# Visualization and Big Data in Official Statistics

## Martijn Tennekes

**In cooperation with Piet Daas, Marco Puts, May Offermans, Alex Priem, Edwin de Jonge**
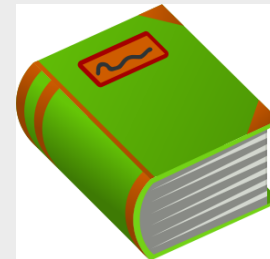
**Statistics Netherlands**

# From a Official Statistics point of view

Three types of data:

1. Survey data = data collected by SN with questionnaires



2. Admin data = administrative (register) data collected by third parties such as the Tax Office



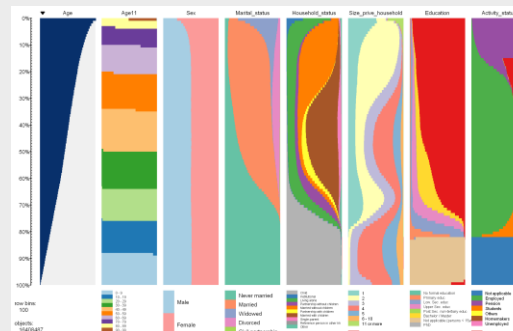3. Big data = machine generated data of events

# Big Data case studies

Big data = machine generated data of events

| Source | Statistics |
|---|---|
| Social media | Sentiment (as indicator for business cycle) |
| Mobile phone metadata | Daytime population, tourism statistics |
| Road sensors | Traffic index statistics |

At the end of this talk:
Visualization methods for Big Data

3

# Big data approach



General Data Science workflow

# Case study 1: Social media

- 3 billion messages as of 2009 gathered from Facebook, Twitter, LinkedIn, Google+ by a Dutch intermediate company Coosto.

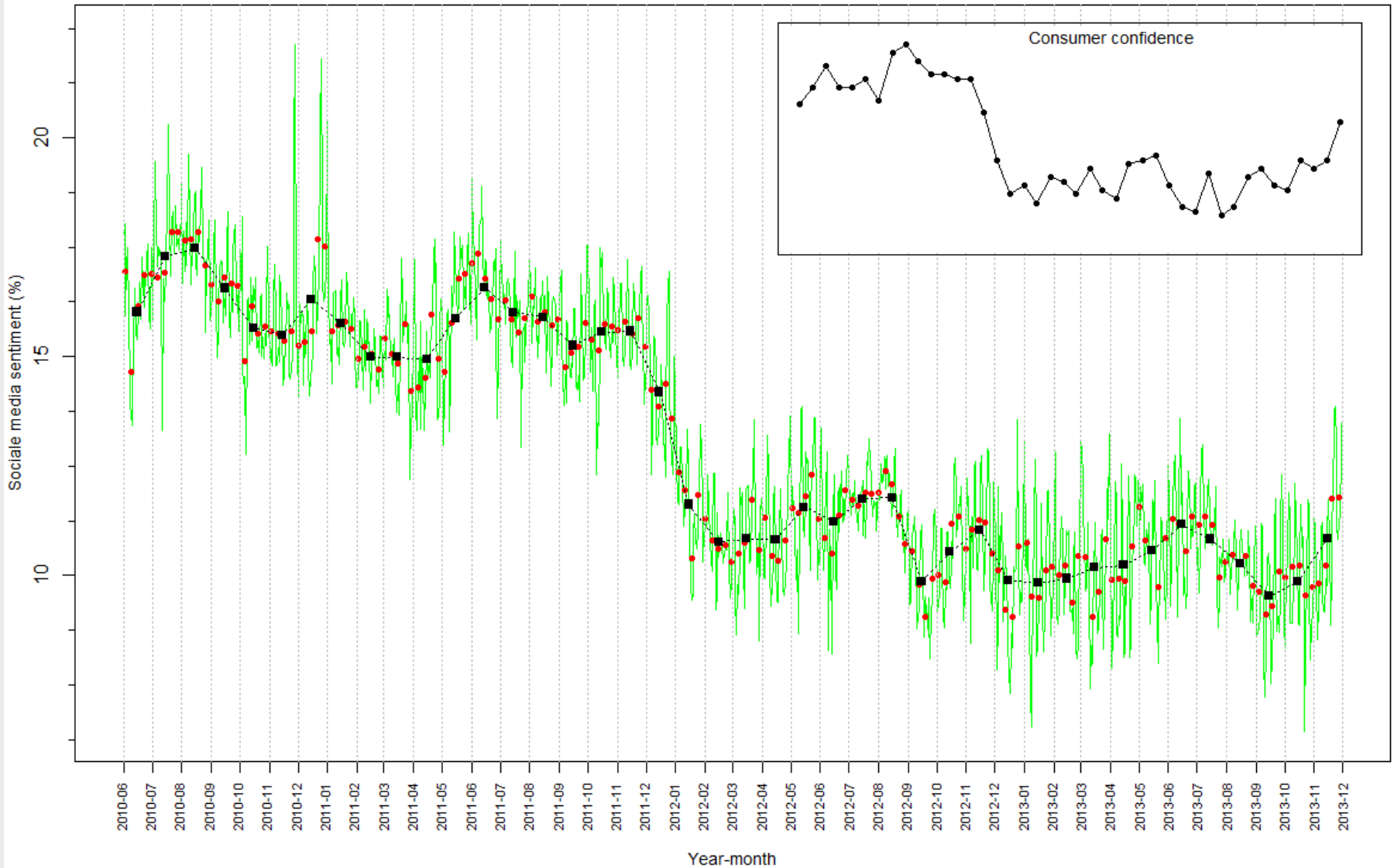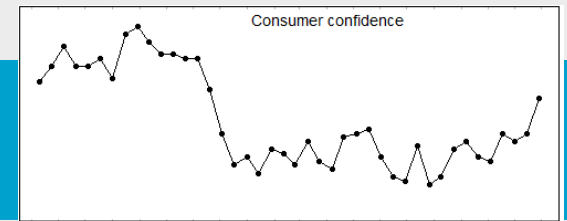- Sentiment per message determined by classifying words as negative or positive.

- Could be used as indicator for the business cycle. Could it be fit to the **consumer confidence**, the leading business cycle indicator?
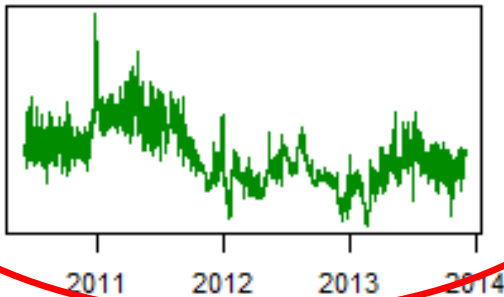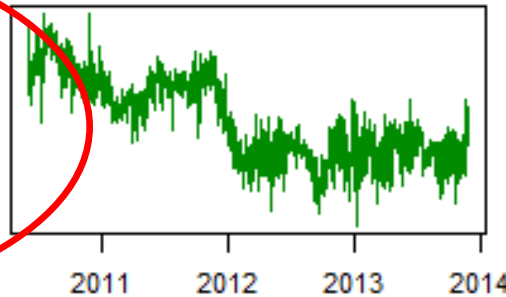
# Sentiment in social media
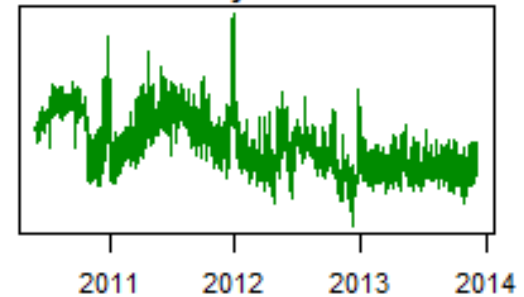
# Platform specific sentiment

# Platform specific results

**Table 1**. Social media messages properties for various platforms and their correlation with consumer confidence

| Social media platform | Number of social media messages[1] | Number of messages as percentage of total (%) | Correlation coefficient of monthly sentiment index and consumer confidence ( $r$ )[2] |
|---|---|---|---|
| All platforms combined | 3,153,002,327 | 100 | 0.75 |
| Facebook | 334,854,088 | 10.6 | 0.81* |
| Twitter | 2,526,481,479 | 80.1 | 0.68 |
| Hyves | 45,182,025 | 1.4 | 0.50 |
| News sites | 56,027,686 | 1.8 | 0.37 |
| Blogs | 48,600,987 | 1.5 | 0.25 |
| Google+ | 644,039 | 0.02 | -0.04 |
| Linkedin | 565,811 | 0.02 | -0.23 |
| Youtube | 5,661,274 | 0.2 | -0.37 |
| Forums | 134,98,938 | 4.3 | -0.45 |

[1]period covered June 2010 untill November 2013

[2]confirmed by visual inspecting scatterplots and additional checks (see text)

*cointegrated

8

# Case study 2: mobile phone metadata

- Pilot study with Vodafone, a provider with market share of 1/3 in the Netherlands.

- Aggregated data is queried by intermediate company Mezuro and delivered to SN. Privacy is guaranteed!

- Applications: daytime population, tourism statistics, economic activity, mobility studies, etcetera.

# Mobile phone population



MPRD (Municipal Personal Records Database) = Dutch population

# Subpopulations model

# Mobile phone metadata

Event Datail Records (EDR) contain metadata on mobile phone events (i.e. call, SMS or data transfer).

Aggregated table: number of unique devices X time period X current region X residential region.

# Weighting method

Example: suppose there are only 3 regions in the Netherlands: Amsterdam, Boskoop and Castricum

| | | Residence | | | |
|---|---|---|---|---|---|
| | | Amsterdam | Boskoop | Castricum | |
| Current region at time $t$ | Amsterdam | 199,000 | 1,000 | 4,000 | |
| | Boskoop | 500 | 3,500 | 0 | |
| | Castricum | 500 | 500 | 16,000 | |
| | | | | | |

# Weighting method (2)

Example: suppose there are only 3 regions in the Netherlands: Amsterdam, Boskoop and Castricum

| | | Residence | | | |
|---|---|---|---|---|---|
| | | Amsterdam | Boskoop | Castricum | |
| Current region at time $t$ | Amsterdam | 199,000 | 1,000 | 4,000 | |
| | Boskoop | 500 | 3,500 | 0 | |
| | Castricum | 500 | 500 | 16,000 | |
| | **MPRD total** | **800,000** | **15,000** | **30,000** | |

# Weighting method (3)

Example: suppose there are only 3 regions in the Netherlands: Amsterdam, Boskoop and Castricum

|  |  | Residence | | | |
| --- | --- | --- | --- | --- | --- |
|  |  | Amsterdam | Boskoop | Castricum |  |
| Current region at time *t* | Amsterdam | 796,000 | 3,000 | 6,000 |  |
|  | Boskoop | 2000 | 10,500 | 0 |  |
|  | Castricum | 2000 | 1,500 | 24,000 |  |
|  | **MPRD total** | **800,000** | **15,000** | **30,000** |  |

# Weighting method (4)

Example: suppose there are only 3 regions in the Netherlands: Amsterdam, Boskoop and Castricum

| | Residence | | | |
|---|---|---|---|---|
| | | Amsterdam | Boskoop | Castricum | DTP total |
| Current region at time $t$ | Amsterdam | 796,000 | 3,000 | 6,000 | 805,000 |
| | Boskoop | 2000 | 10,500 | 0 | 12,500 |
| | Castricum | 2000 | 1,500 | 24,000 | 27,500 |
| | **MPRD total** | **800,000** | **15,000** | **30,000** | |

# Daytime population results

# Day time population (relative)



2013-05-07 00:00:00

- Very sparsely populated
- Sparsely populated
- Normally populated
- Densely populated
- Very densely populated

# Day time population (relative)



**Legend:**
- Very sparsely populated
- Sparsely populated
- Normally populated
- Densely populated
- Very densely populated

# Day time population (relative)

# Day time population (relative)

## City of Eindhoven and surrounding towns

# Day time population – Region profile

## K-means clustering

**Work** = daytime vs. night-time during working weeks

**Weekend** = weekends activity

**Holiday** = May holiday activity



Legend:
- ■ City Centre
- ■ Working region (busy)
- ■ Working region (normal)
- ■ No classification
- ■ Commuting region
- ■ Recreational region

# Case study 3: Road sensors



**Road sensors data**

- Each minute (24/7) the number of passing vehicles is counted in around 20.000 'loops' in the Netherlands (100 million records a day)



- Nice data source for transport and traffic statistics (and more)

23

# Road sensors on main roads



A2

A28

A27

A12

A close look at the highways around Utrecht

# Road sensors on main roads (2)



Traffic loops everywhere…

# Road sensors on main roads (3)



Highways simplified for analysis

Dutch highways by COROP region

# Raw data: Total number of vehicles a day

# Correct for missing data: macro level

**Sliding window of 5 min. Impute missing data.**



Before — Total = ~ 295 million detected vehicles

After — Total = ~ 330 million (+ 12%) detected vehicles

# Data by type of vehicle

Long vehicles (> 12.2 meter)



Small vehicles (<= 5.6 meter)

Medium vehicles (> 5.6 & <= 12.2 meter)

# Selectivity of big data

- Big Data sources may be selective when
  - Only part of the population contributes to the data set (e.g. mobile phone owners)
  - The measurement mechanism is selective (e.g. traffic loops placement on Dutch highways is not random)

- Many Big Data sources contain events
  - How to associate events with units?
  - Number of events per unit may vary.

- Correcting for selectivity
  - Background characteristics – or *features* – are needed (linking with registers; profiling)
  - Use predictive modeling / machine learning to produce population estimates

# Visualization of Big Data

- Large **v**olume:
  - Data binning or aggregation
- High **v**elocity:
  - Animations
  - Dashboard / small multiples
- Large **v**ariety:
  - Interactive interface
  - Advanced visualization methods

# Tableplot: Dutch (Virtual) Census



row bins:
100

objects:
16408487

**Age** | **Age11** | **Sex** | **Marital_status** | **Household_status** | **Size_prive_household** | **Education** | **Activity_status**

Age11 legend:
- 0 - 9
- 10 - 19
- 20 - 29
- 30 - 39
- 40 - 49
- 50 - 59
- 60 - 69
- 70 - 79
- 80 - 89
- 90 - 99
- 100+

Sex:
- Male
- Female

Marital_status:
- Never married
- Married
- Widowed
- Divorced
- Civil partnership

Household_status:
- Child
- Institutional
- Living alone
- Partnership without children
- Married without children
- Partnership with children
- Married with children
- Single-parent
- Reference person in other hh
- Other
- missing

Size_prive_household:
- 1
- 2
- 3
- 4
- 5
- 6 - 10
- 11 or more
- missing

Education:
- No formal education
- Primary educ
- Low. Sec. educ
- Upper Sec. educ
- Post Sec. non-tertiary educ
- Bachelor / Master
- Not applicable (persons < 15yr)
- PhD
- missing

Activity_status:
- Not applicable
- Employed
- Pension
- Students
- Others
- Homemakers
- Unemployed
- missing

# Treemap: Structural Business Statistics



employees

47 - Retail trade, except of
G - Wholesale and retail trade; repair of motor vehicles and motorcycles
motor vehicles and motorcycles

46 - Wholesale trade, except of motor vehicles and motorcycles

45 - Wholesale and retail trade and repair of motor vehicles and motorcycles

49 - Land transport and transport via pipelines

H - Transporting and storage

52 - Warehousing and support activities for transportation

50 - Water transport

53 - Postal and courier activities

51 - Air transport

81 - Services to buildings and landscape activities

N - Administrative and support service activities

82 - Office administrative, office support and other business support activities

77 - Rental and leasing activities

78 - Employment activities

80 - Security and investigation activities

79 - Travel agency, tour operator and other reservation service and related activities

32 - Other manufacturing

25 - Manufacture of fabricated metal products, except machinery and equipment

28 - Manufacture of machinery and equipment n.e.c.

24 - Manufacture of basic metals

22 - Manufacture of rubber and plastic products

23 - Manufacture of other non-metallic mineral products

26 - Manufacture of computer, electronic and optical products

18 - Printing and reproduction of recorded media

33 - Repair and
C - Manufacturing
machinery and equipment

17 - Manufacture of paper and paper products

27 - Manufacture of electrical equipment

13 - Manufacture of textiles

10 - Manufacture of food products

20 - Manufacture of chemicals and chemical products

19 - Manufacture of coke and refined petroleum products

30 - Manufacture of other transport equipment

31 - Manufacture of furniture

71 - Architectural and engineering activities; technical testing and analysis

M - Professional, scientific and technical activities

69 - Legal and accounting activities

70 - Activities of head offices; management consultancy

73 - Advertising and market research

74 - Other professional, scientific and technical activities

72 - Scientific research and development

75 - Veterinary activities

62 - Computer programming, consultancy and related activities

J - Information and communication

61 - Telecommunications

58 - Publishing activities

63 - Information service activities

43 - Specialised construction activities

F - Construction

41 - Construction of buildings

42 - Civil engineering

D - Electricity, gas, steam and air conditioning supply

I - Accommodation and food service activities

E - Water supply; sewerage; waste management and remediation activities

B - Mining and quarrying

-70%  -60%  -50%  -40%  -30%  -20%  -10%  0%  50%  100%  150%

employees.prev

# Heatmap: Income statistics



36

# References

| Topic | Links |
|-------|-------|
| Social Media | Daas, P.J.H., Puts, M.J.H. (2014) Sociale Media Sentiment and Consumer Confidence. Paper for the Workshop on using Big Data for Forecasting and Statistics, Frankfurt, Germany. http://www.ecb.europa.eu/events/pdf/conferences/140407/Daas_Puts_Sociale_media_cons_conf_Stat_Neth.pdf?409d61b733fc259971ee5beec7cedc61 |
| Mobile phone metadata | Paper in progress... |
| Road sensors | Paper in progress... |
| Big Data for Official Statistics | Buelenes, B. et al. (2014) Selectivity of Big Data http://www.cbs.nl/nl-NL/menu/methoden/onderzoek-methoden/discussionpapers/archief/2014/2014-selectivity-of-big-data-pub.htm |
| Visualization | Tennekes, M., Jonge, E. de, Daas, P.J.H. (2013) Visualizing and Inspecting Large Datasets with Tableplots, Journal of Data Science 11 (1), 43-58. http://www.jds-online.com/file_download/379/JDS-1108.pdf<br><br>Tennekes, M., Jonge, E. de, Daas, P.J.H. (2012) Innovative visual tools for data editing. Paper presented at the United Nations Economic Commission for Europe (UNECE) Work Session on Statistical Data Editing, 2012, Oslo, Norway. http://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.44/2012/30_Netherlands.pdf |
| R packages by Statistics Netherlands (all on CRAN) | Visualization:    tabplot, tabplotd3, treemap, geo (in development only)<br>Data editing:    editrules, deducorrect, rspa<br>Large data processing:   ffbase, LaF<br>Other:  extremevalues, stringdist, whisker |