A Taxonomy of Visual Cluster Separation Factors



Michael SedImair, Andrada Tatu, Tamara Munzner, Melanie Tory University of British Columbia, Vancouver, CA <u>http://www.cs.ubc.ca/~msedl</u>



Background

• Phd at LMU, Munich with Andreas Butz: HCI

117

• Automotive Design Studies at BMW











0.73 1-0					
a manual					Inter Tree
		100000	_		140
		-		and the second se	Page State State
		Kalifa Address			Anna and a state
ALC: 107	States States	In Frank State		100	Pillano.
ALC: 12	Manual V	NAMES OF BRIDE	adapted in	and the second	
				1.1	
	SPAN-	Transing .		10.0	and the second second
***	Bright Arty	inter 1	. Ballanteres	tata.	Cours book
440 (2)	Markets.	16	RAUSEN.	10.4	
	1 C C C C C C C C C C C C C C C C C C C			1	
		TopPacket one			
	arteau	8	a. afaana		
	-		aratica		
	artena Kitaketa		s. sranna	-	
	artana Mining		NUMBER OF	sin.	
	artean Managarta		Service	un.	
	Artesar Artesar		New York Contraction	in a	
	Artister		Transver	tan.	
	Artiste References		Terrarella	nan.	Nam 110, Alara (100
 Participation Particip	Alfanan Kalauna Bernata Kalaun		Response	None Concept	Patri International Internatio
Variation of the second	Articular Sectorial Marchines and Articles and Articles Marchines and Articles and Articles Point Articles and Articles and Articles Marchines and Articles and Articles		Tengar Press	Non Barrier Al Lawrence St (2003 1)	100 100 100 100 100 100 100 100
T Cuangi us T Cuangi us Talanta Talanta Talanta Talanta Talanta Talanta Talanta Talanta Talanta Talanta	Artesta Millioner (1990) Artesta Millioner (1990) Artesta A		New York	1 2014	Pater Dis Annual Inst Restor Dis Annual Inst Pater Dis Annual Inst Pater Dis Annual Inst Pater Dis Annual Inst
T Valend Second Second			Transfer Tra	Normal Distance A Constant of A Constant of	No. No. State State State State State State State State State State State State State





Background

- PostDoc at UBC, Vancouver with Tamara Munzner: InfoVis
 - Highdim projects
 - Ethnographic field study: What people are doing with DR?



Data Characterization



 Design Study Methodology (with Miriah Meyer, University of Utah)

A Taxonomy of Visual Cluster Separation Factors

Motivation

Episode I: Once upon a time ...

 What visualization encodings (VisEnc) are good for visualizing dimensionally reduced (DR) data?



What are people doing?

- Informed from ongoing ethnographic field study (not this project!)
- DR techniques
 - PCA, MDS, but also newer ones such as t-SNE
- VisEnc
 - 2d Scatterplot
 - 3d Scatterplot
 - Scatterplot matrices (SPLOMs)



- Cluster identification / verification
- NOT correlation



User Study?

- Idea:
 - Run a user study



- Measuring the human efficiency (time/error) in making a judgment (about clusters)
- to find out which vis technique works best

- Started piloting and found ...
 - What really matters is how the dataset looks like
 - Confounding variable = **dataset characteristic**

Automatic Data Study

- Idea:
 - Take a larger, well-chosen sample of real and synthetic datasets
 - Use recent cluster separation metrics to characterize datasets in terms of when which vis technique is best
 - Use the judgement itself, where the judgement is done by the metric



Centroid and Grid Metric

	Centroid Metric	Grid Metric
Sips et al. EuroVis'09	Distance Consistency	Distribution Consistency
Tatu el al. VAST 2009		2D Histogram Density Measure (2D-HDM)

- Found to be the current cutting edge [Tatu et al., AVI'10]
- Value between 0 (worst) and 100 (best) for a 2D Scatterplot



But...

- Again: Started piloting and found ...
 - The metrics judgment often does not align with human judgement
 - There are hidden assumptions about the data that lead to
 - FN: false negatives (low metric values for sth good)
 - FP: false positives (high metric values for sth poor)
 - Reason: **Data characteristics**

 Current metrics cannot be used for a reliable automatic judgment!

Qualitative Data Study

• Changing the route: **Data Characteristics!**



- Idea:
 - Do a manual qualitative inspection of datasets
 - Make human judgement
 - Compare to metric judgement
- Adaption of research question:
 - I. Identifying **data characteristics** that matter for cluster verification
 - 2. **Evaluating** cluster separation **metrics**

Methods

Methods I: Choosing the variables

• 75 datasets

- 31 real / 44 synthetic
- pre-classified / BUT: Keep non-classified in mind!
- **4 DR** techniques: PCA, MDS, RobPCA, t-SNE
- 3 VisEncs:
 - 2D Scatterplot
 - Interactive 3D Scatterplot
 - SPLOM
- 2 metrics: centroid and grid
 - Extension: 3D, SPLOM
 - Extension: Classwise









• 816 datasets instances (Scatterplots to look at)

Methods 2: Open and Axial Coding

- Coding = Method for qualitative data analysis from social science
- Open Coding = Figuring out codes and iteratively refine *codeset*
- Axial Coding = Analyze relation between codes and categorize *codeset*

• Open Coding

- Two researchers (Andrada and me): Multiple passes over 816 dataset instances
- Coding: Factors affecting visual class separability
- Coding: Failure cases of metrics:
 - ok dubious poor / classwise poor
 - if not ok: reason

Axial Coding

- Merging codesets
- Create and refine taxonomy of factors that matter

Qualitative Data Study: A side contribution of our project

- Qualitative studies (coding) have been used for
 - User analysis: Audio, video, notes, ...
 - Literature analysis
- We: Data analysis

- In Vis: Qualitative Data Study as an **inverse** to a User Study
 - User Study:
 Many users, few datasets
 - Qualitative Data Study:
 Few users, many datasets

Related Work



They

Scagnositcs [Wilkinson 2005]

- General factors
- Informal exploration
- Mathematical depiction
- Gestalt principles
 - Perceptual fundamentals

Clustering (ML)

Automation

We

- Cluster separability
- Systematic and rigorous
- Human perception

Operational guidance

Human judgment

Taxonomy

A taxonomy of visual cluster separation factors

	Within-Class Factors	Between-Class Factors
Scale	Count few :	Class/Point Countfew classes many pointsmany classes few points
	Size	Variance of Countsimilar \sim \sim different
	small o large	Variance of Sizesimilar $\bullet \longrightarrow$ $\bullet \longrightarrow$
Point Distance	Density \therefore sparse \therefore dense dense	Variance of Density similar 🔬 🔅 different
	Clumpiness equidistant uniformly one many dense equidistant random dense spot spots clumpy	Mixture random equidistant interwoven
	Outlier $mone \leftarrow \qquad $	Split contiguous
Shape	Shape $narrow \xrightarrow{Curvature} curvy$ $\downarrow \qquad \uparrow \qquad $	Variance of Shape similar \longrightarrow \longrightarrow different
uo	Centroid evocative	Inner-Outer Position non-existent
ositi		Class full partial Separation <u>overlap</u> adjacent separate distant

nfluence

Top-level structure

	Within-Class Factors Vari	iance Between-Class Factors
	few ··: →	Count many points few points
Scale	Size	Variance of Countsimilar \sim \sim different
	small o — — large	Variance of Size similar ••
e	Density $sparse \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \longrightarrow dense$	Variance of Density similar 🔅 🔅
nt Distan	Clumpiness equidistant uniformly one many dense equidistant random dense spot spots clumpy	Mixture random equidistant interwoven
Poi	Outlier $none \longleftarrow \vdots \qquad \dots \rightarrow many$	Split contiguous
Shape	Shape narrow Curvature curvy No V round	Variance of Shape similar 🔊 — So different
ion	Centroid evocative \longleftarrow \bigstar \longrightarrow misleading	$\begin{array}{c} \texttt{Inner-Outer} \\ \texttt{Position} \end{array} \texttt{non-existent} \longleftarrow \begin{array}{c} \bullet \\ \bullet \end{array} \longrightarrow \texttt{existent} \end{array}$
Positi		Class full partial Separation overlap adjacent separate distant

What can we do with it:

- Mapping **assumptions** of metrics onto taxonomy axes
- Mapping **datasets** onto taxonomy axes



Examples

Fisheries, real, MDS



Centroid: Overall: 29 Problem: FN Reason: "Stringy" shapes

Grid: Black & Red: ~70-80 Others: ~40-50 Problem: FN Reason: Adjacent strings

In terms of taxonomy ...



Gaussian, synthetic, MDS



Centroid:

Red: 77

Problem: **FP**

Reason: Bigger classes

overspread smaller one

In terms of taxonomy ...

		Within-Class Factors Var	iance Between-Class Factors
		Count few	Class/Point Countfew classes many pointsmany classes few points
	Scale	Size	Variance of Count similar $\bigotimes \bigotimes \longrightarrow \bigotimes $ different
		small O large	Variance of Sizesimilar $\bullet \longrightarrow \bullet$ different
	Ce	Density $sparse \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \longrightarrow dense$	Variance of Density similar 🔬 🎆 — — — 🎆 different
	t Distan	Clumpiness uniformly one many dense equidistant random dense spot spots clumpy	Mixture random equidistant interwoven
ence	Poin	Outlier $mone \leftarrow many$	Split contiguous split
Influe	Shape	Shape $\xrightarrow{\text{Curvature}} \text{curvy}$ $\xrightarrow{\text{Adoptog}} (1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,$	Variance of Shape similar in the similar is the second
ľ	ion	Centroid $evocative \longleftarrow misleading$	Inner-Outer Position non-existent \longleftarrow \longrightarrow existent
	Posit		Class full partial Separation overlap adjacent separate distant
	7		

HIV, real, t-SNE



Grid: Overall: 99 Problem: FP Reason: "Equidistant" points

In terms of taxonomy ...

		Within-Class Factors Vari	Between-Class Factors
ence		Count few \cdots many	Class/Point Countfew classes many pointsmany classes few points
	Scale	Size	Variance of Countsimilar \circledast \longrightarrow different
		small	Variance of Sizesimilar $\bullet \longrightarrow$ $\bullet \longrightarrow$ different
	Ce	Density sparse \therefore	Variance of Density similar
	nt Distan	Clumpiness equidistant uniformly one many dense equidistant random dense spot spots clumpy	Mixture random equidistant interwoven
	Poi	Outlier $none \longleftarrow \qquad \vdots \qquad many$	Split contiguous $\implies \longrightarrow \implies$ split
Influe	Shape	Shape $narrow \xrightarrow{Curvature} curvy$ $\downarrow \qquad \qquad$	Variance of Shape $similar \longrightarrow S different$
	ion	Centroid $evocative \longleftarrow$ \bigstar \longrightarrow misleading	Inner-Outer Position non-existent
	Posit		Class full partial Separation overlap adjacent separate distant

Summary of results

High-level results



Summary: Mapping assumptions onto taxonomy axes



Centroid:

Mapping assumptions onto taxonomy axes



Grid: Mapping assumptions onto taxonomy axes



Conclusion

Main Contribution A Taxonomy of Visual Cluster Separation Factors



A Taxonomy of Visual Cluster Separation Factors



Michael SedImair, Andrada Tatu, Tamara Munzner, Melanie Tory University of British Columbia, Vancouver, CA Slides: <u>http://www.cs.ubc.ca/~msedI/talks/konstanz2012.pdf</u>



Appendix

Adopting the Grid Size?





- Developers [Sips'09]: "insensitive to grid size changes"
 - No:Test with 50 instances insensitive only in 16%
- Static: Does not work well with different #points
- Dynamic: No straight forward automatic way

