

Business Intelligence & Process Modelling

Frank Takes

Universiteit Leiden

Lecture 2 — BI & Visual Analytics

Recap

- **Business Intelligence:** anything that aims at providing actionable information that can be used to support business decision making
 - **Business Intelligence**
 - **Visual Analytics**
 - Descriptive Analytics
 - Predictive Analytics
- Process Modelling (April and May)

Business Intelligence

Business Intelligence goals

- Operational intelligence
- Corporate governance
- Risk assessment
- Compliance
- Auditing
 - Sarbanes-Oxley (SOX) — role of IT in corporate governance



Management Approaches in BI

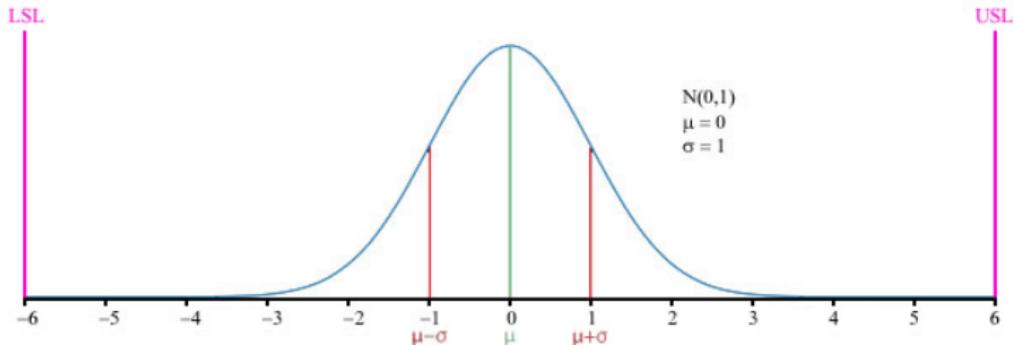
- **Continuous Process Improvement (CPI)**: ongoing effort to improve products, services or processes
 - Incremental improvements vs. Breakthrough improvements
 - Evaluate based on efficiency, effectiveness and flexibility
- **Total Quality Management (TQM)**: improve processes up to the microscopic level, focussing on meeting customer demands and realizing strategic company goals, e.g., **Six Sigma**

Six Sigma

- Originally developed by Motorola in the early 1980s
- Minimize Defective Parts per Million Opportunities (DPMO)
- Mean μ and standard deviation σ

Quality level	DPMO	% broken	% OK
One Sigma	691.462	69	31
Two Sigma	308.538	31	69
Three Sigma	66.807	6,7	93,3
Four Sigma	6.210	0,62	99,38
Five Sigma	233	0,023	99,977
Six Sigma	3,4	0,00034	99,99966

Normal distribution



DMAIC approach

- **Define** the problem and set targets,
- **Measure** key performance indicators (KPI's) and collect data,
- **Analyze** the data to investigate and verify cause-and-effect relationships,
- **Improve** the current process based on this analysis,
- **Control** the process to minimize deviations from the target.

Key Performance Indicators

- **KPI:** measure, variable or metric to analyze the performance of (part of) an organization
- Strategic goals → Measurable variables
- **SMART**
 - Specific
 - Measurable
 - Acceptable
 - Realistic
 - Time-sensitive

KPI examples

- Operational: increasing market share by 10%
- Financial: increase profit by 10%
- Sales: obtain 10 new customers
- Human resources: attract 10 new sales officers that are part of the world's top 1% in the field
- Customer support: forward no more than 10% of the support calls to second line

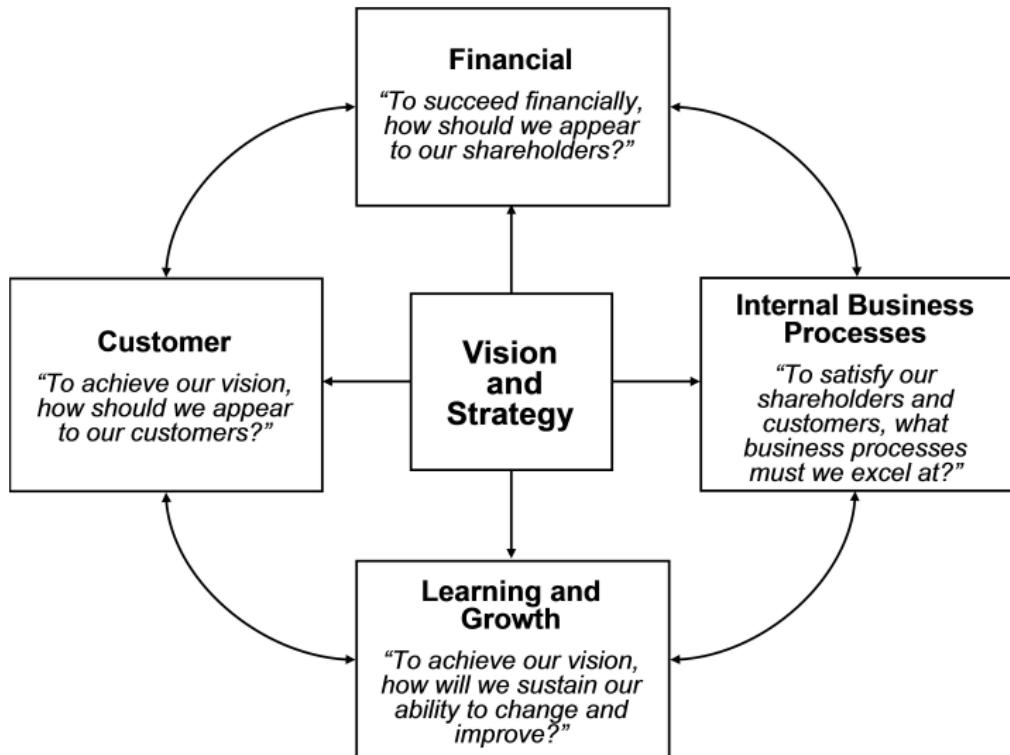
BI in practice

- Codeless reporting
- Instant querying
- Rich visualization
- Dashboards (“Management cockpits”)
- Scorecards

Balanced Scorecards

- R. Kaplan, D. Norton, The balanced scorecard: measures that drive performance, *Harvard business review* 83(7): 172–180, 2005.
- Goal: align business activities to the vision and strategy of the organization
- Financial and nonfinancial goals
- Monitor a relatively small number of summative indicators

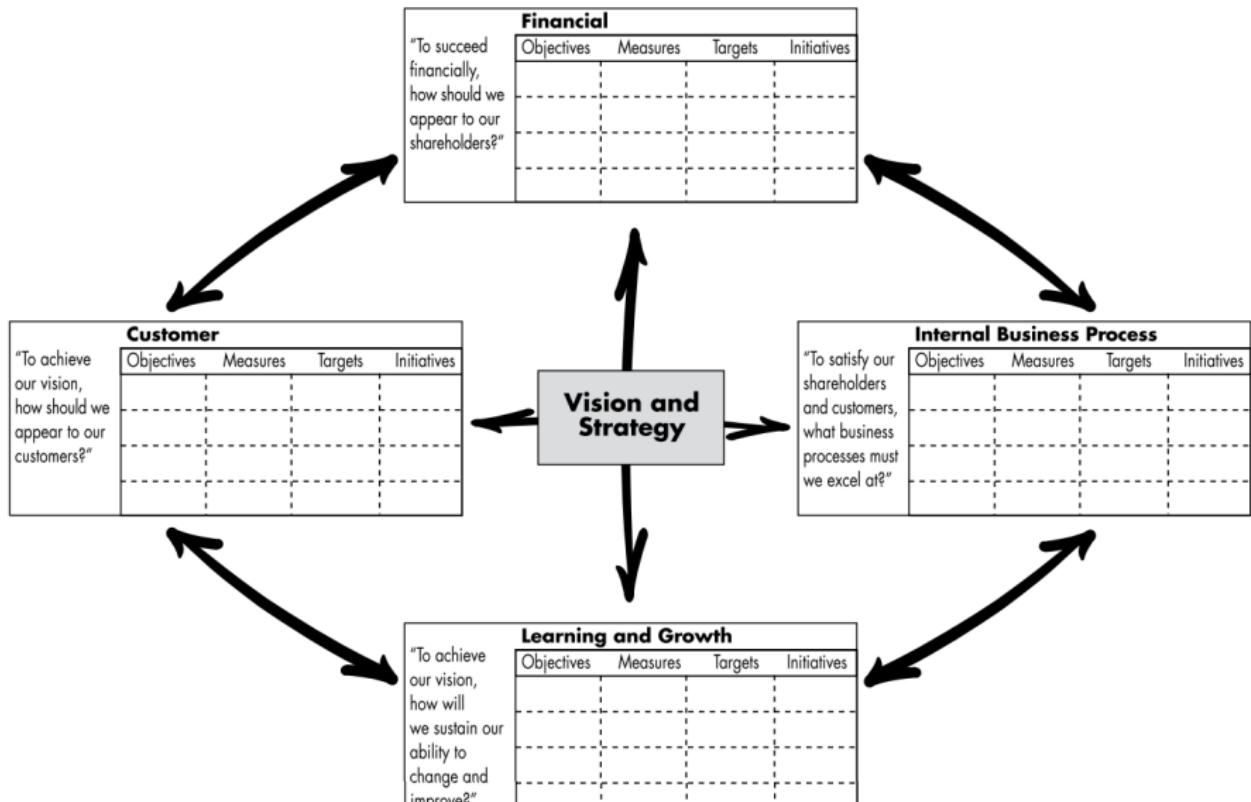
Balanced Scorecard



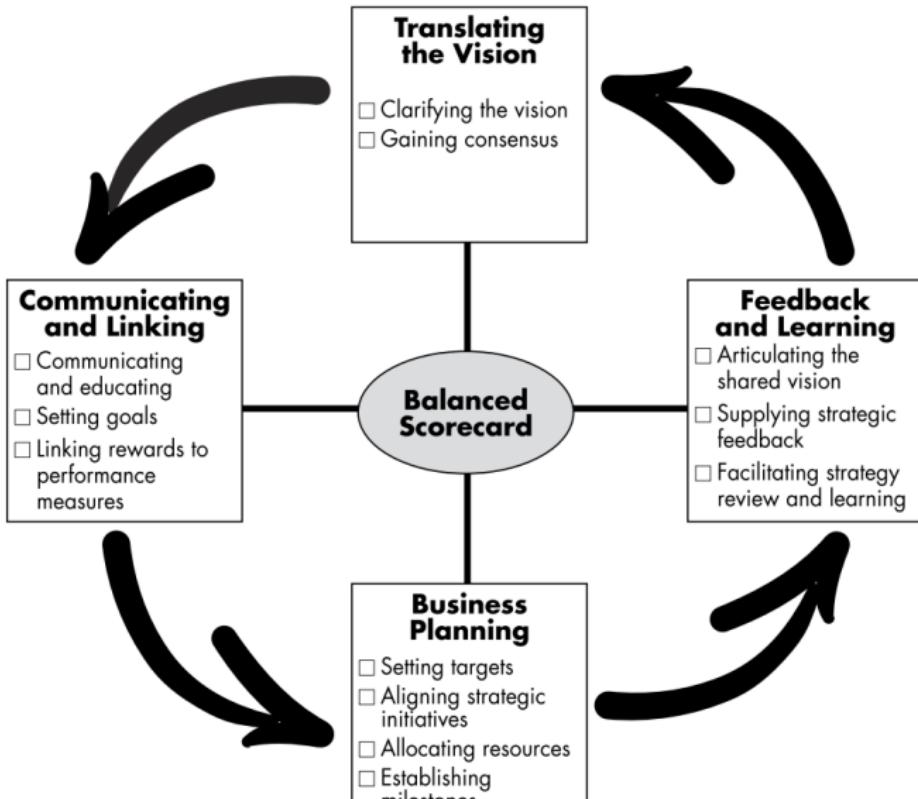
Balanced Scorecard

- Four perspectives:
 - Financial
 - Customer
 - Processes
 - Learning and Growth
- Four elements per perspective
 - Objectives
 - Measures
 - Targets
 - Initiatives

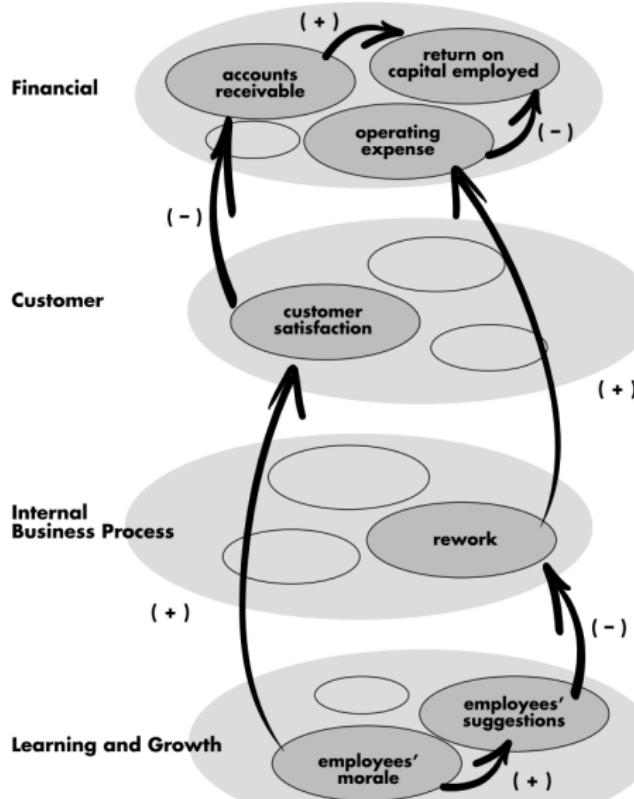
Balanced Scorecard Perspectives



Balanced Scorecard Processes



Relations between Perspectives



Some terms . . .

- **Business Activity Monitoring (BAM)**: insight in operational status and events of a business
- **Complex Event Processing (CEP)**: monitor events and react immediately if a pattern occurs
- **Corporate Performance Management (CPM)**: measuring the (financial) performance of a process or organization

Some systems . . .

- Enterprise Resource Planning (ERP) Systems
- Enterprise Information Systems (EIS)
- Business Information Systems (BIS)
- Management Information System (MIS)
- Executive Information System (EIS)

ETL

- Extract data from source systems: generate dumps, exports, etc.
- Transform data: aggregating, linking, sorting, joining, etc.
- Loading data into target system into desired (reporting) format

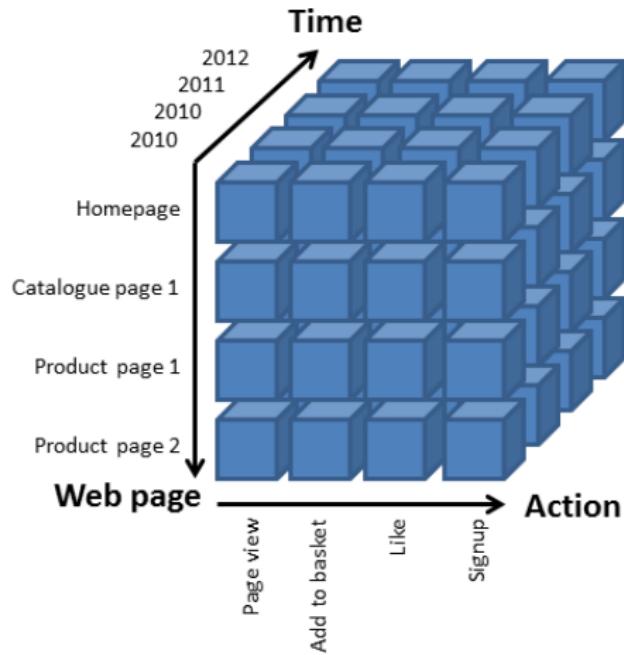
OLAP

- **OnLine Analytical Processing (OLAP)**
- Given a data table with n attributes:
 - Dimensions of an $(n - 1)$ -dimensional cube represent $n - 1$ attributes of the data
 - Value in a cell of the cube represents the remaining attribute
- Use a **slice** or **dice** to get the desired information
- Suitable for, e.g., star schema data

OLAP Example

- Example: website visitor logs, storing:
 - 1 Time
 - 2 Web page
 - 3 Action
 - 4 Conversion

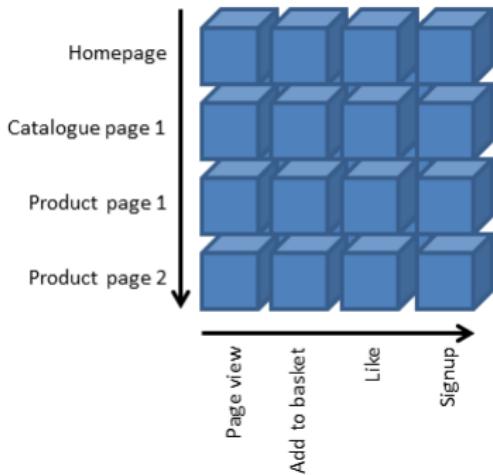
OLAP Cube Example



<http://snowplowanalytics.com>

OLAP Cube Example Slice

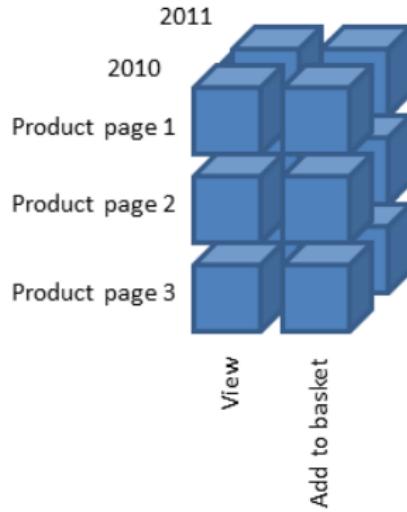
2010 slice



<http://snowplowanalytics.com>

OLAP Cube Dice

Product action dice



<http://snowplowanalytics.com>

OLAP Formalized

- **OnLine Analytical Processing (OLAP)**
- Given a data table D with 4 attributes W, X, Y and Z
- An OLAP cube can be characterized as a function
 $f : (X, Y, Z) \rightarrow W$
- An example of a **slice** is a function $g : (Y, Z) \rightarrow W$
- Given subsets $X' \subseteq X$ and $Z' \subseteq Z$ a **dice** is a function
 $h : (X', Y, Z') \rightarrow W$

Break?

Visual Analytics

What is Visualization?

- Intuition: data is more than its raw bits and bytes
 - **Visualization:** making something visible to the eye (Oxford dictionary)
 - *All visualizations share a common “DNA” — a set of mappings between data properties and visual attributes such as position, size, shape, and color — and customized species of visualization might always be constructed by varying these encodings.*
- Heer et al., A Tour through the Visualization Zoo, CACM 53(6): 59–67, 2010.
- **Visual Analytics:** knowledge discovery (DIKW) based on visualization

What is Visualization?

- **Data properties:** attributes of (groups of) data objects
Name; Age; City

What is Visualization?

- **Data properties:** attributes of (groups of) data objects
Name; Age; City
Frank; 28; "Niels Bohrweg 1, Leiden"

What is Visualization?

- **Data properties:** attributes of (groups of) data objects
Name; Age; City
Frank; 28; "Niels Bohrweg 1, Leiden"
- **Visual attributes:** e.g., position, size, shape, label, color, etc.
Label; Size; Position
- **Visualization:** mapping data properties to visual attributes

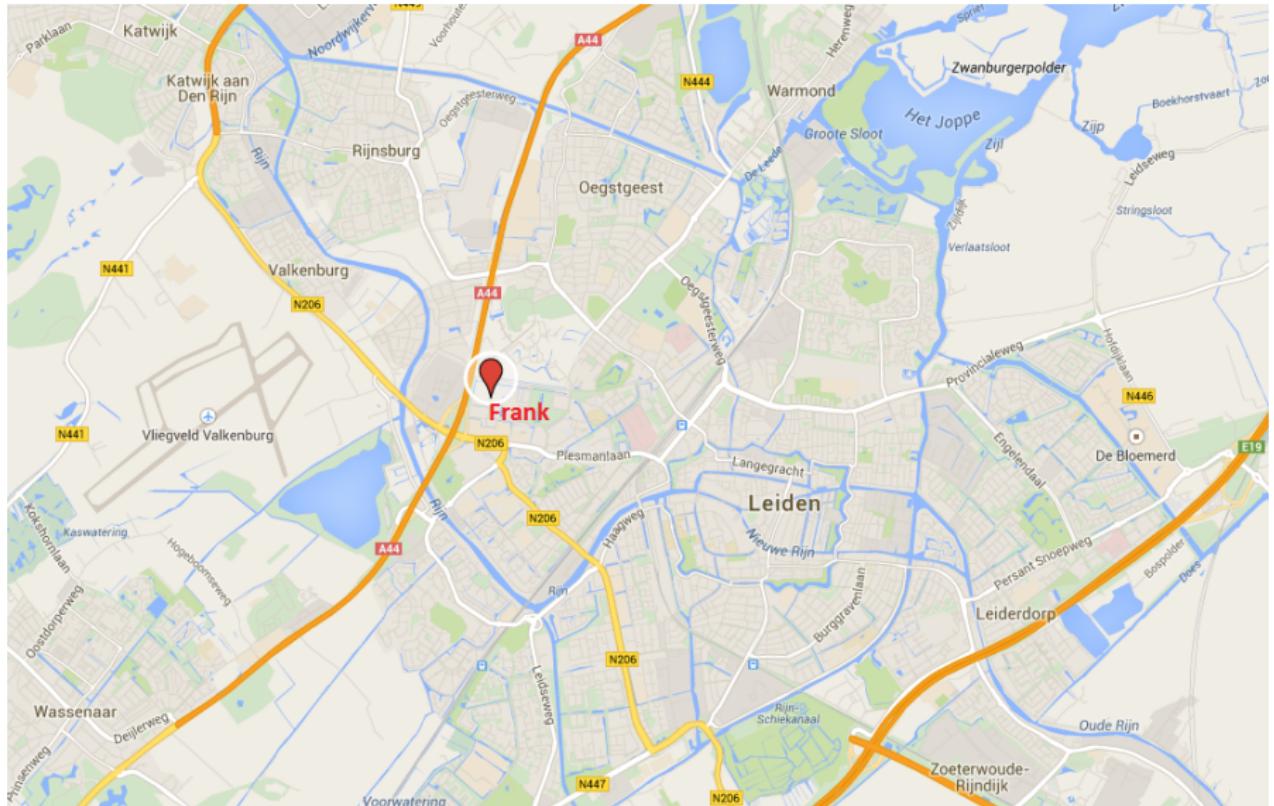
What is Visualization?

- **Data properties:** attributes of (groups of) data objects
Name; Age; City
Frank; 28; "Niels Bohrweg 1, Leiden"
- **Visual attributes:** e.g., position, size, shape, label, color, etc.
Label; Size; Position
- **Visualization:** mapping data properties to visual attributes
Name → Label
Age → Size
City → Position

What is Visualization?

- **Data properties:** attributes of (groups of) data objects
Name; Age; City
Frank; 28; "Niels Bohrweg 1, Leiden"
- **Visual attributes:** e.g., position, size, shape, label, color, etc.
Label; Size; Position
- **Visualization:** mapping data properties to visual attributes
Name → Label
Age → Size
City → Position
"Frank", $\log_2(28)$, (52.1603216, 4.4939262)

What is Visualization?



Why visualization?

Which of the following are the most important business benefits that your organization seeks to gain from deploying data visualization and visual analysis technologies? (Please select all that apply.)



SAS, Data Visualization: Making Big Data Approachable and Valuable, 2014

Why visualization?

Top Benefits of Data Visualization Tools

Improved decision-making **77%**

Better ad-hoc data analysis **43%**

Improved collaboration/
information sharing **41%**

Provide self-service
capabilities to end users **36%**

Increased ROI **34%**

Time savings **20%**

Reduced burden on IT **15%**

SAS, Data Visualization: Making Big Data Approachable and Valuable, 2014

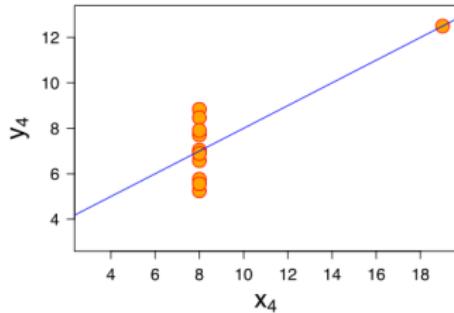
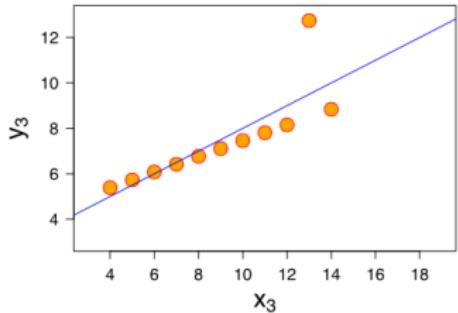
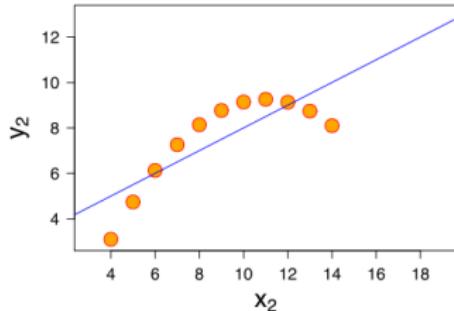
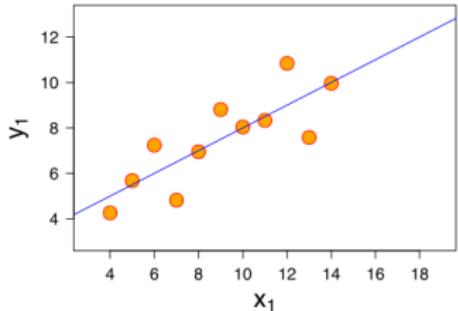
Visualization theory

- Discrete vs. continuous data
- Categorical vs. quantitative data
- Mean or median?
- Variance?
- Correlations? Regression?
- Normal distribution or power law?
- The **correct visualization** depends on the data itself!

Listen to the data to ...

- Catch mistakes
- See patterns
- Find violations of statistical assumptions
- Generate hypotheses
- Do outlier detection

Anscombe's quartet

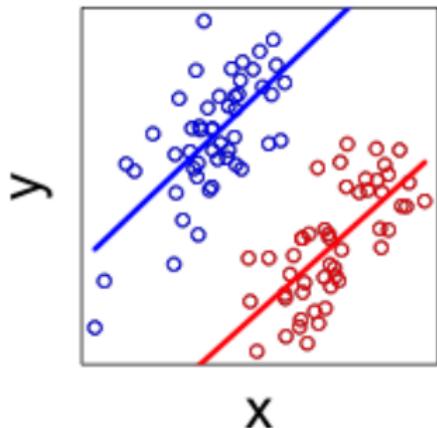
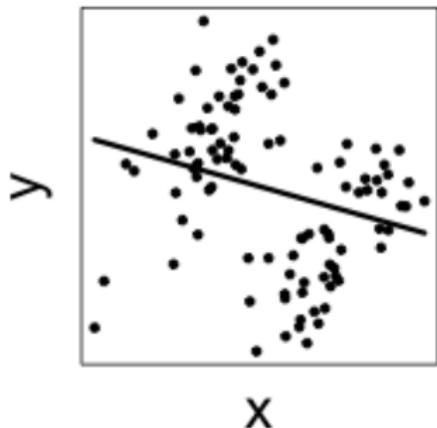


F. J. Anscombe. Graphs in Statistical Analysis. *American Statistician* 27 (1): 1721, 1973.

Anscombe's quartet

Property	Value	Accuracy
Mean of x	9	exact
Sample variance of x	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of y	4.125	plus/minus 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.0 + 0.50x$	to 2 decimal places

Simpson's Paradox



Visualization Quality

- When is a certain visualization “**good**”?

Visualization Quality

- When is a certain visualization “**good**”?
- “Proper mapping of data properties to visual attributes” ?

Visualization Quality

- When is a certain visualization “**good**”?
- “Proper mapping of data properties to visual attributes” ?
- The number of data properties (variables) that is visualized?

Visualization Quality

- When is a certain visualization “**good**”?
- “Proper mapping of data properties to visual attributes” ?
- The number of data properties (variables) that is visualized?
- The number of visual attributes that is utilized?

Visualization Quality

- When is a certain visualization “**good**”?
- “Proper mapping of data properties to visual attributes” ?
- The number of data properties (variables) that is visualized?
- The number of visual attributes that is utilized?
- Aesthetics?
- ...

Visualization Quality

- When is a certain visualization “**good**”?
- “Proper mapping of data properties to visual attributes” ?
- The number of data properties (variables) that is visualized?
- The number of visual attributes that is utilized?
- Aesthetics?
- ...
- Hard to answer objectively!

Infographic of infographics

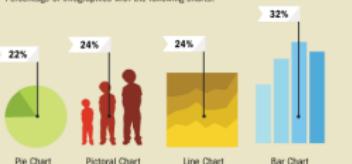
INFOGRAPHIC OF INFOGRAPHICS

Data visualization is a popular new way of sharing research. Here is a look at some of the visual devices, informational elements, and general trends found in the modern day infographic.

DESIGN

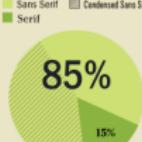
CHART STYLE

Percentage of infographics with the following charts:



FONT

Sans Serif Condensed Sans Serif
Serif



KEY INFO

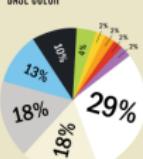
Percentage of infographics with key:



Average number of symbols per key: 5.1

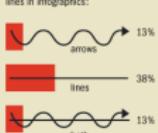


BASE COLOR



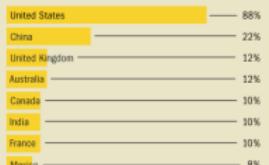
NAVIGATIONAL ICONOGRAPHY

Frequency of arrows & connecting lines in infographics:



CONTENT

COUNTRIES FEATURED

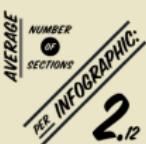


THEME

Relative popularity of different infographic themes:



SECTIONS



CREDITED SOURCES

Average number of sources per infographic: 2.29



TITLE

Average number of words per infographic title: 4.36

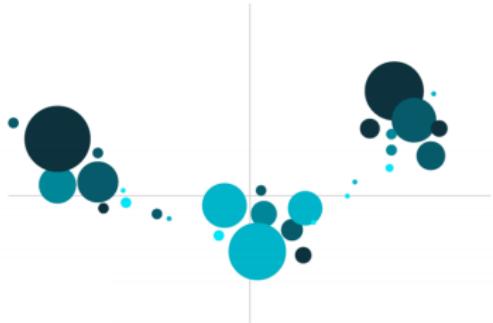
"RICHEST AND POOREST AMERICAN NEIGHBORHOODS"

Visualization Metaphors

- Important is Big
- Happy is Up
- More is Up
- Categories Are Containers
- Organization is Physical Structure
- Similarity is Closeness
- Control is Up

Visualization Metaphors

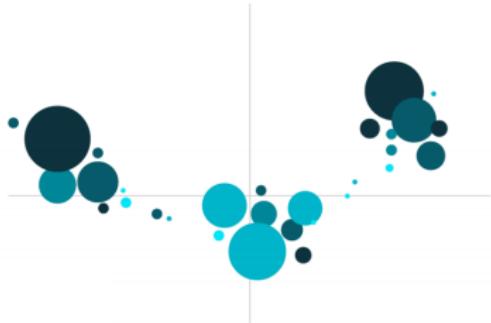
- Important is Big
- Happy is Up
- More is Up
- Categories Are Containers
- Organization is Physical Structure
- Similarity is Closeness
- Control is Up



http://www.bostondatafest.com/wp-content/uploads/2013/11/big_data_viz.pdf

Visualization Metaphors

- **Important is Big**
- Happy is Up
- More is Up
- Categories Are Containers
- Organization is Physical Structure
- **Similarity is Closeness**
- Control is Up



http://www.bostondatafest.com/wp-content/uploads/2013/11/big_data_viz.pdf

Visualization Metaphors

- Important is Big
- Happy is Up
- More is Up
- Categories Are Containers
- Organization is Physical Structure
- Similarity is Closeness
- Control is Up



http://www.bostondatafest.com/wp-content/uploads/2013/11/big_data_viz.pdf

Visualization Metaphors

- Important is Big
 - Happy is Up
 - More is Up
 - Categories Are Containers
 - **Organization is Physical Structure**
 - **Similarity is Closeness**
 - **Control is Up**



http://www.bostondatafest.com/wp-content/uploads/2013/11/big_data_viz.pdf

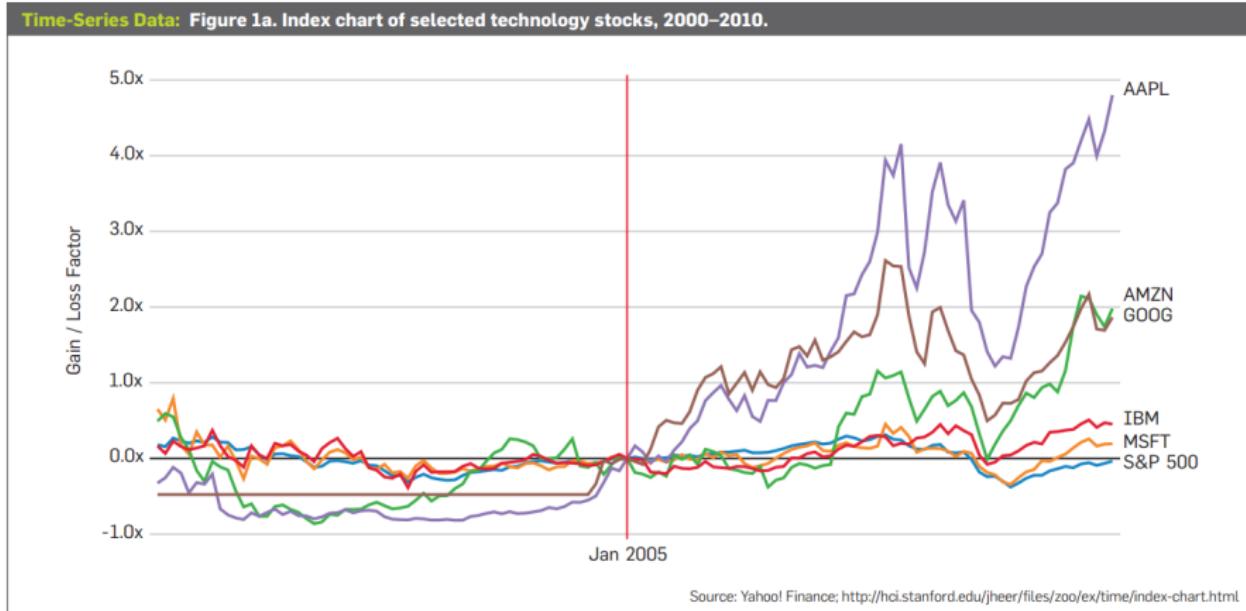
Examples

- Time-series data
- Statistical data
- Geographical data
- Hierarchical data
- Network data

Heer et al., A Tour through the Visualization Zoo, CACM 53(6): 59–67, 2010.

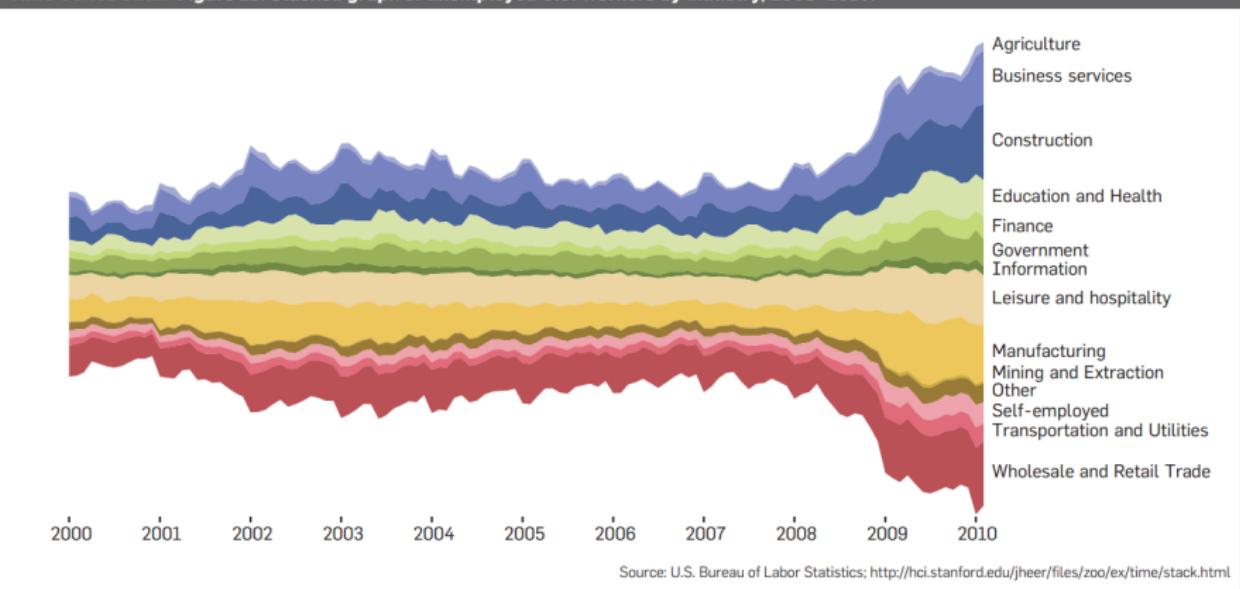
Index chart

Time-Series Data: Figure 1a. Index chart of selected technology stocks, 2000–2010.



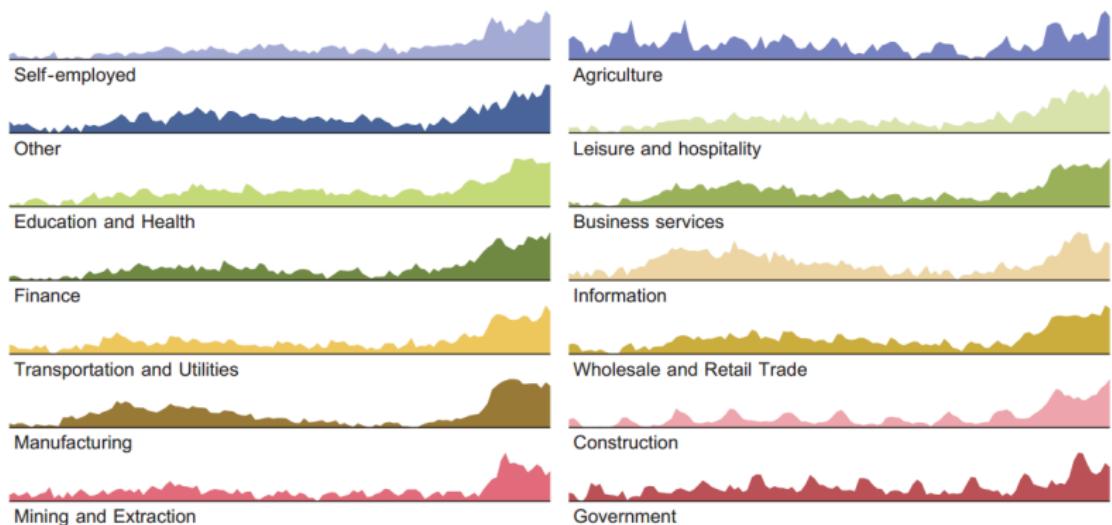
Stacked graph

Time-Series Data: Figure 1b. Stacked graph of unemployed U.S. workers by industry, 2000–2010.



Small multiples

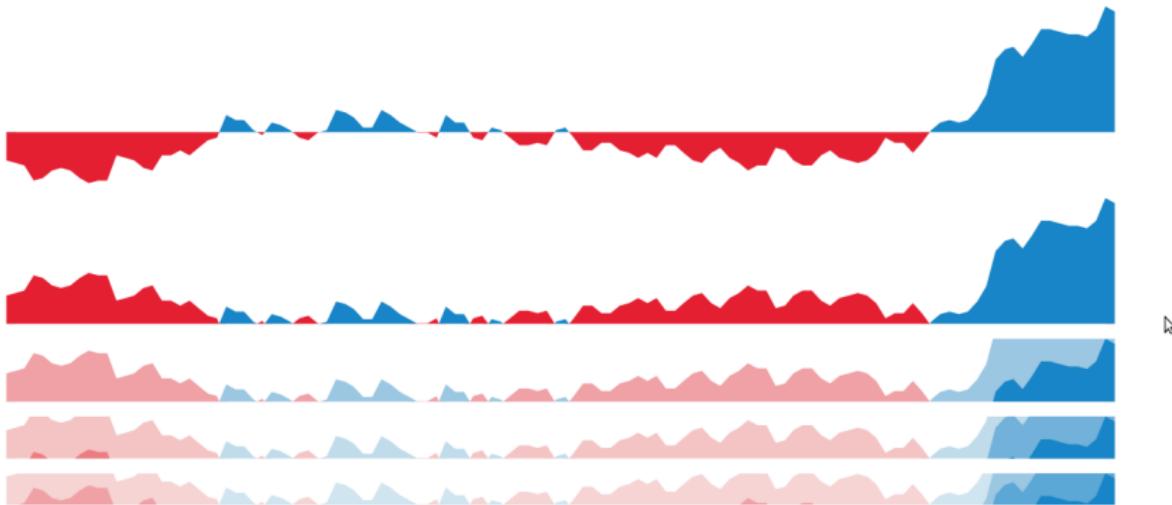
Time-Series Data: Figure 1c. Small multiples of unemployed U.S. workers, normalized by industry, 2000–2010.



Source: U.S. Bureau of Labor Statistics; <http://hci.stanford.edu/jheer/files/zoo/ex/time/multiples.html>

Horizon graphs

Time-Series Data: Figure 1d. Horizon graphs of U.S. unemployment rate, 2000–2010.

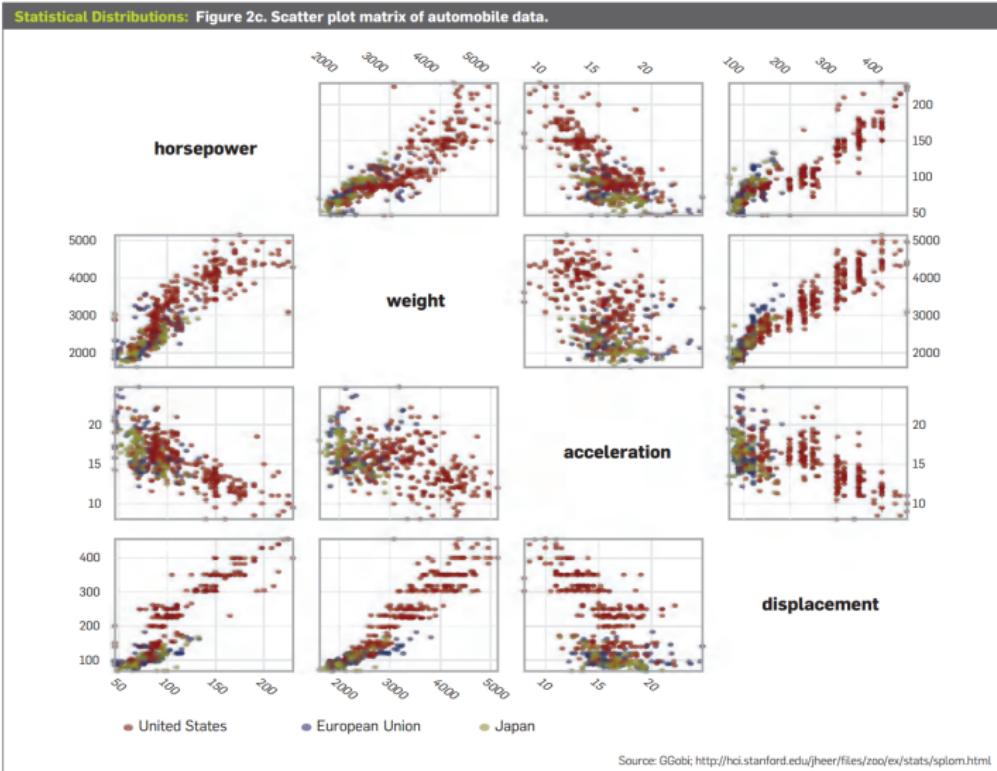


Source: U.S. Bureau of Labor Statistics; <http://hci.stanford.edu/jheer/files/zoo/ex/time/horizon.html>

Examples

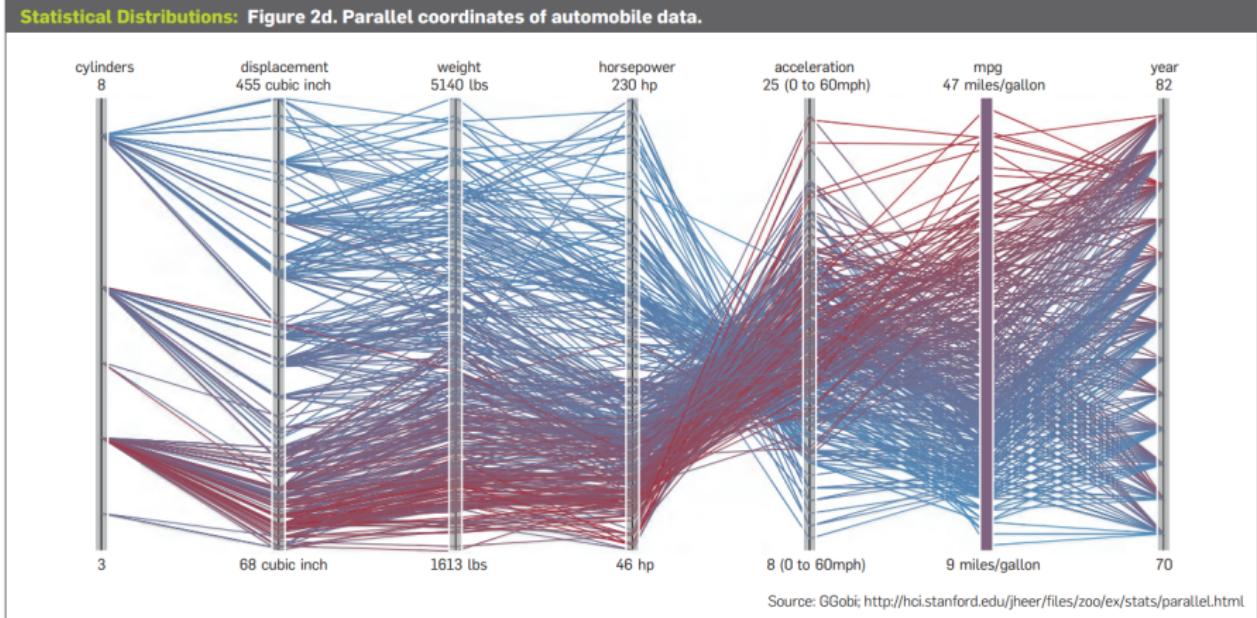
- Time-series data
- **Statistical data**
- Geographical data
- Hierarchical data
- Network data

Scatter plot matrix



Parallel coordinates

Statistical Distributions: Figure 2d. Parallel coordinates of automobile data.



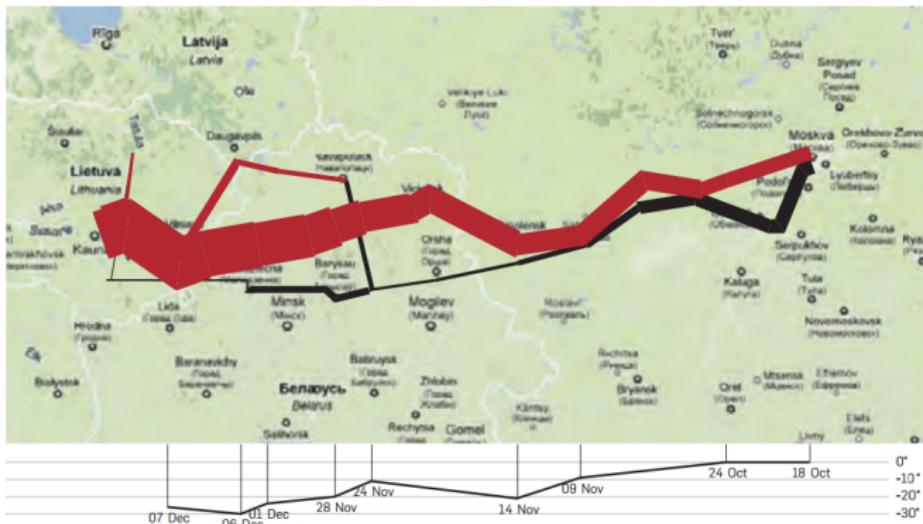
Source: GGobi; <http://hci.stanford.edu/jheer/files/zoo/ex/stats/parallel.html>

Examples

- Time-series data
- Statistical data
- **Geographical data**
- Hierarchical data
- Network data

Flow maps

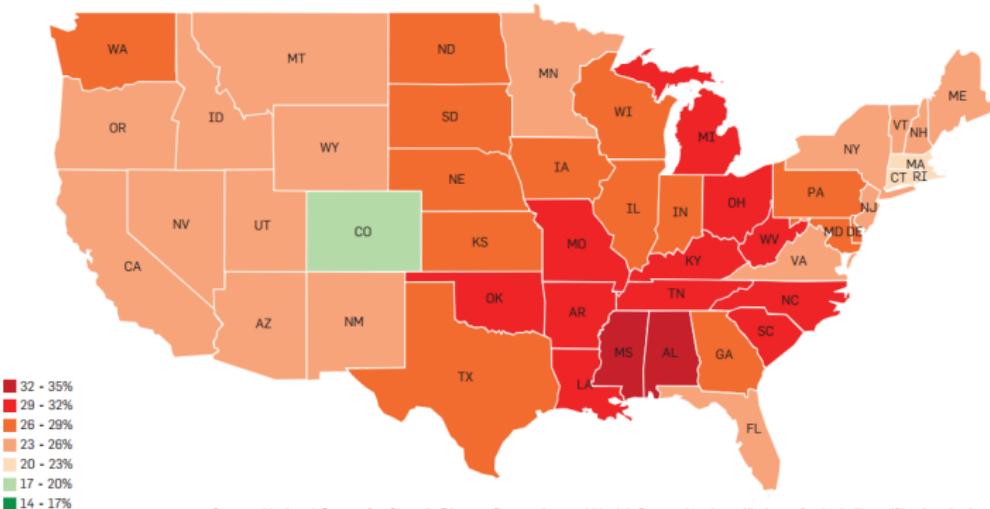
Maps: Figure 3a. Flow map of Napoleon's March on Moscow, based on the work of Charles Minard.



<http://hci.stanford.edu/jheer/files/zoo/ex/maps/napoleon.html>

Choropleth maps

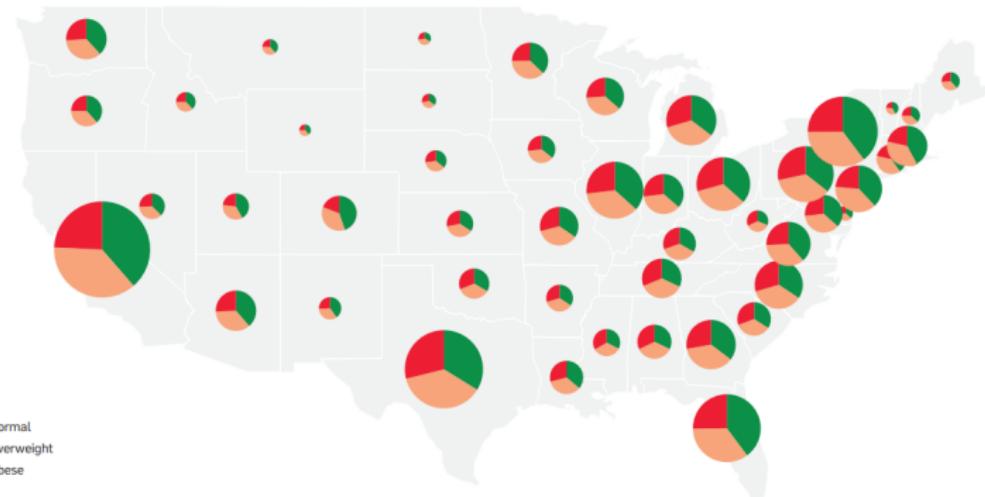
Maps: Figure 3b. Choropleth map of obesity in the U.S., 2008.



Source: National Center for Chronic Disease Prevention and Health Promotion; <http://hci.stanford.edu/jheer/files/zoo/ex/maps/choropleth.html>

Graduated symbol maps

Maps: Figure 3c. Graduated symbol map of obesity in the U.S., 2008.



Source: National Center for Chronic Disease Prevention and Health Promotion; <http://hci.stanford.edu/jheer/files/zoo/ex/maps/symbol.html>

Cartogram

Maps: Figure 3d. Dorling cartogram of obesity in the U.S., 2008.



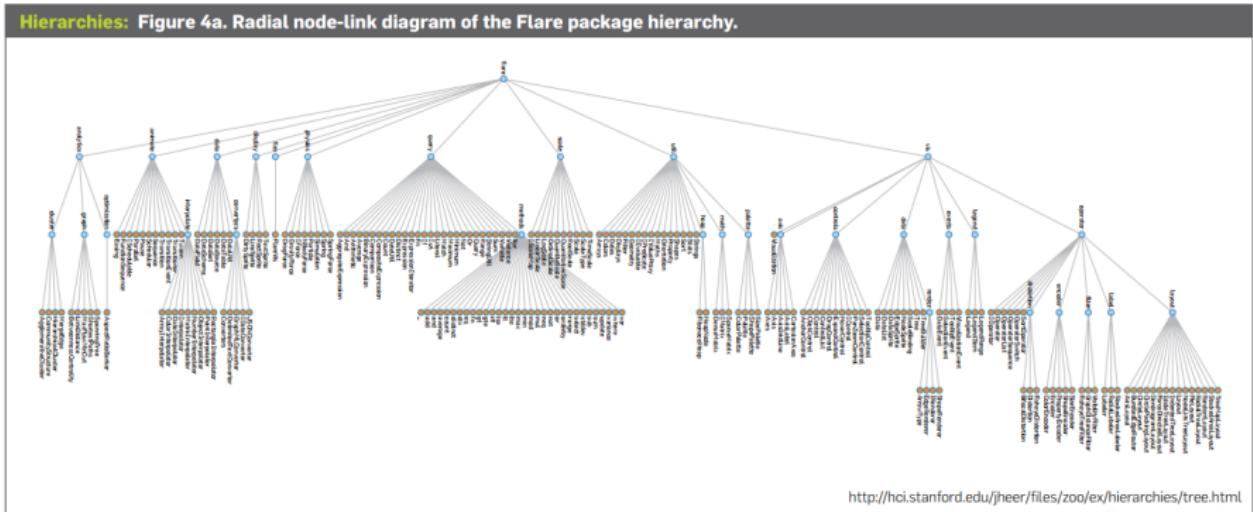
Source: National Center for Chronic Disease Prevention and Health Promotion: <http://hcip.stanford.edu/iheer/files/zoo/ex/maps/cartogram.html>

Examples

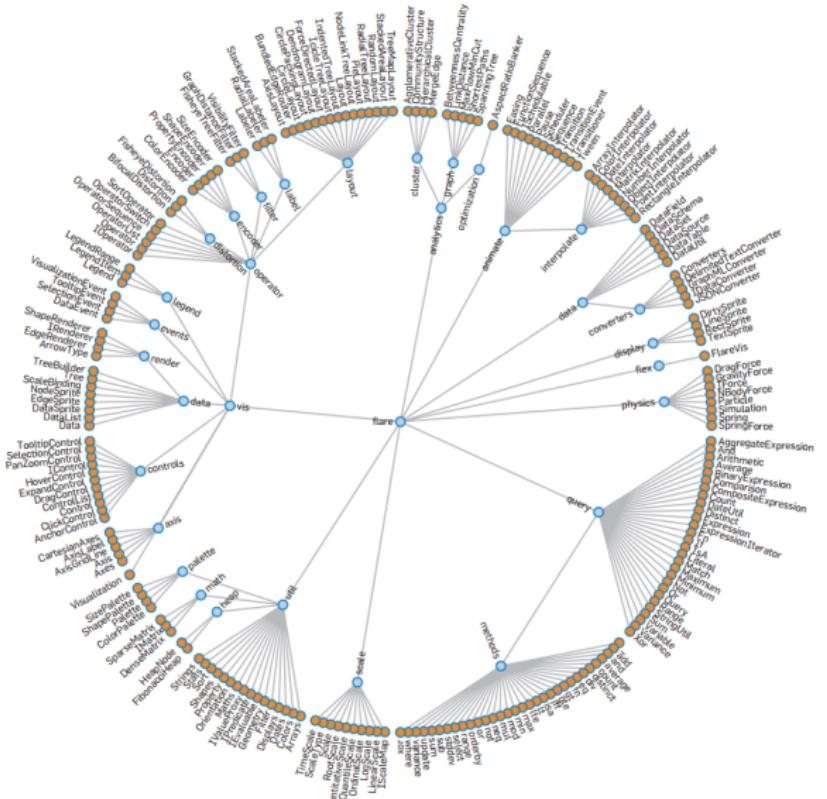
- Time-series data
- Statistical data
- Geographical data
- **Hierarchical data**
- Network data

Node-link diagram

Hierarchies: Figure 4a. Radial node-link diagram of the Flare package hierarchy.



Circular dendrogram

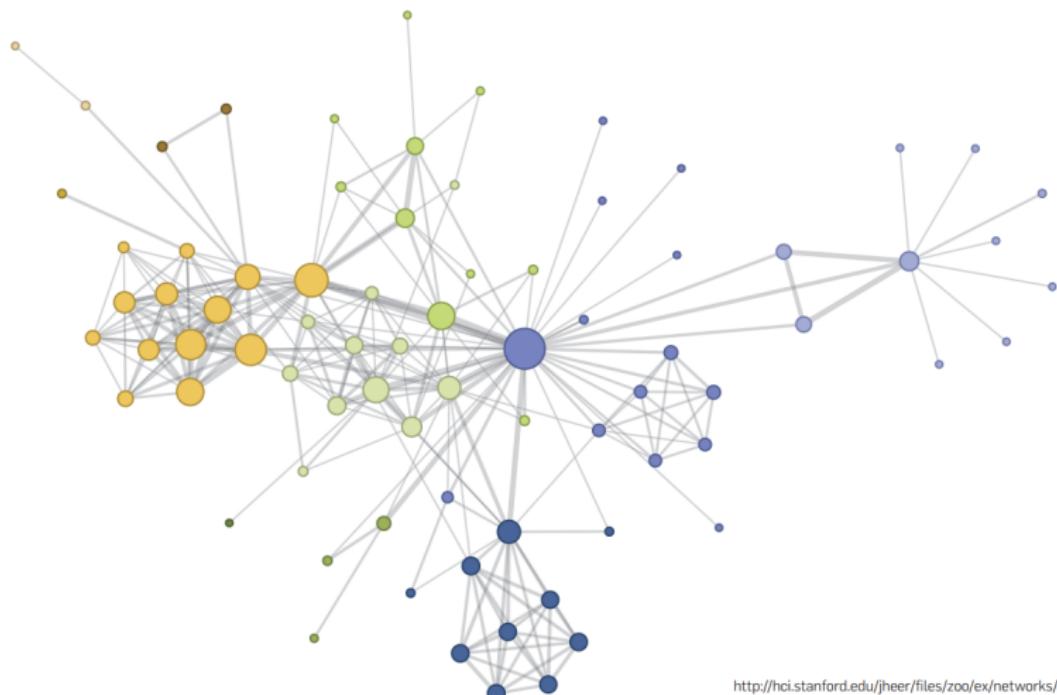


Examples

- Time-series data
- Statistical data
- Geographical data
- Hierarchical data
- **Network data**

Force-directed layout

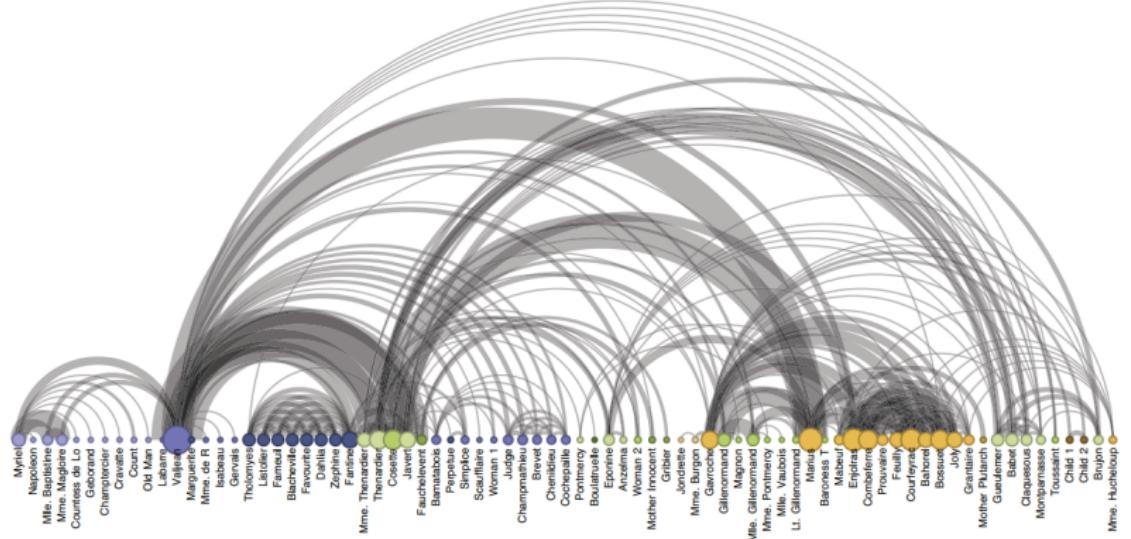
Networks: Figure 5a. Force-directed layout of *Les Misérables* character co-occurrences.



<http://hci.stanford.edu/jheer/files/zoo/ex/networks/force.html>

Arc diagram

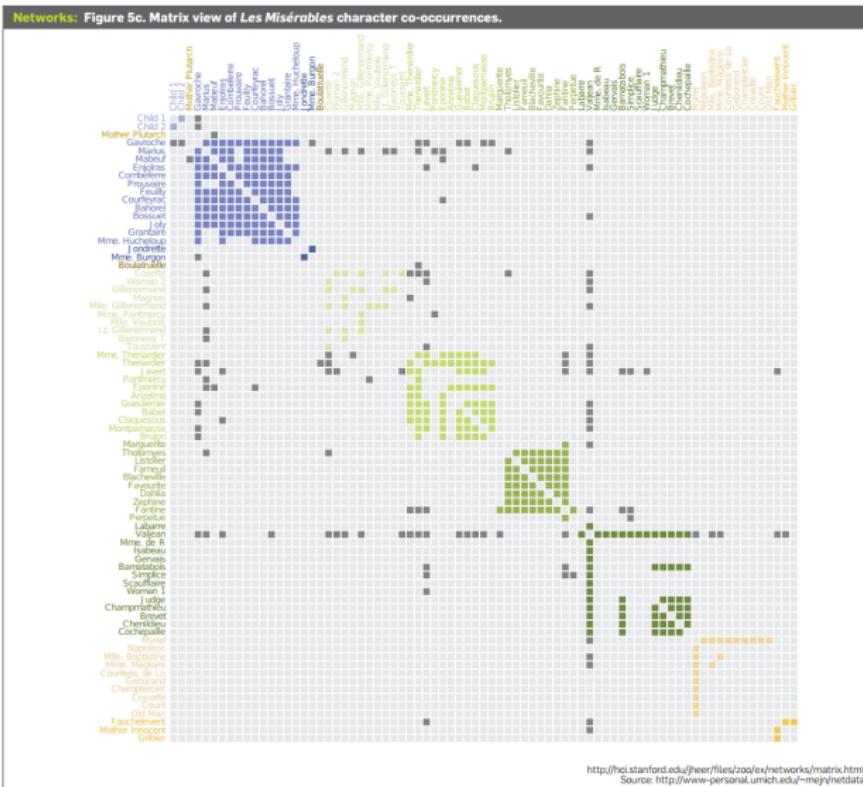
Networks: Figure 5b. Arc diagram of *Les Misérables* character co-occurrences.



<http://hci.stanford.edu/jheer/files/zoo/ex/networks/arc.html>

Annotated matrix

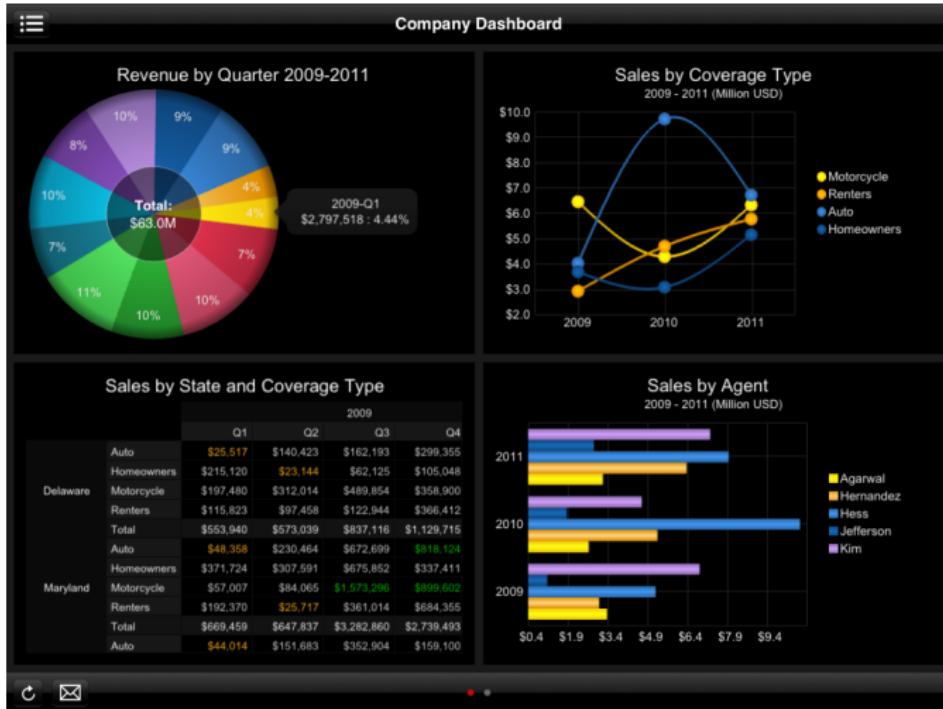
Networks: Figure 5c. Matrix view of *Les Misérables* character co-occurrences.



Dashboards

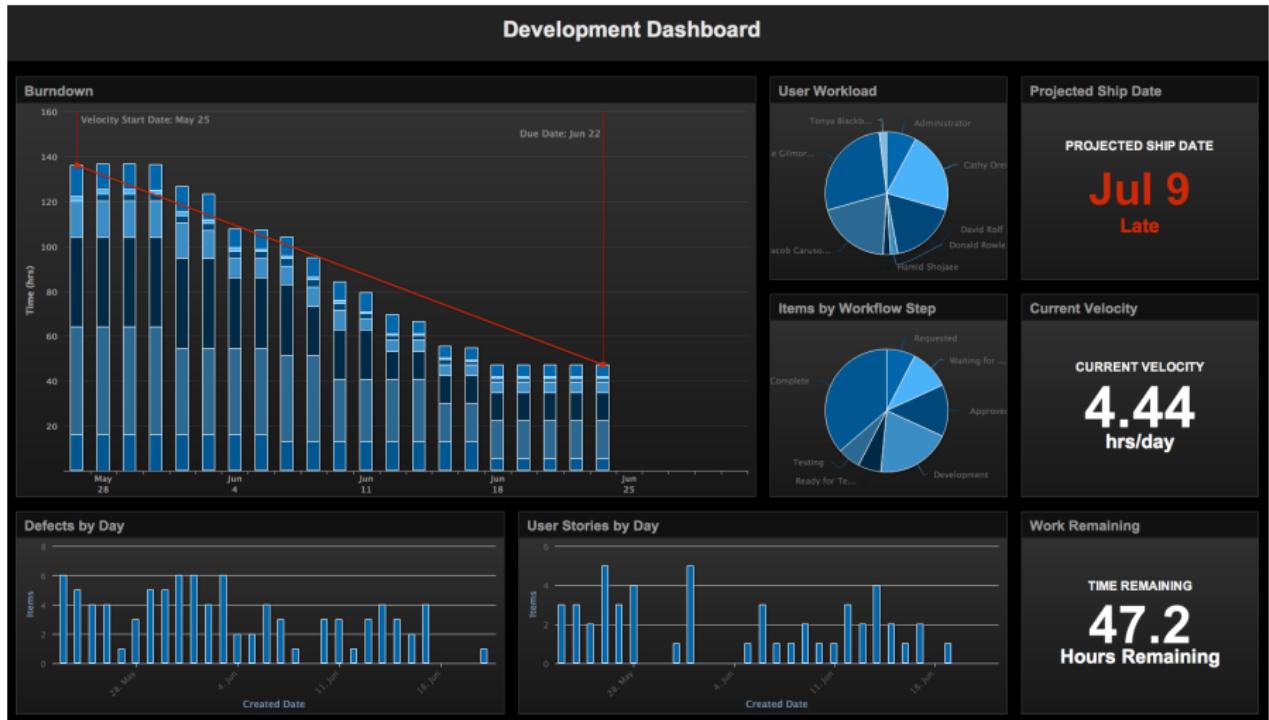
- Multiple **widgets** on one page
- A widget can contain:
 - OLAP slice
 - KPI metric
 - Data mining results
 - ...
- Codeless reporting
- BI in the blink of an eye!

Dashboard (1)


<http://blog.jinfonet.com/>

Dashboard (2)

Development Dashboard


<http://www.axosoft.com/>

Dashboard (3)

Cyfe Social Media Mar 11, 2014 - Apr 9, 2014

FACEBOOK OVERVIEW

450 +97% REACH 109 +36% VIEWS 60 ENGAGED 43 +71% CLICKS 30 +17% LIKES



TWITTER ANALYTICS

644 TWEETS 16 FOLLOWING 9,353 FOLLOWERS 992 LISTED 0 FAVORITES



FACEBOOK LIKES MAP



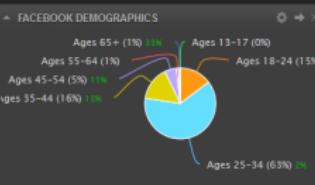
FACEBOOK TRAFFIC SOURCES

EXTERNAL REFERER	VIEWS
1. cyfe.com	37 +1%
2. google.com	7 +1%
3. google.com.br	4 +1%

LATEST TWEETS

- Element Three**
Day 2 of #PartnerDay got you feeling a bit tired? Revitalize yourself with some @HubSpot bingo! <http://it.co/gRvzSgPyop>
Wednesday, April 09, 2014 9:59:54 AM
- Heather Sutton**
I hate that @Android isn't supported yet. | 9 Ways to Use Twitter's New Photo Collages in Your Marketing <http://it.co/jcKXWGNtmd> via @hubspot
Wednesday, April 09, 2014 9:59:23 AM
- RealTimePersonalise**
How Dynamic Content Makes Your Marketing a Helluva Lot More Personal <http://it.co/Skw44bh2Eb> via @hubspot
Wednesday, April 09, 2014 9:59:12 AM
- Kyle Rumble**
9 Ways to Use Twitter's New Photo Collages in Your Marketing <http://it.co/d9rmWGq8ogh>
Wednesday, April 09, 2014 9:59:12 AM
- GoatCloud**
via @Hubspot: 9 Ways to Use Twitter's New Photo Collages

FACEBOOK DEMOGRAPHICS



AGE GROUP	PCT
Ages 13-17 (0%)	0%
Ages 18-24 (15%)	15%
Ages 25-34 (63%)	63%
Ages 35-44 (16%)	16%
Ages 45-54 (5%)	5%
Ages 55-64 (1%)	1%
Ages 65+ (1%)	1%

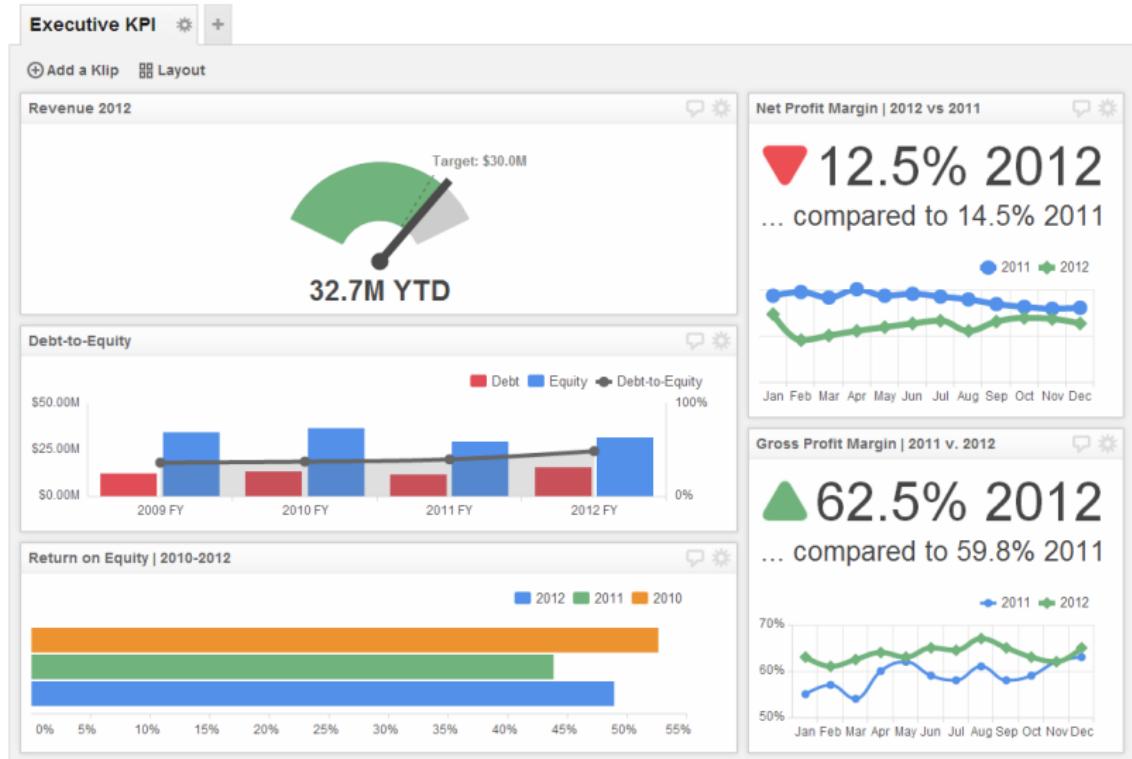
LINKEDIN ANALYTICS

61,044 FOLLOWERS



<http://www.cyfe.com/>

Dashboard (4)



Dashboard (5)

The Music Industry Dashboard


<http://insideanalysis.com/>

Assignment 1

- Gaming industry context
- Sales log spanning 4 years of sales
- Apply and compare BI techniques
- Inspect, visualize, aggregate, segment, score ...
- Deliverables:
 - 1 Web-based BI Dashboard
 - 2 Short assignment report in L^AT_EX

Assignment 1 — Hints

- Model: MySQL database containing the data
- View: HTML page using Javascript that reads JSON
- Controller: PHP outputs relevant data in JSON

Lab session February 17

- Make progress with Assignment 1
- Read paper by Kooti et al.
- Setup a framework for your dashboard
- Load some data into your framework
- Investigate visualization options

Credits

Lecture based on (slides of the (previous edition of the)) course book:
W. van der Aalst, *Process Mining: Data Science in Action*, 2nd edition,
Springer, 2016.

