

STAT 380: Week 1

Instructor: Murali Haran Professor of Statistics
TA: John Ensley, PhD Student

- ▶ Use the computer expressively to prepare, explore, and analyze data
- ▶ Work closely with original raw data
- ▶ Use existing software rather than build routines from the ground up.
- ▶ Focus on aspects of computing to conduct statistical analysis, NOT the computational aspects of statistical methods (For that: STAT 440, Computational Statistics)
- ▶ Book:
 - ▶ Data Technologies and Computational Reasoning by D. Nolan and D. Temple Lang (pdf files will be posted weekly).
 - ▶ Supplement: *Data Science in R: A Case Studies Approach to Computational Reasoning* by Nolan and Temple Lang.

(With thanks to Professor Nolan for lecture notes)

1 / 17

2 / 17

What are data?

- ▶ Data are recorded/measured observations together with context.
- ▶ By context we mean the details of who, what, where, when, and/or how the observations were obtained, aka "metadata".

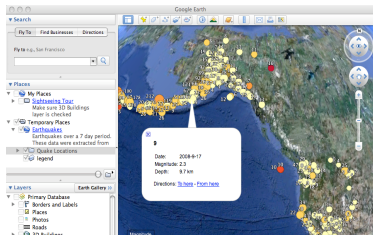
Tables of Numbers

Traffic on I-80

[illegible]

Geographic Information and Time

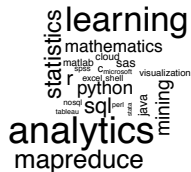
Earthquake Location, Date, and Magnitude



17

Text

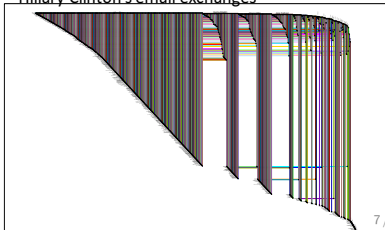
Kaggle Job Postings for a Data Scientist



6 / 17

Graph

Hillary Clinton's email exchanges



7 / 17

Meta-data:

Information about Spotify playlists

```
{
  "href": "https://api.spotify.com/v1/users/spotify_espa%C3%B1a/playlists/21T8aBj9TaSGuKYNBUSTac/tracks",
  "items": [ {
    "added_at": "2014-08-18T20:16:08Z",
    "added_by": {
      "external_urls": {
        "spotify": "http://open.spotify.com/"
      },
      "href": "https://api.spotify.com/v1/users/spotify_espa%C3%B1a",
      "id": "spotify_espa%C3%B1a",
      "type": "user",
      "uri": "spotify:user:spotify_espa%C3%B1a"
    }
  } ]
}
```

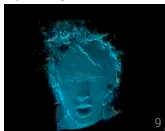
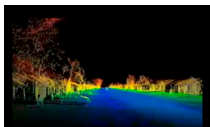


8 / 17

Images, Video, or Audio

Radiohead House of Cards

Yorke: "I liked the idea of making a video of human beings and real life and time without using any cameras, just lasers, so there are just mathematical points – and how strangely emotional it ended up being."



9 / 17

AIG job posting on Kaggle for senior data scientist:

- ▶ Build predictive models utilizing both traditional statistical methods and modern machine learning techniques
- ▶ Extract, clean, and manipulate large datasets (structured and unstructured) for model building.
- ▶ Communicate (written and verbal) insights from quantitative analyses to technical and non-technical audiences.
- ▶ Stay current on the latest machine learning and big data trends.
- ▶ Work with business sponsors and IT teams to implement analytic solutions.
- ▶ Serve as a technical expert on one or more domains (e.g. Time Series Analysis, Text Mining, etc.)

10 / 17

What Skills does a Data Scientist need?

AIG job postings on Kaggle

- ▶ Expertise in at least one modeling/machine learning platform such as R, Python, or SAS.
- ▶ Knowledge of an additional general purpose programming language such as C++ or Java.
- ▶ Advanced SQL skills and experience with No SQL technologies.
- ▶ Built several predictive models that have been put into live production.
- ▶ Obsess over sample bias, over-fitting, variable selection, missing values, etc.
- ▶ Understand the need to balance predictive power, interpretability, and ease of implementation

11 / 17

Data analysis cycle

- ▶ Data ACQUISITION Input/output, regular expressions
- ▶ Data CLEANING verification, manipulation
- ▶ Data ORGANIZATION data frames, data bases, XML
- ▶ Data EXPLORATION search for interesting patterns
- ▶ Data VISUALIZATION create statistical graphs
- ▶ Data ANALYSIS fit and assess statistical models
- ▶ Data SIMULATION studies of random behavior
- ▶ Data REPORTING report findings from analysis

12 / 17

- ▶ Basic numeracy: Variability, Patterns, comparisons
- ▶ Exploratory Data Analysis
- ▶ Graphics: Elements and principles of graphing
- ▶ Computationally intensive methods, e.g., Classification and Regression trees, multi-dimensional scaling, nearest neighbor method
- ▶ Simulation tools: Monte Carlo, bootstrap, cross-validation

13 / 17

- ▶ R statistical software
- ▶ SQL structured query language for relational databases
- ▶ XML Extensible Markup Language (and HTML) and XPath
- ▶ Unix shell commands

15 / 17

- ▶ Programming concepts Control flow trees functions
- ▶ Regular expressions and text manipulation
- ▶ Relational databases
- ▶ Random number generation
- ▶ Representation of information in the computer

14 / 17

- ▶ Homework + projects: add up to 50%. Exact proportion may change. *Tentatively:*
 - ▶ Homework = 30% + Projects = 20%
- ▶ Homework due in class. After class, before 3:30pm (in my mailbox in Thomas 326): 20% off. After that, 0 credit no matter what.
- ▶ Drop two lowest homework scores.
- ▶ Midterm: 20%
- ▶ Final: 30%

16 / 17

Academic integrity

- ▶ Free to discuss course matters with instructor, TA, and fellow students
- ▶ DO NOT SHARE CODE
- ▶ Make significant contribution to your groups work
- ▶ If you are uncertain as to whether something may be a violation of the code, ask the instructor
- ▶ Writing a program is like writing a paper your code should be your original work.
- ▶ A violation will result in at least one of the following: 0 on the assignment, F for the course grade, Report to the Office of Student Conduct