



# An Empirical Study on the Quality of Algorithms for Dimensionality Reduction and Visualization of High-Dimensional Data

#### **Bachelor's Thesis Final Presentation**

Lyubomir Stoykov 21.03.2015

Software Engineering for Business Information Systems (sebis) Department of Informatics Technische Universität München, Germany

wwwmatthes.in.tum.de

## Structure



- 1. Motivation & Context
- 2. Thesis requirements
- 3. Approach
- 4. Evaluation Results
- 5. Conclusion
- 6. Future Work
- 7. Web Application Architecture & Live Demo
- 8. Timeline

#### **Motivation & Context**



**Enable visual analysis & pattern recognition.** 

Exponentially faster data processing with less resources.

Dimensionality reduction techniques mostly used "as is", without evaluation.

No tool for easy appliance of dimensionality reduction techniques.

## Thesis Requirements



#### First requirement – answer research questions

- 1. What kind of dimensionality reduction techniques do exist?
- 2. What are some other means to explore and visualize patterns in high-dimensional data sets?
- 3. In what extent do dimensionality reduction techniques and means for pattern exploration in high-dimensional data sets reveal information about the underlying data?
- 4. What results do they yield if used to visualize categorical data after a user-defined mapping f: String → Number has been applied?

Second requirement – develop a web application to apply dimensionality reduction

## Approach



#### Overview of dimensionality reduction techniques → select best for evaluation

- 1) Linear (PCA)
- 2) Nonlinear, local properties (LLE)
- 3) Nonlinear, global properties (Isomap)
- 4) Recent algorithm (T-SNE)

Other techniques for high-dimensional data visualization → select best for evaluation

1) Cluster Analysis (DBSCAN)

#### **Evaluation**



#### **Cluster Analysis**

- 3 Real-world datasets tested Business data, MNIST, Iris.
- 1) Almost no clusters identified.
- 2) Rand Index nearly 0.
- 3) Applied also PCA prior to reduce curse of dimensionality no improvement.

#### **Implications**

- 1) Real world data not really clustered.
- 2) Cluster analysis doesn't deal well with real world data.
- 3) Cluster analysis doesn't deal well with high-dimensional data.



#### **Dimensionality Reduction – Metric data**

#### 8 Data sets

1. Assign scores (Local vs Global structure)

```
1 - Very Poor Visualization – ....
```

2. Average scores

#### Final averaged scores

PCA - 4.2

LLE - 3.11

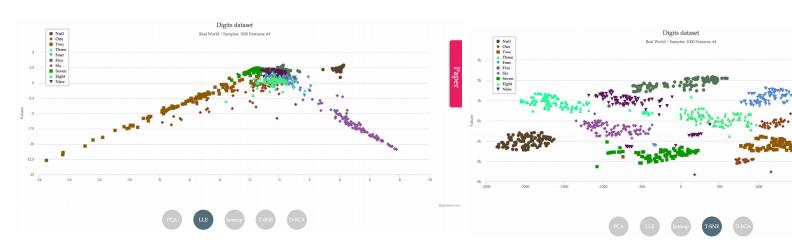
Isomap - 4.11

T-SNE - 3.88



#### **Dimensionality Reduction – Metric data.**

#### **Example: MNIST dataset – Local vs Global Structure**



Local structure: well preserved, digits

grouped in clusters

Global structure: very poor

representation, no distances between

sub-manifolds

Score: 3 (or 2?)

Local structure: very well preserved,

digits grouped in clusters

Global structure: very good

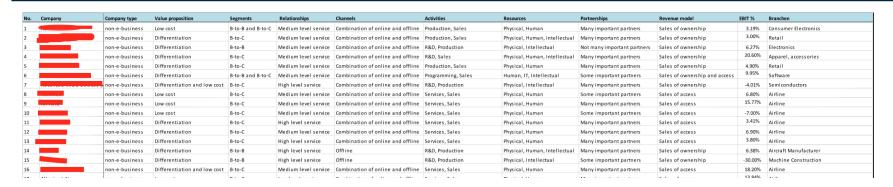
representation, clear separation

between sub-manifolds

Score: 5



#### **Dimensionality Reduction – Categorical data**



#### Ask 11 students between 20 and 26 y.o.

- 1) PCA 4. Patterns and clusters recognized.
- 2) Isomap 4.2. Patterns and clusters recognized.
- 3) LLE 3. Patterns and clusters recognized, but also a lot of inconsistencies.

4) T-SNE – 2. A lot of inconsistencies.



### **Dimensionality Reduction – Categorical data**

#### **Example: Isomap**

	62:2016 6:3 Empirical study on the quality of dimensionality reduction algorithms for visualization of categorical data.  6. How would you classify the visual representation?  Mark only one ovail.	What I would categorise is all
pirical study on the quality of dimensionality uction algorithms for visualization of	Very poor - Most of the patterns that I recognize contradict reality (for instance, companies that are similar such as BMW and Audi are not in the same neighborhood or companies that are very different such as Turkish Airlines and Spotify are grouped together)	What I would categorize is the I would classify Spottify and
	Poor - I can't really recognize similarities, patterns or structures	Her to the group of ouline
egorical data	Somehow good, somehow poor - I can recognize some patterns compliant with	P I Alvania
pal of this evaluation is to explore suitability of dimensionality reduction algorithms for nizing patterns in complex categorical data structures by visual analysis of the reduced in order to achieve this goal, dimensionality reduction was applied on high sional companies data. After that, the reduced, low - dimensional representation has	reality (for instance, companies that are similar such as BMW and Audi are in the same neighborhood or companies that are very different such as Turkish Airlines and Spotify are far away from each other), but I also recognize a fair amount of inconsistencies that contradict reality	Services such as Netflix, ebay. Over all I think that similar
visualized in a two - dimensional plot. survey people are asked to elaborate on their interpretation of the visual	<ul> <li>Good - Most of the patterns that I recognize are compliant with reality, but there are some visible inconsistencies</li> </ul>	oupanies are ground together
entation. After that their answers are examined for compliance with the patterns igated in the initial data set.	Very good - Most of the patterns that I recognize are compliant with reality without visible inconsistencies	without much information lost.  t is also easy to identify Clusters.
	7. Did 3D representation enhance visual analysis?	
lgorithm	Mark only one oval.	t is also easy to identify Clusters
lark only one oval.	Yes	S - (4.5143.
PCA	No	The second secon
. LIE	Neutral	
$\exists$		
T-SNE	Please explain your interpretation of the visualization (Patterns, Structures,     Outliers, Relationships, Inconsistencies, Accuracies)	
	The industry feelds	
ickname	are clearly separated	
	ors I would expect by	and the relative countries and the relative contribution of the second
Wser 7		
	the services they provided	
ield of study	(food, cosmolics, health care, transport of eparated in 3	
BWL	what I found contradictory is that airlines	
	What I tound controlled by the other groups of	The state of the s
ge	powered by companies that I've could see and belong to	
O C	Google Forms this same group are so there, hardware electronic	8
<i>V</i>		
ender	1 (Airlines)	
lark only one oval.	97	and the second s
	Pharmadiculicals (8 Productions It	
Male	The state of the s	
Female	1 100 and 100 Marian	Ч
	Control Cosmetics cheer cleaning the	
	Loon Alden	<b>d</b>

#### Conclusion



#### High-Dimensional Real-World Data Visualization... (Research Questions 3 & 4)

1. Use Cluster Analysis?: No

2. Use Dim. Reduction for Metric Data?: Yes, with a grain of salt

3. Use Dim. Reduction for Categorical Data?: Yes, but see 2)

4. 3D visualization better than 2D?: No

5. Best algorithm?: PCA

## **Future Work**



**Dimensionality Reduction for Categorical Data** 

**Combination of Techniques?** 

**Supervised vs Unsupervised?** 

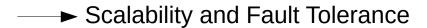
**Curse of dimensionality?** 

## Web Application Architecture & Live Demo





- 1) Integrate Meteor and Scikit
- 2) Most computationally expensive calculations on server-side





























**MathJax** 

## Timeline



START: 10.15.2015

OCTOBER

**NOVEMBER** 

**DECEMBER** 

**JANUARY** 

**FEBRUARY** 

15.10.2015

**Literature Review** 

15.11.2015

App Concept

20.12.2015

**App Implementation** 

20.12.2015

**Thesis Outline & Concept** 

15.01.2016

**Evaluation and Thesis completion** 

15.02.2016



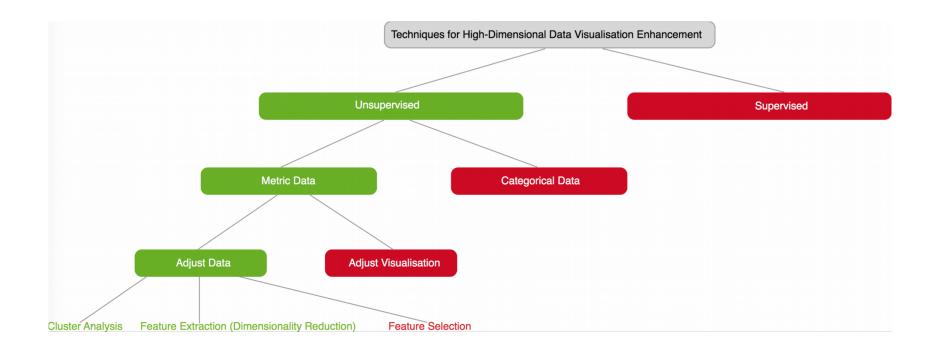
# Thank you



# Backup

## Scope





## Map Function



```
# visual_mapping = {
    "Company type":{
     "e-business":"1",
     "non-e-business":"75"
   "Segments":{
     "B-to-B":"1",
     "B-to-B and B-to-C":"75",
     "B-to-C":"150"
   },
    "Value Proposition":{
     "Low cost":"1",
     "Differentiation and low cost": "75",
     "Differentiation":"150"
    "Relationships":{
      "Low level service":"1",
     "Medium level service": "75",
     "High level service":"150"
   "Channels":{
      "Offline":"1",
      "Combination of online and offline": "75",
     "Online":"150"
#
   },
# }
```