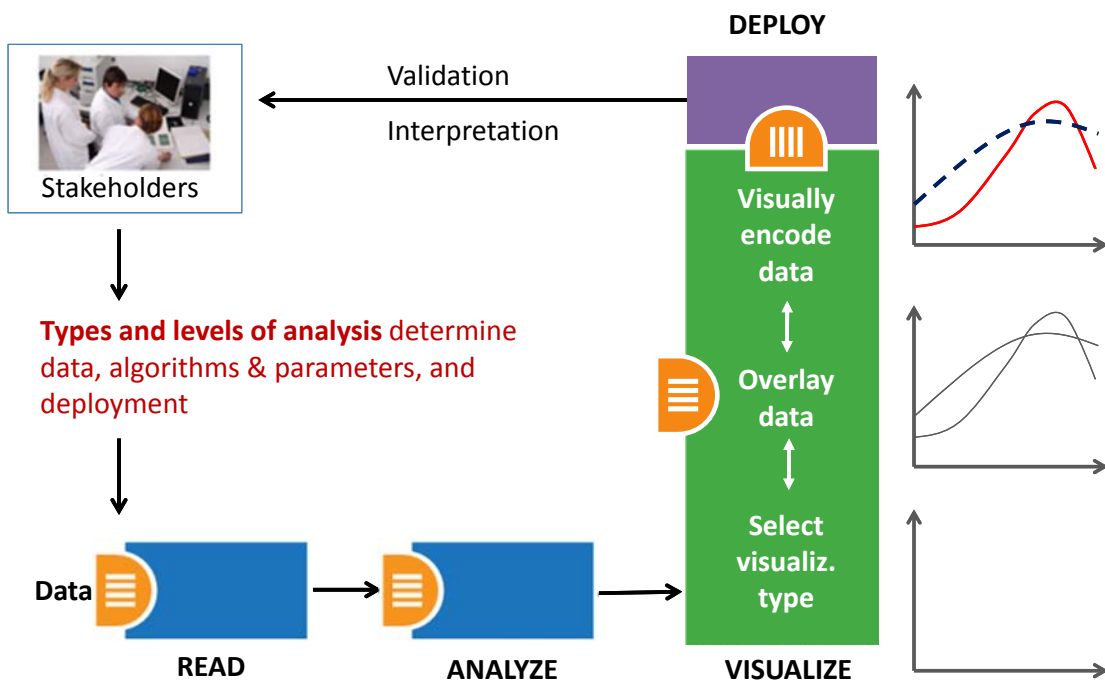


# Needs-Driven Workflow Design



25

## Information Visualization MOOC

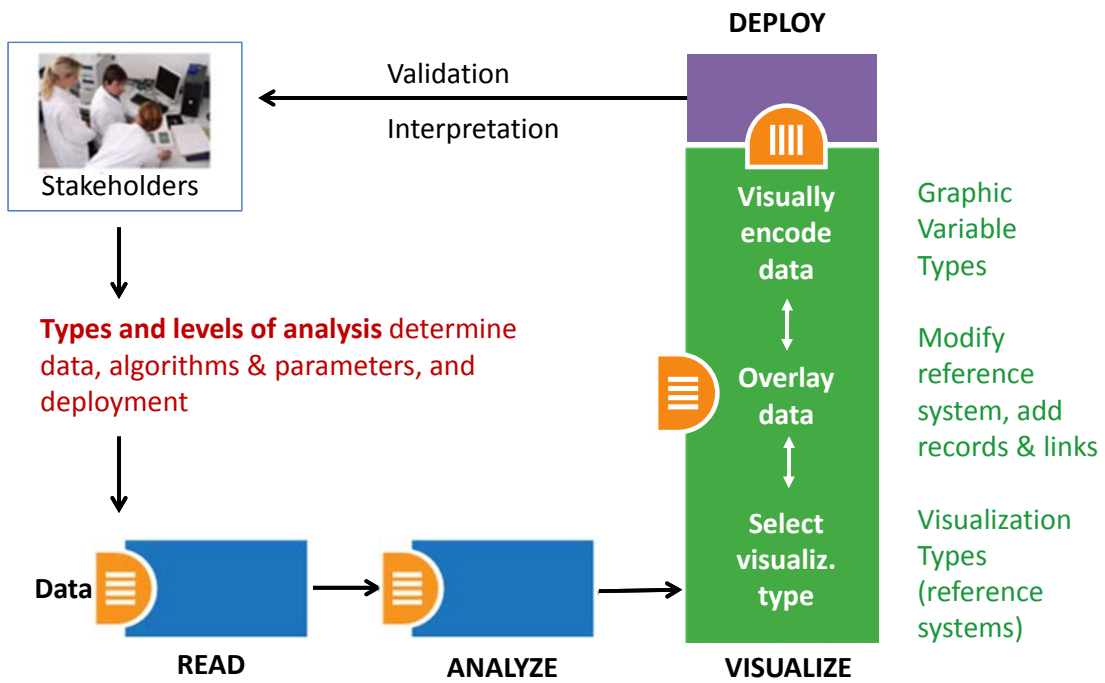
### Unit 2 – “When”: Temporal Data

#### Workflow Design

#### Relevant Research Disciplines:

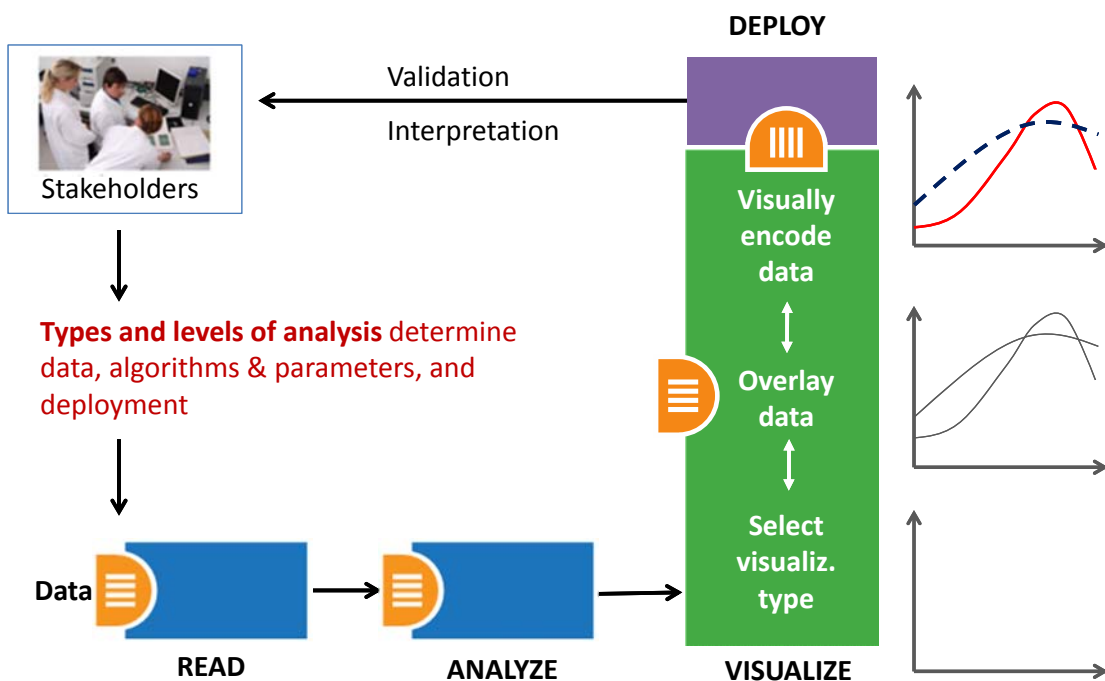
Mathematics, Statistics, Information Visualization

# Needs-Driven Workflow Design



27

# Needs-Driven Workflow Design



28

## Read Data

Data Repositories:

- Gapminder data, <http://www.gapminder.org/data/>
- UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml> (Time Series)
- Scholarly Database, <http://sdb.cns.iu.edu>
- IBM ManyEyes Datasets, <http://www-958.ibm.com/software/data/cognos/manyeyes/datasets> (350,976 )
- Eurostat Data Market, <http://datamarket.com>

Data Formats:

- TXT
- XLS, CSV
- Databases

29

## Data Formats

Time-series events are commonly represented by lists—e.g., emails

<b>Subject header</b>	<b>Date</b>	<b>Time</b>
Meeting to get key	1/1/2012	2:01 PM
CNS Talk–Meet Speaker?	1/1/2012	2:05 PM
Review paper draft	1/1/2012	3:01 PM
Happy New Year!	1/1/2012	3:01 PM
Action List Reminder	1/1/2012	3:04 PM

or Amazon book rankings over time:

<b>Month</b>	<b>Day</b>	<b>Year</b>	<b>Time</b>	<b>Rank</b>
10	23	2011	22:00	63,290
10	23	2011	23:00	76,128
10	24	2011	0:00	76,128
10	24	2011	1:00	95,578
10	24	2011	2:00	104,756

Each event is represented by a data record (row) that has temporal attribute(s) and possibly other attributes.

30

# Data Preprocessing

Filtering—e.g., time slicing (see next slide)

Test for and remove large spikes in the data, but report this.

Normalization

- Deduplication
- Unit conversion
- Adjust (dollars) for inflation
- Adjust for time zones

Integration/interlinkage of different data sources

Classification/aggregation

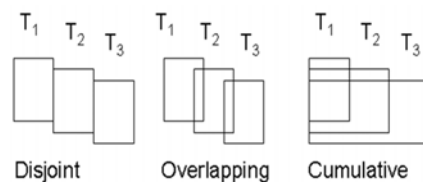
31

## Data Preprocessing – Time Slicing

**Resolution:** Milliseconds, seconds, minutes, hours, days, weeks, fortnights (fourteen days /two weeks), months, quarters, years, decades, and centuries.

**Type:**

- Disjoint: Every row in the original table is in exactly one time slice.
- Overlapping: Selected rows are in multiple time slices.
- Cumulative: Every row in a time slice is in all later time slices.



**Alignment with calendar:** If first event is June 7th, 2006, and yearly slices are chosen, then the first slice will be from

- June 7th, 2006, to June 6th, 2007 (if not aligned)
- January 1st, 2006, to December 31st, 2006 (if aligned)

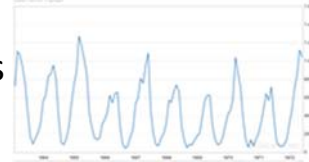
32

# Analyze Data

**Trends**—e.g., in baby names

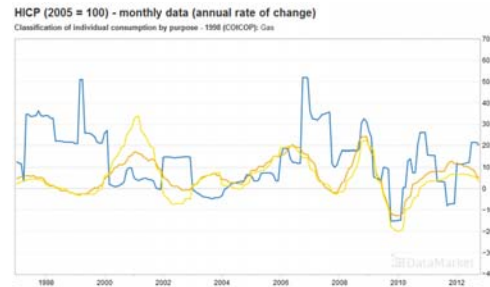


**Patterns**—e.g., seasonality in chickenpox cases

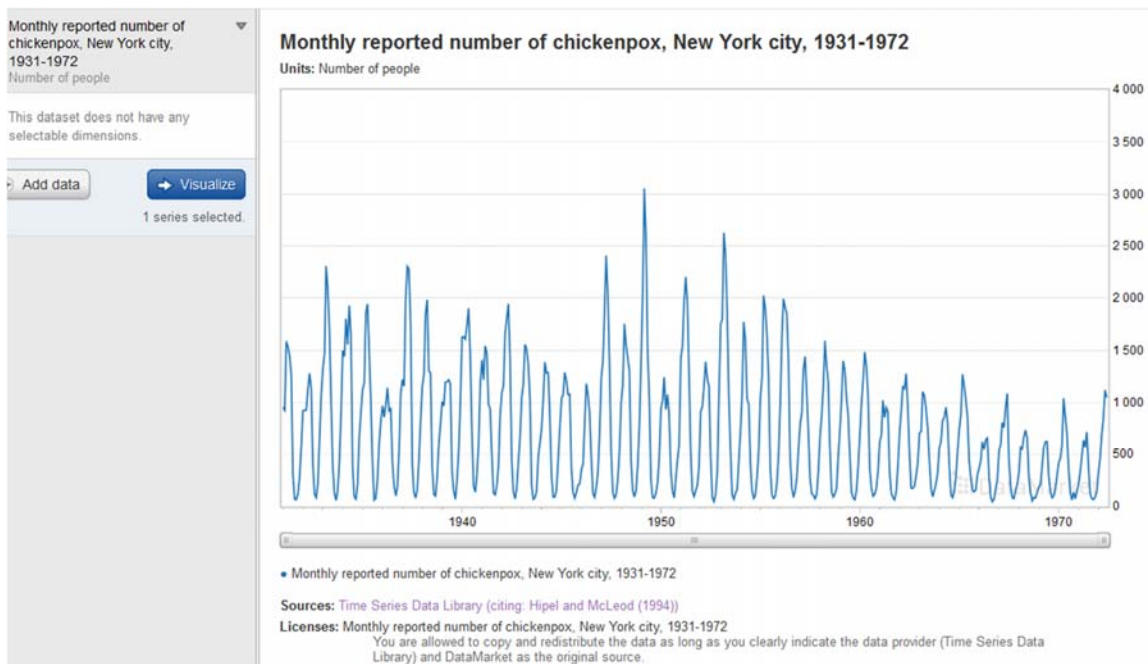


**Correlations**—e.g., in gas prices or movie ratings

**Bursts**—see final [Unit 3 slide set](#).



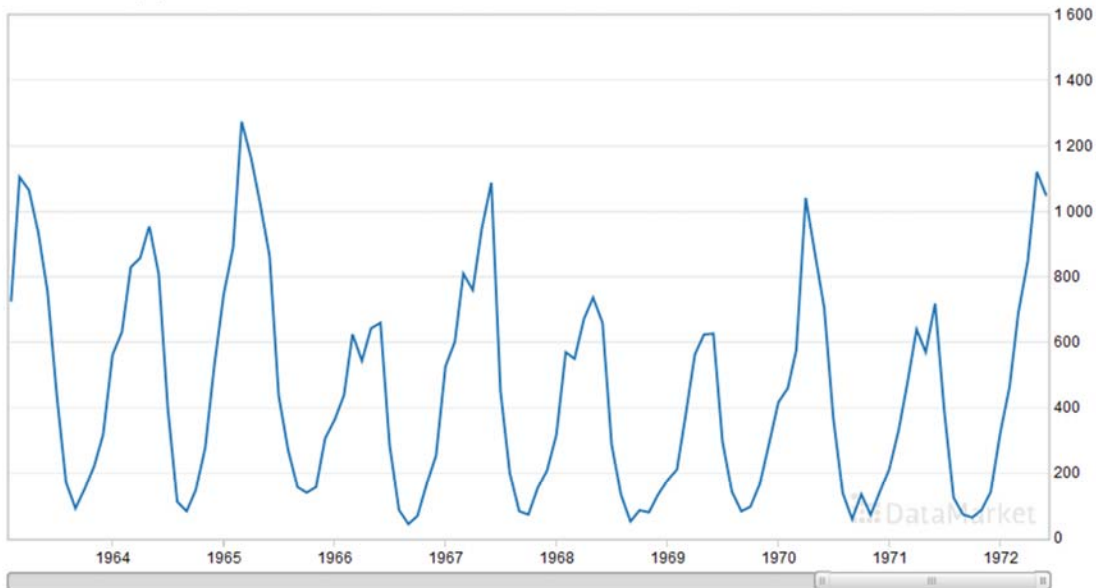
DataMarket: Monthly reported number of chickenpox, New York City, 1931-1972  
<http://datamarket.com/data/set/22v7/monthly-reported-number-of-chickenpox-new-york-city-1931-1972#display=line&ds=22v7=provider:tsdl>



DataMarket: Monthly reported number of chickenpox, New York City, **Feb. 1963** -1972

<http://datamarket.com/data/set/22v7/monthly-reported-number-of-chickenpox-new-york-city-1931-1972#display=line&ds=22v7=provider:tsdl>

Units: Number of people



• Monthly reported number of chickenpox, New York city, 1931-1972

Sources: Time Series Data Library (citing: Hipel and McLeod (1994))

Licenses: Monthly reported number of chickenpox, New York city, 1931-1972

You are allowed to copy and redistribute the data as long as you clearly indicate the data provider (Time Series Data Library) and DataMarket as the original source.

35

DataMarket: Gas Prices from Jan. 1997 to Oct. 2012

<http://datamarket.com/en/data/set/1a6e/#!ds=1a6e!qvc=4b:qvd=m.e&display=line>

▼ HICP (2005 = 100) - monthly data (annual rate of change)

Search in dimensions

Classification of individual consumption by purpose - 1998 (COICOP) Clear

Geopolitical entity (declaring) Clear

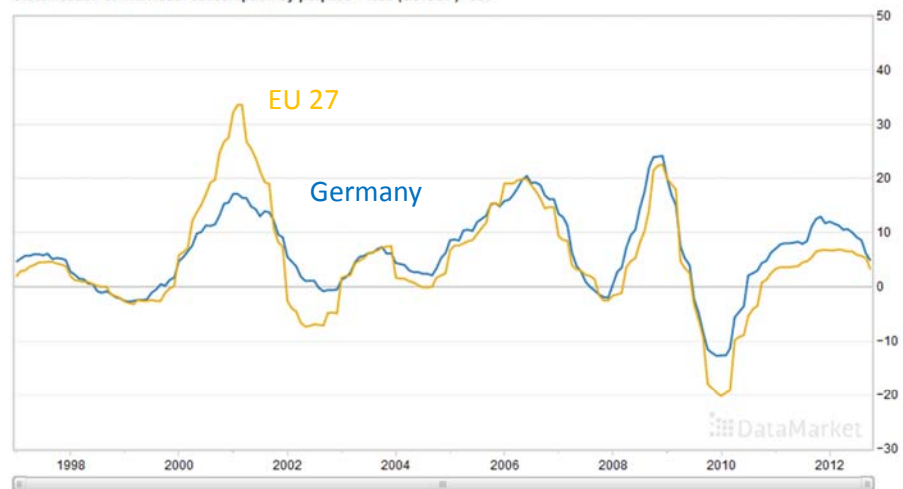
- EA15-2008, EA16-2010, EA17)
- European Economic Area (EEA18-2004, EEA28-2006, EEA30)
- European Union (27 countries)
- European Union (EU6-1972, EU9-1980, EU10-1985, EU12-1994, EU15-2004, EU25-2006, EU27)
- Finland
- France
- Germany (including former GDR from 1991)
- Greece
- Hungary
- Iceland

+ Add data Visualize

2 series selected.

HICP (2005 = 100) - monthly data (annual rate of change)

Classification of individual consumption by purpose - 1998 (COICOP): Gas



• European Union (27 countries) • Germany (including former GDR from 1991)

Sources: Eurostat

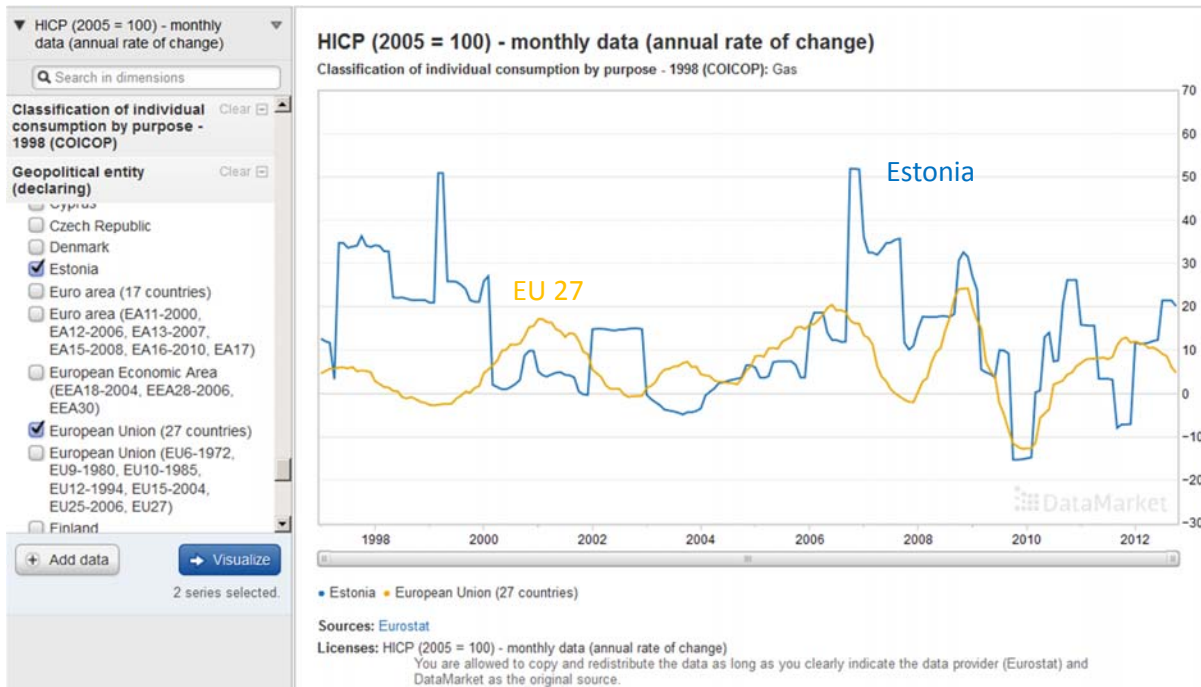
Licenses: HICP (2005 = 100) - monthly data (annual rate of change)

You are allowed to copy and redistribute the data as long as you clearly indicate the data provider (Eurostat) and DataMarket as the original source.

36

## DataMarket: Gas Prices from Jan. 1997 to Oct. 2012

<http://datamarket.com/en/data/set/1a6e/#!ds=1a6e!qvc=4b:qvd=m.e&display=line>



37

## DataMarket: Gas Prices from Jan. 1997 to Oct. 2012

<http://datamarket.com/en/data/set/1a6e/#!ds=1a6e!qvc=4b:qvd=m.e&display=line>

Is there a correlation? Germany and Estonia are part of EU27.

Geopolitical entity	Estonia	European	Germany (including former GDR from 1991)
Month			
1997-01	12.6	4.6	1.9
1997-02	12	5.2	2.9
1997-03	11.7	5.7	3
1997-04	3.3	5.7	3.7
1997-05	34.6	6	4
1997-06	34.6	6	4.5
1997-07	33.5	5.8	4.5
1997-08	33.8	6.1	4.6

Also relevant for calculating if:

- Twitter activity correlated to stock market behavior.
- Paper download counts correlated to citation counts.

38

## Correlation

Use MS Excel or other statistics program to calculate

$$\text{Correl}(X, Y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}}, \text{ where } \bar{x} \text{ and } \bar{y} \text{ are the means.}$$

E15		fx =CORREL(C14:C203,D14:D203)					
	A	B	C	D	E	F	G
11							
12	Geopolitical en	Estonia	European	Germany (including former GDR from 1991)			
13	Month						
14	1997-01	12.6	4.6	1.9			
15	1997-02	12	5.2	2.9	0.905826		
16	1997-03	11.7	5.7	3	0.049288		
17	1997-04	3.3	5.7	3.7			

Note: Sequence of data pairs does not impact correlation result.

39

## Correlation

Values run from -1 (no correlation) to 1 (identical).

D1		fx =CORREL(A1:A7,B1:B7)				
	A	B	C	D	E	F
1	2	2	1	1		
2	1	1	2	-1		
3	2	2	1			
4	1	1	2			
5	2	2	1			
6	1	1	2			
7	2	2	1			

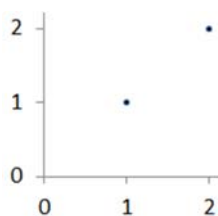
**Positive association:**

Upward trend (positive slope).

**Negative association:**

Downward trend (negative slope).

**Non-associated:** No trends.



40

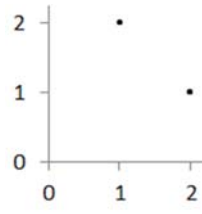
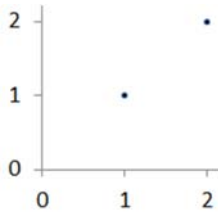


# Correlation

Values run from -1 (no correlation) to 1 (identical).

D1		fx: =CORREL(A1:A7,B1:B7)				
	A	B	C	D	E	F
1	2	2	1	1		
2	1	1	2	-1		
3	2	2	1			
4	1	1	2			
5	2	2	1			
6	1	1	2			
7	2	2	1			

D2		fx: =CORREL(B1:B7,C1:C7)				
	A	B	C	D	E	F
1	2	2	1	1		
2	1	1	2	-1		
3	2	2	1			
4	1	1	2			
5	2	2	1			
6	1	1	2			
7	2	2	1			



41

DataMarket: Gas Prices from Jan. 1997 to Oct. 2012

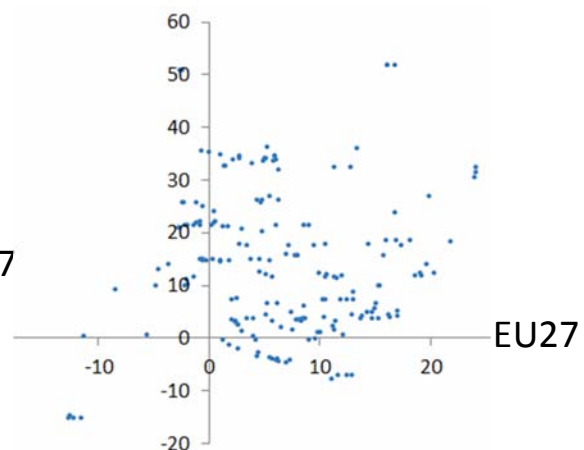
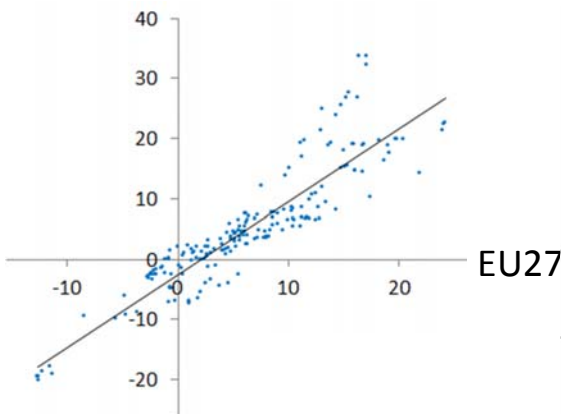
<http://datamarket.com/en/data/set/1a6e/#!ds=1a6e!qvc=4b:gvd=m.e&display=line>

Is there a correlation? Germany and Estonia are part of EU27.

Visualization via scatterplot:

Germany (0.9058)

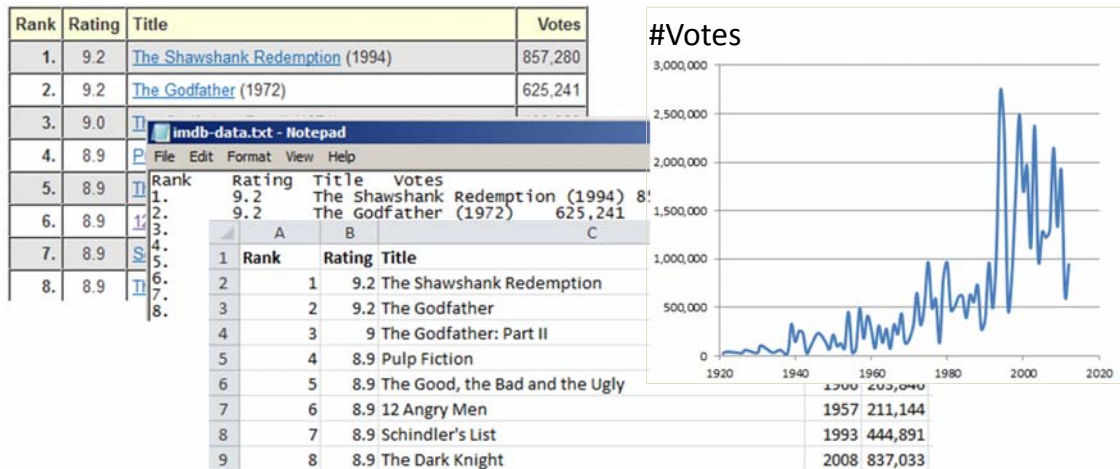
Estonia (0.04929)



42

## Top-250 Movies from IMDb

Copy from web page at <http://www.imdb.com/chart/top> (on Nov 15, 2012), save as text file, open in Excel or other table editing program:

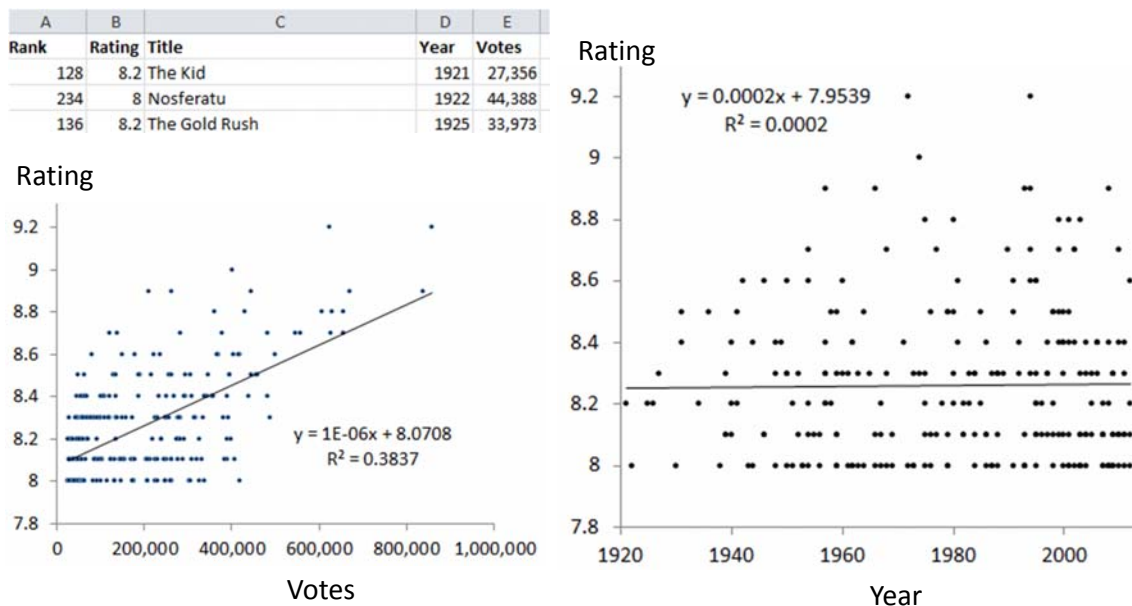


Is there a correlation between Rating and #Votes?

Any trends over time? Recent movies more popular?

43

## Top-250 Movies from IMDb



Is there a correlation between Rating and #Votes?

Any trends over time? Recent movies more popular?

44

# Visualize Data

**Spreadsheet programs** such as MS Excel or the free Apache Open Office (<http://www.openoffice.org>) support the easy generation of diverse charts and graphs.

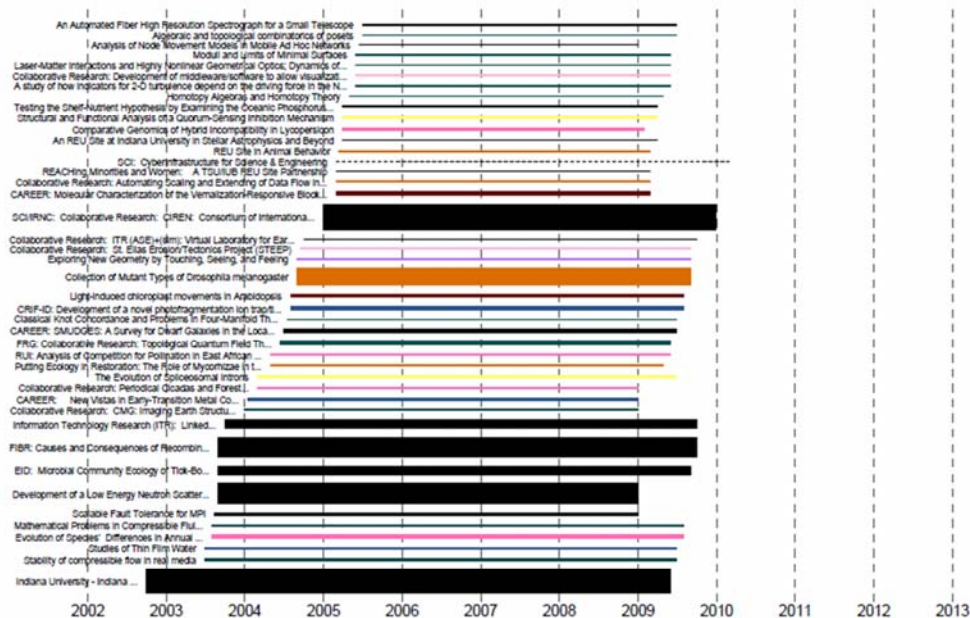
**Temporal bar graphs**, see next slide and Unit 2: Hands-on.

**Time slice data** to generate, e.g.,

- Evolving geomaps, see Unit 3.
- Evolving networks, see Unit 7: Hands-on.

## Temporal Visualization

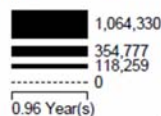
Generated from NSF csv file: Indiana.nsf  
July 18, 2012 | 8:46 AM EDT



### Legend

Area size: Awarded Amount to Date  
Minimum = 0  
Maximum = 6,402,330  
Text label: Title  
Color: NSF Organization  
See end of PDF for color legend.

### Area



### How To Read This Map

This temporal bar graph visualization represents each record as a horizontal bar with a specific start and end date and a text label on its left side. The area of each bar encodes a numerical attribute value, e.g., total amount of funding. Bars may be colored to present categorical attribute values of records.

## Relevant Tools

- TimeSearcher from HCIL supports the visual exploration of time-series data <http://www.cs.umd.edu/hcil/timesearcher/>
- Tableau, <http://www.tableausoftware.com>

Please post your favorite to Twitter or Flickr using tags “#ivmoooc” and “#timetools.”

