

Lecture 11: High Dimensionality Information Visualization

CPSK 533C, Fall 2009

Tamara Munzner

UBC Computer Science

Wed, 21 October 2009

Readings Covered

Highdimensional Data Analysis Using Parallel Coordinates. Edward J. Wegman. *Journal of the American Statistical Association*, Vol. 85, No. 411. (Sep. 1990), pp. 664-675.

Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Ying-Hui Fu, Matthew O. Ward, and Elke A. Rundensteiner. *IEEE Visualization '99*.

Clmmer: Multiload MOS on the GPU. Stephan Ingram, Tamara Munzner and Mac Olano. *IEEE TVCG*, 15(2):240-261, Mar/Apr 2009.

Cluster Stability and the Use of Noise in Interpretation of Clustering. George S. Davidson, Brian N. Wyllie, Kevin W. Boyack, *Proc InfoVis 2001*.

Interactive Hierarchical Dimension Ordering, Spacing and Filtering for Exploration Of High Dimensional Datasets. Jing Yang, Wei Peng, Matthew O. Ward and Elke A. Rundensteiner. *Proc. InfoVis 2003*.

Further Reading

Visualizing the non-visual: spatial analysis and interaction with information from text documents. James A. Wise et al. *Proc. InfoVis 1995*

Parallel Coordinates: A Tool for Visualizing Multi-Dimensional Geometry. Alfred Inselberg and Bernard Dimsdale, *IEEE Visualization '90*.

A Data-Driven Reflectance Model. Wojciech Matusik, Hanspeter Pfister, Matt Brandt, and Leonard McMillan. *SIGGRAPH 2003*. graphics.cs.msl.edu/~wojciech/pubs/sg03.pdf

Parallel Coordinates

- only 2 orthogonal axes in the plane
- instead, use parallel axes!



[Highdimensional Data Analysis Using Parallel Coordinates. Edward J. Wegman. *Journal of the American Statistical Association*, 85(411), Sep 1990, p. 664-675.]

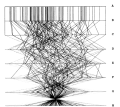
11-08

11-09

11-08

11-08

PC: Correlation



[Highdimensional Data Analysis Using Parallel Coordinates. Edward J. Wegman. *Journal of the American Statistical Association*, 85(411), Sep 1990, p. 664-675.]

11-08

11-09

11-08

11-08

PC: Duality

- rotate-translate
- point-line
 - pencil: set of lines coincident at one point



[Parallel Coordinates: A Tool for Visualizing Multi-Dimensional Geometry. Alfred Inselberg and Bernard Dimsdale, *IEEE Visualization '90*.]

11-08

11-09

11-08

11-08

PC: Axis Ordering

- geometric interpretations
 - hyperplane, hypersphere
 - points do have intrinsic order
- infovis
 - no intrinsic order, what to do?
 - indeterminate/arbitrary order
 - weakness of many techniques
 - downside: human-powered search
 - upside: powerful interaction technique
- most implementations
 - user can interactively swap axes
- Automated Multidimensional Detective
 - Inselberg 99
 - machine learning approach

11-08

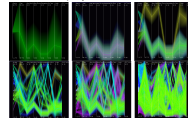
11-09

11-08

11-08

Hierarchical Parallel Coords: LOD

- variable-width opacity bands



[Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Fu, Ward, and Rundensteiner. *IEEE Visualization '96*.]

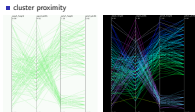
11-08

11-09

11-08

11-08

Proximity-Based Coloring



[Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Fu, Ward, and Rundensteiner. *IEEE Visualization '96*.]

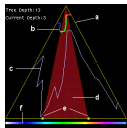
11-08

11-09

11-08

11-08

Structure-Based Brushing



[Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Fu, Ward, and Rundensteiner. *IEEE Visualization '96*.]

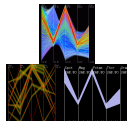
11-08

11-09

11-08

11-08

Dimensional Zooming



[Hierarchical Parallel Coordinates for Visualizing Large Multivariate Data Sets. Fu, Ward, and Rundensteiner. *IEEE Visualization '96*.]

11-08

11-09

11-08

11-08

Critique

- not easy for novices
- now used in many apps
- hier: major scalability improvements
 - combination of encoding, interaction

11-08

11-09

11-08

11-08

Dimensionality Reduction

- mapping multidimensional space into space of fewer dimensions
 - filter subset of original dimensions
 - generate new synthetic dimensions
- why is lower-dimensional approximation useful?
 - assume true/intrinsic dimensionality of dataset is (much) lower than measured dimensionality!
- why would this be the case?
 - only indirect measurement possible
 - fisheries ex: want spawn rates, have water color, air temp, catch rates...
 - sparse data in verbose space
 - documents ex: word occurrence vectors 10K+ dimensions, want dozens of topic clusters

11-08

11-09

11-08

11-08

Dimensionality Reduction: Isomap

- 4096 D: pixels in image
- 2D: wrist rotation, fingers extension



[A Global Geometric Framework for Nonlinear Dimensionality Reduction. J. B. Tenenbaum, V. de Silva, and J. C. Langford. *Science* 290(5501), pp 2316-2322, Oct 22 2000]

11-08

11-09

11-08

11-08

Goals/Tasks

- goal: keep/explain as much variance as possible
- find clusters
 - or compare/evaluate vs. previous clustering
- understand structure
 - absolute position not reliable
 - arbitrary rotations/reflections in lowD map
 - fine-grained structure not reliable
 - coarse near/far positions safer

11-08

11-09

11-08

11-08

Dimensionality Analysis Example

- measuring materials for image synthesis
 - BRDF measurements: 4M samples × 103 materials
 - goal: lowD model where can interpolate

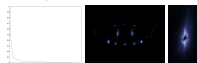


[A Data-Driven Reflectance Model, SIGGRAPH 2003, W. Matusik, H. Pfister, M. Brand and L. McMillan, graphics.toronto.edu/~wajtech/pubs/tg2003.pdf]

01-18

Dimensionality Analysis: Linear

- how many dimensions is enough?
 - could be more than 2 or 3
 - find kink in curve: error vs. dims used
- linear dim reduct: PCA, 25 dims
 - physically impossible intermediate points when interpolate



[A Data-Driven Reflectance Model, SIGGRAPH 2003, W. Matusik, H. Pfister, M. Brand and L. McMillan, graphics.toronto.edu/~wajtech/pubs/tg2003.pdf]

01-18

Dimensionality Analysis: Nonlinear

- nonlinear dim reduct (charting): 10-15
 - all intermediate points physically possible

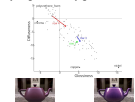


[A Data-Driven Reflectance Model, SIGGRAPH 2003, W. Matusik, H. Pfister, M. Brand and L. McMillan, graphics.toronto.edu/~wajtech/pubs/tg2003.pdf]

01-18

Meaningful Axes: Nameable by People

- red, green, blue, specular, diffuse, glossy, metallic, plastic-y, roughness, rubbery, greasiness, dustiness...



[A Data-Driven Reflectance Model, SIGGRAPH 2003, W. Matusik, H. Pfister, M. Brand and L. McMillan, graphics.toronto.edu/~wajtech/pubs/tg2003.pdf]

01-18

MDS: Multidimensional scaling

- large family of methods
 - minimize differences between interpoint distances in high and low dimensions
 - distance scaling: minimize objective function
 - stress(D, Δ) = $\sqrt{\sum_{i,j} \frac{(d_{ij} - \Delta_{ij})^2}{d_{ij}^2}}$
 - D : matrix of lowD distances
 - Δ : matrix of hD distances d_{ij}

01-18

Spring-Based MDS: Naive

- repeat for all points
 - compute spring force to all other points
 - difference between high dim, low dim distance
 - move to better location using computed forces
- compute distances between all points
 - $O(n^2)$ iteration, $O(n^2)$ algorithm



01-18

Faster Spring Model: Stochastic

- compare distances only with a few points
 - maintain small local neighborhood set



01-18

Faster Spring Model: Stochastic

- compare distances only with a few points
 - maintain small local neighborhood set
 - each time pick some randoms, swap in if closer



01-18

Faster Spring Model: Stochastic

- compare distances only with a few points
 - maintain small local neighborhood set
 - each time pick some randoms, swap in if closer



01-18

Faster Spring Model: Stochastic

- compare distances only with a few points
 - maintain small local neighborhood set
 - each time pick some randoms, swap in if closer
 - small constant: 6 locals, 3 randoms typical
 - $O(n)$ iteration, $O(n^2)$ algorithm



01-18

Glimmer Algorithm

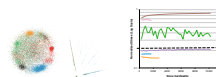


- multilevel, designed to exploit GPU
 - restriction to decimate
 - relaxation as core computation
 - relaxation to interpolate up to next level
- GPU stochastic as subsystem
 - poor convergence properties if run alone
 - low-pass-filter stress approx. for termination

[Glimmer: Multilevel MDS on the GPU, Ingham, Munster and Otazo. IEEE TVCG, 15(2):249-260, Mar/Apr 2008.]

Glimmer Results

- sparse document dataset: 28K dims, 28K points



[Glimmer: Multilevel MDS on the GPU, Ingham, Munster and Otazo. IEEE TVCG, 15(2):249-260, Mar/Apr 2008.]

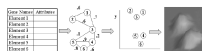
01-18

Cluster Stability

- display
 - also terrain metaphor
- underlying computation
 - energy minimization (springs) vs. MDS
 - weighted edges
- do same clusters form with different random start points?
 - "ordination"
 - spatial layout of graph nodes

01-18

Approach



- normalize within each column
- similarity metric
 - discussion: Pearson's correlation coefficient
- threshold value for marking as similar
 - discussion: finding critical value

01-18

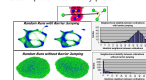
Graph Layout

- criteria
 - geometric distance matching graph-theoretic distance
 - vertices one hop away close
 - vertices many hops away far
 - insensitive to random starting positions
 - major problem with previous work!
 - tractable computation
- force-directed placement
 - discussion: energy minimization
 - others: gradient descent, etc
 - discussion: termination criteria

01-18

Barrier Jumping

- same idea as simulated annealing
 - but compute directly
 - just ignore repulsion for fraction of vertices
- solves start position sensitivity problem

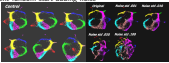


01-18

Results

- efficiency
 - naive approach: $O(V^2)$
 - approximate density field: $O(V)$
- good stability
 - rotation/reflection can occur

different random start adding noise



01-16

Critique

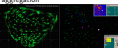
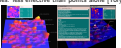
01-16

Critique

- real data
 - suggest check against subsequent publication
- give criteria, then discuss why solution fits
- visual + numerical results
 - convincing images plus benchmark graphs
- detailed discussion of alternatives at each stage
- specific prescriptive advice in conclusion

01-16

MDS Beyond Points

- galaxies: aggregation
 - 
- themescapes: terrain/landscapes
 - studies: less effective than points alone [Tory 07, 09]
 - 

[www.pd4.com/infvis/graphics.html] Visualizing the non-visual: spatial analysis and interaction with information from text documents. James A. Wise et al. Proc. InfoVis 2006

01-16

Dimension Ordering

- in NP: heuristic, like most interesting infovis problems
- divide and conquer
 - iterative hierarchical clustering
 - representative dimensions
- choices
 - similarity metrics
 - importance metrics
 - variance
- ordering algorithms
 - optimal
 - random swap
 - simple depth-first traversal

01-16

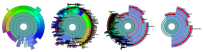
Spacing, Filtering

- same idea: automatic support
- interaction
 - manual intervention
 - structure-based brushing
 - focus+context

01-16

Results: InterRing

- raw, order, distort, rollup (filter)

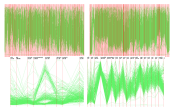


[Interactive Hierarchical Dimension Ordering, Spacing and Filtering for Exploration of High Dimensional Datasets. Yang Peng, Ward, and Rundenstein. Proc. InfoVis 2003]

01-16

Results: Parallel Coordinates

- raw, order/space, zoom, filter




[Interactive Hierarchical Dimension Ordering, Spacing and Filtering for Exploration of High Dimensional Datasets. Yang Peng, Ward, and Rundenstein. Proc. InfoVis 2003]

01-16

Results: Star Glyphs

- raw, order/space, distort, filter

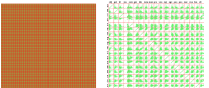


[Interactive Hierarchical Dimension Ordering, Spacing and Filtering for Exploration of High Dimensional Datasets. Yang Peng, Ward, and Rundenstein. Proc. InfoVis 2003]

01-16

Results: Scatterplot Matrices

- raw, filter



[Interactive Hierarchical Dimension Ordering, Spacing and Filtering for Exploration of High Dimensional Datasets. Yang Peng, Ward, and Rundenstein. Proc. InfoVis 2003]

01-16

Critique

01-16

Critique

- pro
 - approach on multiple techniques,
 - real data!
- con
 - always show order then space then filter
 - hard to tell which is effective
 - show ordered vs. unsorted after zoom/filter?

01-16

Reminders

- meet with me before end of week!
- presentation topics also due Friday
 - your call whether presentation and project topics match
 - submit: 3 topic choices, web day
- project data/task ideas on resources page
 - VAST/InfoVis Contest!

01-16

Readings Next Week

Graph Visualization in Information Visualization: a Survey. Ivan Herman, Cor Meinel, M. Scott Marshall. IEEE Transactions on Visualization and Computer Graphics, 6(1), pp. 24-44, 2000. <http://citeseer.ri.ac.com/herman01graph.html>

change: Configuring Hierarchical Layouts to Address Research Questions. Adrian Stingle, Jason Dwyer, and Jo Wood. IEEE Transactions on Visualization and Computer Graphics, 15 (8), Nov-Dec 2009 (Proc. InfoVis 2009). <http://ieeexplore.ieee.org/xml/document/stingle09vis2009.pdf>

Multiscale Visualization of Small World Networks. David Adler, Yves Chikina, Fabien Jordan, Guy Melancon. Proc. InfoVis 2003. <http://ieeexplore.ieee.org/xml/document/adler03vis2003.pdf>

Topological Fish-eye Views for Visualizing Large Graphs. Emrah Gonen, Yehuda Koren and Stephen North. IEEE TVCG 11(4), p. 467-486, 2005. <http://www.research.att.com/~emrah/visualizations/papers/atlant.pdf/DLDP-conf-infovis-CasconK05a.pdf>

ISep-CiL: An Incremental Procedure for Separation Constraint Layout of Graphs. Tim Dwyer, Kim Marriott, and Yehuda Koren. Proc. InfoVis 2006, published as IEEE TVCG 12(5), Sep 2006, p. 621-630. <http://www.research.att.com/~yehuda/pubs/bayer.pdf>

01-16

01-16

01-16