Visualization through Textualization

Ye Zhao Kent State University



Motivation

- Gap exist between what users want and what information visualization can provide
 - Visual metaphors unfamiliar to novice users
 - Cluttered views
 - Complicated interactions
- Interactive visual exploration of large, complex datasets can only be performed by a small number of visualization experts and domain experts collaborating with them



Power of Text

- Text is an intelligent tool
 - Text is used to refer to something that carries its interpretation within itself
 - it can render abundant information with its placement, font, color, feeling, and style
- Text helps people convey information and form opinions promptly









Textualization

- "to put into text" of oral epics, blogs, emotions ...
- Textualization: externalizing and refining associated contextual information of the original data
- A variety of abstract data can be processed and analyzed in semantic-rich forms



DE GRUYTER

Leard Honko (Ed.) TEXTUALIZATION OF ORAL EPICS

THENDS IN LINGUISTICS. STUDIES AND MONOGRAPHS [TILSH]

Visualization through Textualization

- Develop visualization techniques to fully utilize the power of text
 - Expose semantics and contextual information associated with the original data
 - Discover patterns and trends of the data with text mining techniques
 - Explore discoveries with interactive visualization



Our Scheme





Stage 1: Data processing

- A multilevel textualization scheme to manifest the contextual and semantic information of the original data
 - Semantic transformation: convert data values into textual descriptions from domain knowledge
 - User Input: incorporate users' input of annotation, tagging, comments
 - Domain ontology: use the vocabulary to denote given data with its types, implicit properties, and interrelationships



Multilevel Textualization: An example

• Taxi trajectory data: two GPS samples



Challenge: Handling Different Data Types

- How can we handle different types of data with the textualization scheme?
- How can we preserve the data features and relationships?



Geospatial data

- A variety of "georeferenced" information
 - demographic (e.g. census and real estate)
 - environmental (e.g. weather and climatological records)
 - geology (e.g. land features)
- Geospatial-temporal data
 - severe weather systems
 - population movement of plant and animal species
 - epidemics of flu
 - human and vehicle mobility trajectories

Images from the Web Copyright belongs to original owners

• Map geospatial positions to a meaningful text representation







Tabular and Relational Data

- A fundamental type incorporates numeric, ordinal, categorical, and textual variables, as well as unstructured metadata
 - Each data record can be turned into a document
 - A value in the record is converted into a keyword
- Relational datasets may contain a set of tables
 - Tables can be processed through controlled denormalization



Images from the Web Copyright belongs to original owners





Challenge: Interactively Incorporating User Input

- How can we employ user input for a large dataset?
- How can we design an effective interface for users?
- Users tag or comment over instance data, then apply to remaining data
- Systems developed in semantic web and information extraction
 - Mostly dedicated to unstructured text
- A good visual interface needed

Images from the Web Copyright belongs to original owners





Challenge: Organizing Data

- How can we store and organize data?
- How can we make the organization efficient for text analysis and interactive visualization?
- To support:
 - fast query and retrieval by text
 - query and retrieval by original data items
 - fast computation promoting interactive visualization



Stage 2: Text Data Processing

- A set of techniques in natural language processing, statistical, and machine learning that extract and analyze the information of textual data
 - Filtering
 - Clustering
 - Classification
 - Query and Search



Studying Textualized Data

- Find information patterns from textualized data
 - Utilize the term (keyword) vector model to represent data items and datasets
- Fast computing and easy user interaction



Visualization

- Significant discoveries can be expressed using intuitive visualizations and textual explanations written in natural language
- Challenge 1: Expand existing text visualization techniques
- Challenge 2: Combine text visualization with abstract data visualization



Case study: Explore Taxi Trajectory with Semantic Transformation

- Convert each trajectory as a *document* consisting of the taxi-traversed streets
- Enable analysis of massive taxi datasets as document corpora with text mining tools
- Use LDA Topic modeling to infer hidden patterns of moving taxi populations
- Visualization based on the *taxi topics*



Shenzhen Data

- Daily trajectories of 21,360 taxis in Shenzhen
 - A big city in southern China bordering with Hong Kong
 - Fifteen million residents in a condensed area
 - taxis are a major means of passenger transportation
- Each taxi reports nearly three thousand GPS sample positions per day
- Each sample consists of taxi plate, time, status, speed, direction, and latitude and longitude
- A total of 59,087,230 samples recorded in one day.

Shenzhen plan for city-wide SEZ





Images from the Web Copyright belongs to original owners



Тахі Торіс

• Reveal typical traveling patterns of city cabs



- The topics approximate the city's functional regions
 - A large portion of taxis can accomplish their movement inside a district
- Topics are more than geometrical divisions
 - An airport highway (Green) is an important component of several topics, connecting different regions



Eight Topics





PCP-based Analysis of Topic 2

Minimum Values -

+ Num

+ Number of Roads bigger than Minimum Value: 2 + Change!

Selected Roads(Click the lines in the plot!)

2.561061

BinheAve,BeihuanAve, G 2 0 5,

Plot roads on map Clear selected roads





Topics and Trajectories







Visualize Street Changes







(c) Topic2: 6pm-9pm

 Visualizing disappearing (brown) and emerging (orange) streets



Case Study: Explore Categorical Datasets as Documents

- Map each record in a categorical dataset to a document represented by a bag of categories
- Convert a categorical dataset into a document corpus
- Apply text-based cluster analysis (LDA) to discover subspace clusters of textual category values
- Use associate rule mining discovers optimal risk rules describing multivariate relationships in the textualized topical subspaces



MovieLens dataset

- 1,000,209 data items representing rankings from 6,040 users for 3,883 movies
- Use word clouds, word tree, and fingerprinting, are then used to visualize the rules and data items for interactive knowledge discovery
- Topic Cards:

https://www.youtube.com/watch?v=W2Kt7WKIMTI



Mushroom Data



Working on Evaluation

- The effectiveness and efficiency of algorithms of textualization methodology?
- Compare them with other visualizations?
- User study over different applications



Use in More Visualization Tasks?



KENT STATE.

Conclusion

- Text can be of interest in abstract data analysis and visualization
- Text analysis tools are enabled
- Text visualization to be compactly integrated with existing visualizations
- Specific approaches for different data types and applications



Thanks!

- Collaborator
 - Jing Yang, UNC-Charlotte
- Acknowledgements
 - Wei Chen, Zhejiang University
 - George Chen, Maogong Zheng, Shenzhen Institute of Advanced Technologies
 - Blake Stringer, Aeronautics Program, Kent State U.
 - Xiaoling Pu, Dept. of Finance, Kent State U.
 - Davis Sheets, Ding Chu, Xiaoke Huang, Wendy Wu, Shamal Aldohuki, Farah Kamw, Yang Chen, Yueqi Hu, Chong Zhang, Scott Barlowe

