# Methods for unstructured data

Introduction
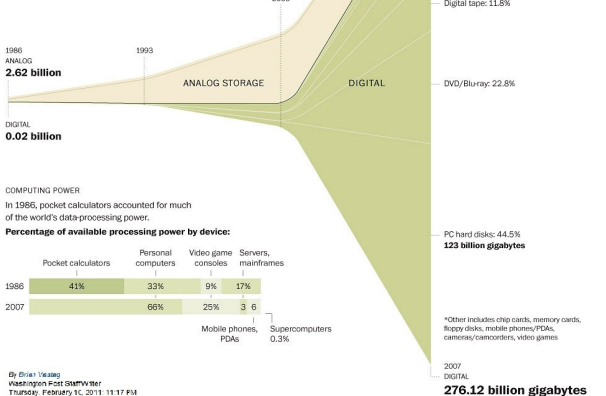
Helge Liebert

# Worldwide data storage capacity

**2007**
ANALOG
## 18.86 billion gigabytes

Paper, film, audiotape and vinyl: 6.2%
Analog videotapes: 93.8%

ANALOG
DIGITAL

Other digital media: 0.8%*
Portable media players, flash drives: 2%
Portable hard disks: 2.4%

CDs and minidisks: 6.8%

Computer servers and
mainframe hard disks: 8.9%

Digital tape: 11.8%

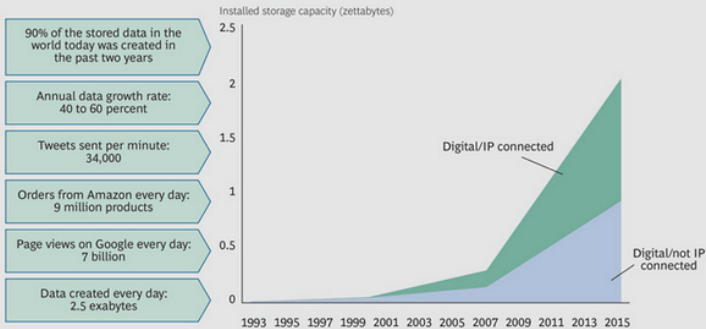DVD/Blu-ray: 22.8%

THE WORLD'S CAPACITY TO STORE INFORMATION
This chart shows the world's growth in storage capacity for both
analog data (books, newspapers, videotapes, etc.) and digital
(CDs, DVDs, computer hard drives, smartphone drives, etc.)
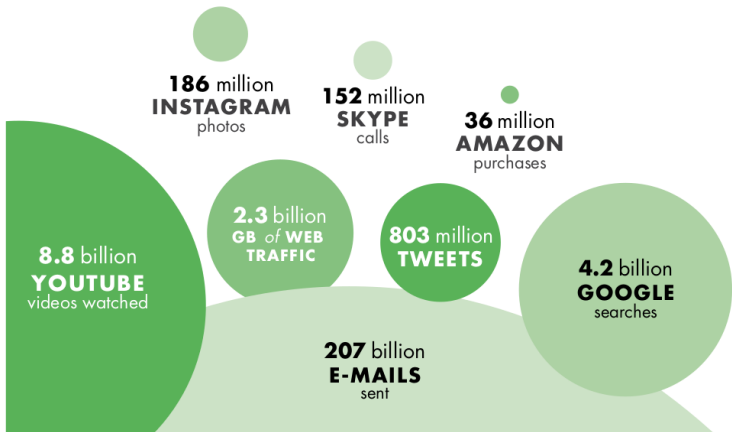
**In gigabytes or estimated equivalent**

2000

1986
ANALOG
**2.62 billion**

1993

ANALOG STORAGE          DIGITAL

DIGITAL
**0.02 billion**

PC hard disks: 44.5%
**123 billion gigabytes**

COMPUTING POWER
In 1986, pocket calculators accounted for much
of the world's data-processing power.

**Percentage of available processing power by device:**

|      | Pocket calculators | Personal computers | Video game consoles | Servers, mainframes |
|------|-----|-----|-----|-----|
| 1986 | 41% | 33% | 9%  | 17% |
| 2007 |     | 66% | 25% | 3   | 6 |

Mobile phones,
PDAs

Supercomputers
0.3%

*Other includes chip cards, memory cards,
floppy disks, mobile phones/PDAs,
cameras/camcorders, video games

2007
DIGITAL
## 276.12 billion gigabytes

By Brian Vastag
Washington Post Staff Writer
Thursday, February 10, 2011 11:17 PM

**EXHIBIT 2 | In 2015, More Than Half of All Data Will Have an IP Address**

90% of the stored data in the world today was created in the past two years

Annual data growth rate: 40 to 60 percent

Tweets sent per minute: 34,000

Orders from Amazon every day: 9 million products

Page views on Google every day: 7 billion

Data created every day: 2.5 exabytes

Installed storage capacity (zettabytes)

Digital/IP connected

Digital/not IP connected

**Sources:** Martin Hilbert and Priscilla Lopez, "The World's Technological Capacity to Store, Communicate, and Compute Information," *Science*, February 2011; BCG estimates.

# Data, then and now

b. A typical day in the life of the internet



**186** million **INSTAGRAM** photos

**152** million **SKYPE** calls

**36** million **AMAZON** purchases

**8.8** billion **YOUTUBE** videos watched

**2.3** billion **GB** *of* **WEB TRAFFIC**

**803** million **TWEETS**

**4.2** billion **GOOGLE** searches

**207** billion **E-MAILS** sent

*Sources:* World Development Indicators (World Bank, various years); WDR 2016 team; http://www.internetlivestats.com/one-second/ (as compiled on April 4, 2015). Data at http://bit.do/WDR2016-FigO_4.

*Note:* In panel a, for some years data for electricity are interpolated from available data. GB = gigabytes.

# Introduction

- 90% of data today has been created in the last two years.
- 235 million emails sent per day.
- 3.3 million Facebook posts created every minute.
- 3.8 million Google searches performed each minute.
- 1.7 megabytes of new information created every second, per person.

# Introduction

- 90% of data today has been created in the last two years.
- 235 million emails sent per day.
- 3.3 million Facebook posts created every minute.
- 3.8 million Google searches performed each minute.
- 1.7 megabytes of new information created every second, per person.

➡ An immense amount of data, new and old, is recorded as text.

➡ More generally, much of this data is unstructured.

# Structured vs. unstructured

### Structured data

- Adheres to a defined data model.
- Examples: Tables, relational/hierarchical databases, …

### Unstructured data

- Does *not* adhere to a defined data model.
- Typically text-heavy.
- Examples: Text feeds, speech transcripts, audio, images …

# Structured vs. unstructured

### Structured data

- Adheres to a defined data model.
- Examples: Tables, relational/hierarchical databases, ...

### Semi-structured data

- Does *not* adhere to a formal data model,
- ... *but* contains tags or semantic mark-up.
- Examples: JSON, XML, emails, tagged text, ...

### Unstructured data

- Does *not* adhere to a defined data model.
- Typically text-heavy.
- Examples: Text feeds, speech transcripts, audio, images ...

# Text as data

- Text differs from other, traditional forms of data.
- Text is inherently *unstructured* and *high-dimensional*.
- One of the major fields of application of machine learning methods.
- Fast-growing field. Many new techniques developed in industry.
- Recent applications in economics and other social sciences.

## This lecture

This lecture covers techniques for unstructured data.

- Methods for wrangling data.
- ➡ When unstructured ≈ dirty (or differently structured).

# This lecture

This lecture covers techniques for unstructured data.

- Methods for wrangling data.
- ➥ When unstructured $\approx$ dirty (or differently structured).

- Methods for analyzing data which are naturally unstructured.
- ➥ No rectangular (or graph) structure, no well-defined relations between data elements.

# Focus points

Focus on three main points.

1. Processing and transforming un-/semi-structured data.
2. Representing inherently unstructured text data.
3. Analyzing text data and using models to discover structure.
   (Supervised and unsupervised learning.)

# Outline

# Dates

| | | | |
|---|---|---|---|
| Tuesday | 15.09.2020 | 09.15-12.00 | PC-Lab S18 HG.37 |
| Wednesday | 16.09.2020 | 14.15-18.00 | PC-Lab S18 HG.37 |
| Friday | 18.09.2020 | 09.15-12.00 | PC-Lab S18 HG.37 |
| Tuesday | 22.09.2020 | 09.15-12.00 | PC-Lab S18 HG.37 |
| Wednesday | 23.09.2020 | 14.15-18.00 | PC-Lab S18 HG.37 |
| Friday | 25.09.2020 | 09.15-12.00 | PC-Lab S18 HG.37 |
| Wednesday | 30.09.2020 | 14.15-18.00 | PC-Lab S18 HG.37 |

## Technical requirements: Lab sessions

- All class material is available online:
  https://github.com/hliebert/course-unstructured-data.
- All material will run on the Windows computers in the lab.
- The lab materials can also be accessed online:
  Jupyter notebooks
  Rstudio server

- Alternatively, a linux virtual machine is available on the shared drive in the lab. Copy it to your computer and run it using VirtualBox. Has all course dependencies pre-installed. Extend the VM to your liking.
- Feel free to set up your own computer. Please ask if you need help.
- Clone/download the course repository to get started.

## Required programs

- R.
- An editor or GUI (RStudio, VScode or Atom with R plugins, ...)
- If you want to use Jupyter notebooks or Python, install Anaconda (or its smaller miniconda version). On Linux, Anaconda is not necessary (but may be preferred to using pip and virtualenv).
- Run the R install script provided with the class material to install the R package dependencies and the R Kernel for Jupyter notebooks.
- A shell (bash or zsh) is pre-installed on Linux and MacOS. On MacOS, iterm2 is more convenient than Terminal. Use Windows Subsystem for Linux (WSL) on Windows, or Git Bash for a minimal set of features.
- To use Selenium from R, install Docker.
- To use Git (and minimal Bash on Windows), install Git.

## How to install them

- Linux: Use your distribution's package manager.
- Mac: Use installer packages or set up and use homebrew as a package manager.
- Windows: Use installer packages or look into scoop or chocolatey as native package managers for Windows. WSL, Cygwin, Msys2 provide access to Unix functions on Windows.

# Assignment

1. Web scraping assignment (20%)
2. Text analysis and research proposal (80%)

- Deadline: 25.10.2020.
- Grading is pass/fail.
- More details during the course of the lecture.

# Primary references

- The course covers relatively broad and diverse ground, no single reference.
- Primary and secondary references below. Seminal references in the slides,
- Books are available for free online (use newest draft of Jurafsky & Martin).

📄 Gentzkow, M., B. Kelly, and M. Taddy (2019). Text as Data. *Journal of Economic Literature* 57(3), 535–574. DOI: 10/gf7rd5.

📕 Hastie, T., R. Tibshirani, and J. Friedman (2001). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.* Ed. by R. Tibshirani and J. H. ( H. Friedman. New York.

📕 Jurafsky, D. and J. H. Martin (2009). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.* 2nd ed. Prentice Hall Series in Artificial Intelligence. Upper Saddle River, N.J: Pearson Prentice Hall.

# Secondary references

- Reference material, applied or introductory text books.

Baumer, B., D. Kaplan, and N. Horton (2017). *Modern Data Science with R*. CRC.

Casella, G. and R. L. Berger (2001). *Statistical Inference*. Second. Duxbury Press.

James, G., D. Witten, T. Hastie, and R. Tibshirani (2015). *An Introduction to Statistical Learning with Applications in R*. Springer.

Matloff, N. (2011). *The Art of R Programming: A Tour of Statistical Software Design*. No Starch Press.

Mitchell, R. E. (2018). *Web Scraping with Python: Collecting More Data from the Modern Web*. Second edition. Sebastopol, CA: O'Reilly Media.

Munzert, S. (2014). *Automated Data Collection with R: A Practical Guide to Web Scraping and Text Mining*. Chichester, West Sussex, United Kingdom: Wiley.

Shotts, W. E. (2019). *The Linux Command Line: A Complete Introduction*. Second edition. San Francisco: No Starch Press.

Silge, J. and D. Robinson (2017). *Text Mining with R: A Tidy Approach*. First edition. Boston: O'Reilly.

Wasserman, L. (2006). *All of Nonparametric Statistics*. Springer.

Wasserman, L. (2010). *All of Statistics: A Concise Course in Statistical Inference*. Springer Texts in Statistics. New York, NY: Springer.