. A (Raw data) \rightarrow C = C^T (Similarity matrix)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ via similarity metric of your choice

Step 1. **A** (Raw data) \rightarrow **C** = **C**^{*T*} (Similarity matrix) $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ via similarity metric of your choice

Step 2. Sinkhorn–Knopp Scaling

C (Similarity matrix) \rightarrow **P** = **P**^T (Doubly stochastic)

Step 1. A (Raw data) \rightarrow C = C^T (Similarity matrix)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ via similarity metric of your choice

Step 2. Sinkhorn–Knopp Scaling C (Similarity matrix) \rightarrow **P** = **P**^T (Doubly stochastic)

Step 3. Determine k (the number of clusters)
k = # eigenvalues of P closest to 1 (Determined by largest eigengap)

Step 1. A (Raw data) \rightarrow C = C^T (Similarity matrix)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ via similarity metric of your choice

Step 2. Sinkhorn–Knopp Scaling

C (Similarity matrix) \rightarrow **P** = **P**^T (Doubly stochastic)

Step 3. Determine k (the number of clusters)

k = # eigenvalues of P closest to 1 (Determined by largest eigengap)

Step 4. Initialize Markov Chain

Pick $\pi(0)$ to be significantly different than uniform

Step 1. A (Raw data) \rightarrow C = C^T (Similarity matrix)

 $c_{ij} = d(\mathbf{a}_i, \mathbf{a}_j)$ via similarity metric of your choice

Step 2. Sinkhorn–Knopp Scaling

C (Similarity matrix) \rightarrow **P** = **P**^{*T*} (Doubly stochastic)

Step 3. Determine k (the number of clusters)

k = # eigenvalues of **P** closest to 1 (Determined by largest eigengap)

Step 4. Initialize Markov Chain Pick $\pi(0)$ to be significantly different than uniform

Step 5. Observe The Chain For A Few Steps

 $\pi(t + 1) = \pi(t)\mathbf{P}$ t = 1, 2, ..., a (Until short-run stabilization) Order $\pi(t > a)$ and partition at the k - 1 largest gaps States in each of these k segments define the k clusters

Allows for dynamic visualization tools for cluster analysis

Allows for dynamic visualization tools for cluster analysis

Allows the ability to select & analyze clusters in isolation

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation

Allows tracking & analysis of selected pieces of data

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data

Works well for determining the number of clusters

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data Works well for determining the number of clusters

Several possibilities for variations & innovations

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data Works well for determining the number of clusters Several possibilities for variations & innovations

Remaining Questions

Scalability

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data Works well for determining the number of clusters Several possibilities for variations & innovations

Remaining Questions

Scalability

Application dependency

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data Works well for determining the number of clusters Several possibilities for variations & innovations

Remaining Questions

Scalability

Application dependency

Sensitivity to similarity metric

Allows for dynamic visualization tools for cluster analysis Allows the ability to select & analyze clusters in isolation Allows tracking & analysis of selected pieces of data Works well for determining the number of clusters Several possibilities for variations & innovations

Remaining Questions

Scalability

Application dependency

Sensitivity to similarity metric

Effectiveness of application directly to raw data

Thanks For Your Attention!